# Interdisciplinary Applied Mathematics

## Volume 23

Problems in engineering, computational science, and the physical and biological sciences are using increasingly sophisticated mathematical techniques. Thus, the bridge between the mathematical sciences and other disciplines is heavily traveled. The correspondingly increased dialog between the disciplines has led to the establishment of the series: *Interdisciplinary Applied Mathematics.*

The purpose of this series is to meet the current and future needs for the interaction between various science and technology areas on the one hand and mathematics on the other. This is done, firstly, by encouraging the ways that mathematics may be applied in traditional areas, and well as point towards new and innovative areas of applications; and, secondly, by encouraging other scientific disciplines to engage in a dialog with mathematicians outlining their problems to both access new methods and suggest innovative developments within mathematics itself.

The series will consist of monographs and high-level texts from researchers working on the interplay between mathematics and other fields of science and technology.

# Interdisciplinary Applied Mathematics

Volumes published are listed at the end of this book.

**Springer**

*New York*
*Berlin*
*Heidelberg*
*Hong Kong*
*London*
*Milan*
*Paris*
*Tokyo*

Muhammad Sahimi

# Heterogeneous Materials

## Nonlinear and Breakdown Properties and Atomistic Modeling

With 119 Illustrations

Springer

Muhammad Sahimi
Department of Chemical Engineering
University of Southern California
Los Angeles, CA 90089-1211
USA
moe@iran.usc.edu

*Editors*

J.E. Marsden
Control and Dynamical Systems
Mail Code 108-81
California Institute of Technology
Pasadena, CA 91125
USA
marsden@cds.caltech.edu

L. Sirovich
Division of Applied Mathematics
Brown University
Providence, RI 02912
USA
chico@camelot.mssm.edu

S. Wiggins
School of Mathematics
University of Bristol
Bristol, BS8 1TW
United Kingdom
s.wiggins@bristol.ac.uk

S.S. Antman
Department of Mathematics
*and*
Institute of Physical Science and Technology
University of Maryland
College Park, MD 20742
USA
ssa@math.umd.edu

*Cover illustration:*

*To children of the third world*
*who have the talent but not the means to succeed*
*and to*
*the memory of my father, Habibollah Sahimi,*
*who instilled in me, a third world child, the love of reading*

# Preface

Disorder plays a fundamental role in many natural and man-made systems that are of industrial and scientific importance. Of all the disordered systems, heterogeneous materials are perhaps the most heavily utilized in all aspects of our daily lives, and hence have been studied for a long time. With the advent of new experimental techniques, it is now possible to study the morphology of disordered materials and gain a much deeper understanding of their properties. Novel techniques have also allowed us to design materials of morphologies with the properties that are suitable for intended applications.

With the development of a class of powerful theoretical methods, we now have the ability for interpreting the experimental data and predicting many properties of disordered materials at many length scales. Included in this class are renormalization group theory, various versions of effective-medium approximation, percolation theory, variational principles that lead to rigorous bounds to the effective properties, and Green function formulations and perturbation expansions. The theoretical developments have been accompanied by a tremendous increase in the computational power and the emergence of massively parallel computational strategies. Hence, we are now able to model many materials at molecular scales and predict many of their properties based on first-principle computations.

In this two-volume book we describe and discuss various theoretical and computational approaches for understanding and predicting the effective macroscopic properties of heterogeneous materials. Most of the book is devoted to comparing and contrasting the two main classes of, and approaches to, disordered materials, namely, the continuum models and the discrete models. Predicting the effective properties of composite materials based on the continuum models, which are based on solving the classical continuum equations of transport, has a long history and goes back to at least the middle of the nineteenth century. Even a glance at the literature on the subject of heterogeneous materials will reveal the tremendous amount of work that has been carried out in the area of continuum modeling. Rarely, however, can such continuum models provide accurate predictions of the effective macroscopic properties of *strongly* disordered multiphase materials. In particular, if the contrast between the properties of a material's phases is large, and the phases form large clusters, most continuum models break down. At the same time, due to their very nature, the discrete models, which are based on a lattice representation of a material's morphology, have the ability for providing accurate predictions for the effective properties of heterogeneous materials, even when the heterogeneities are strong, while another class of discrete models, that represent a material as a collection of its constituent atoms and molecules, provides accurate predictions of

the material's properties at mesoscopic scales, and thus, in this sense, the discrete models are complementary to the continuum models. The last three decades of the twentieth century witnessed great advances in discrete modeling of materials and predicting their macroscopic properties, and one main goal of this book is to describe these advances and compare their predictions with those of the continuum models. In Volume I we consider characterization and modeling of the morphology of disordered materials, and describe theoretical and computational approaches for predicting their *linear* transport and optical properties, while Volume II focuses on nonlinear properties, and fracture and breakdown of disordered materials, in addition to describing their atomistic modeling. Some of the theoretical and computational approaches are rather old, while others are very new, and therefore we attempt to take the reader through a journey to see the history of the development of the subjects that are discussed in this book. Most importantly, we always compare the predictions with the relevant experimental data in order to gain a better understanding of the strengths and/or shortcomings of the two classes of models.

A large number of people have helped me gain deeper understanding of the topics discussed in this book, and hence have helped me to write about them. Not being able to name them all, I limit myself to a few of them who, directly or indirectly, influenced the style and contents of this book. Dietrich Stauffer has greatly contributed to my understanding of percolation theory, disordered media, and critical phenomena, some of the main themes of this book; I am deeply grateful to him. For their tireless help in the preparation of various portions of this book, I would like to thank two of my graduate students, Sushma Dhulipala and Alberto Schroth. Although they may not be aware of it, Professors Pedro Ponte Castañeda of the University of Pennsylvania and Salvatore Torquato of Princeton University provided great help by guiding me through their excellent work, which is described in this book; I would like to thank them both. Some of my own work described in this book has been carried out in collaboration with many people; I am pleased to acknowledge their great contributions, especially those of Dr. Sepehr Arbabi, my former doctoral student. The constant encouragement and support offered by many of my colleagues, a list of whom is too long to be given here, are also gratefully acknowledged. I would like particularly to express my deep gratitude to my former doctoral student Dr. Jaleh Ghassemzadeh, who provided me with critical help at all stages of preparation of this book. Several chapters of this book have been used, in their preliminary versions, in some of the courses that I teach, and I would like to acknowledge the comments that I received from my students.

My wife, Mahnoush, and son, Ali, put up with the countless hours, days, weeks, and months that I spent in preparing this book and my almost complete absence during the time that I was writing, but never denied me their love and support without which this book would have never been completed; I love and cherish them both.

Muhammad Sahimi
Los Angeles, California, USA
May 2002

# Contents

## III  Atomistic and Multiscale Modeling of Materials    455

# Abbreviated Contents for Volume I

Preface

Abbreviated Contents for Volume II

# Introduction to Volume II

In Volume I of this book, we presented a self-contained analysis of the morphology of heterogeneous materials and their effective linear properties. Some of the properties of heterogeneous materials that were studied in Volume I were,

(1) the effective (electrical, thermal, hopping, and Hall) conductivity;
(2) the effective dielectric constant and optical properties, and
(3) the effective elastic moduli.

In addition, we also considered some aspects of the classical (as opposed to quantum-mechanical) superconductivity of composite materials. Both steady-state and time- and frequency-dependent properties were considered, and the most significant theoretical developments for modelling the morphology of heterogeneous materials and predicting their effective linear properties were described in detail. In addition, we also described the techniques for computer simulations of transport processes in disordered materials, and compared their predictions with the theoretical ones and also the relevant experimental data.

In the present Volume, we continue our study of transport processes in heterogeneous materials, except that we consider their effective *nonlinear* properties. After the introductory Chapter 1 in which we study characterization of surface structure of materials when the surface is rough, we embark on studying various nonlinear processes in heterogeneous materials. To do this, we divide nonlinear transport processes into two groups, which are as follows.

## A. Constitutive Nonlinearity

Materials of this type *always* behave nonlinearly. For example, if in a composite material the relation between the current $I$ and voltage $V$ is given by

$$I = gV^n$$

where $g$ is a generalized conductance of the material, then, as far as the electrical conductivity is concerned, for $n \neq 1$ the material *always* behaves nonlinearly. We will study such nonlinear phenomena in Chapters 2–4, and describe various approaches for predicting and estimating the effective nonlinear conductivity, dielectric constant, optical properties, and elastic moduli and rigidity.

## B. Threshold Nonlinearity

In this class of materials are those for which the nonlinearity arises as a result of imposing on them an external field of sufficient intensity. Brittle fracture and dielectric breakdown of composite solids are two important examples of such nonlinear transport processes. In brittle fracture, for example, the elastic response of a solid material is governed by the equations of linear elasticity until the external stress or strain that has been imposed on the material exceeds a critical value, at which time the material breaks down and microcracks begin to emerge. A list of all possible nonlinear transport processes of this type is very long. This type of nonlinearity will be studied in Chapters 5–8, and will include electrical and dielectrical breakdown, brittle fracture, and the transition between brittle fracture and ductile behavior.

One important point to remember is that, the interplay between a nonlinear transport process and the disordered morphology of a composite material gives rise to a rich variety of phenomena that are usually far more complex than what one usually must deal with in linear processes. Over the past 15 years, an increasing number of investigations have been devoted to such nonlinear transport processes, and deeper insight into their properties has been acquired. A major goal of Volume II is to describe this progress and compare various properties of nonlinear transport processes in heterogeneous materials with their linear counterparts.

## C. Theoretical Approaches

Although the analysis of transport processes in composite materials has a long history, it is only in the past three decades that this analysis has been extended to include detailed structural properties of the materials, and in particular the distribution of their heterogeneities. Deriving *exact* results for the effective properties of composite materials with anything but the simplest morphologies is extremely difficult, if not impossible, and thus one must resort to various approximate techniques. At the same time, however, the advent of powerful computers and development of efficient computational algorithms have allowed us to estimate various properties of heterogeneous materials to practically any desired or affordable accuracy.

To describe the theoretical approaches for estimating the effective properties of composite materials, we divide them into two classes. In the first class of models are what we refer to as the *continuum models*, while the second class is made of the *discrete models*. Both types of models are described and analyzed in this Volume, and what follows is a brief description of the general features of each class of models.

## C.1 The continuum models

The physical laws that govern the transport processes at the microscopic level are well understood. Thus, one can, in principle, write down the differential equations that describe transport of energy, charge, or stress in a material and specify the associated initial and boundary conditions. However, as the morphology of most real composite materials is very irregular, practical and economically fea-

sible computations for exact estimation of the effective properties are still very difficult—even in the event that one knows the detailed morphology of the material. Thus, it becomes essential to adopt a macroscopic description at a length scale much larger than the dimension of the individual phases of a composite material. The governing equations are then discretized and solved numerically, provided that the effective properties that appear in the transport equations are either supplied as the inputs (through, for example, experimental measurements), or else a model for the morphology of the material is assumed so that the effective transport properties can be somehow estimated, so that the numerical solution yields other quantities of interest, such as the potential distribution in the material. We refer to various models associated with this classical description as the *continuum models*. These models have been widely used because of their convenience and familiarity to the engineers and materials scientists. Their limitations will be described and discussed in the subsequent chapters.

In addition to deriving the effective macroscopic equations and obtaining their solution by numerical calculations, one may also derive exact results in terms of rigorous upper and lower bounds to the properties of interest. Hence, powerful tools have been developed for deriving accurate upper and lower bounds and estimates. Finally, various approximations, such as the mean-field and effective-medium approximations, have also been developed in the context of the continuum models. We will describe most of these theoretical approaches throughout both this book and Volume II.

**C.2 The discrete models**
The second class of models, the *discrete models*, are free of many limitations of the continuum models. They themselves are divided into two groups.

(1) In the first class of discrete models, a material is represented by a discrete set of atoms and molecules that interact with each other through interatomic potentials. In a solid material, the distance between the atoms is fixed. One then carries out atomistic simulations of the materials' behavior under a variety of conditions. Several types of such simulations have been developed over the past few decades. With the advent of massively-parallel computational algorithms, atomistic simulations have increasingly become a viable and *quantitative* method of predicting the effective properties of materials. We will describe such approaches in Chapters 9 and 10.

(2) In the second class are lattice models of composite materials. The bonds of the lattices represent microscopic elements of the material. For example, they represent a conducting or insulating elements, or an elastic or a plastic region. They do *not* represent molecular bonds, and therefore such lattice models are appropriate for length scales that are much larger than molecular scales. These models have been advanced to describe various phenomena at the microscopic level and have been extended in the last several years to also describe them at the macroscopic length scales. We will describe both classes of discrete models in this volume.

The main shortcoming of both groups of the discrete models, from a practical point of view, is the large computational effort required for a realistic discrete representation of the material and simulating its behavior, although the ever-increasing computational power is addressing this difficulty.

**D. The Organization of the Book**

What we intend to do in this Volume, similar to Volume I, is describing the most important developments in predicting the effective nonlinear properties of composite materials, and comparing the predictions with the relevant experimental data. To accomplish our goal, for each effective property we describe and discuss the continuum and discrete models separately. Then, in Chapter 10 we describe recent advances in *multiscale* modelling of materials' properties—a method that combines a discrete approach with a continuum model. Similar to Volume I, the structure of each chapter is as follows.

(1) The main problem(s) of interest is (are) introduced.
(2) The problem(s) is (are) then analyzed by several methods, each of which provide valuable insight into the solution of the problem(s) and the physical phenomena that it (they) represent. Typically, each chapter starts with exact and rigorous results, then describes analytical approximations, and finally discusses the numerical and computer simulation methods. The weakness and strengths of each method are also pointed out. In this way, the most important progress in understanding the physical phenomena of interest is described and discussed.
(3) When possible (which is almost always the case), we compare the theoretical predictions with the experimental data and/or high-resolution computer simulation results.

Characterization of surface morphology materials, which are directly relevant to most of what is discussed in this Volume, will be described in Chapter 1. Aside from this chapter, this Volume is divided into three parts. In Part I (Chapters 2–4) we study transport processes in heterogeneous materials that are characterized by constitutive nonlinearities. Part II (Chapters 5–8) contains the description and discussion of transport processes with threshold nonlinearity, including electrical and dielectric breakdown, and brittle fracture of disordered materials. Finally, in Part III we will describe (in Chapters 9 and 10) advances in atomistic modelling of materials, and how a powerful new approach that combines atomistic simulations with the continuum description—in effect a combination of a discrete approach with a continuum model—promises to provide much deeper understanding of materials, and deliver quantitative predictions for their effective properties.

Let us emphasize that, as in Volume I, although every attempt has been made to discuss and cite the relevant literature on every subject that we consider, what we do cite and bring to the attention of the reader represents what was known to us at the time of writing this book, and/or what we considered to be the most relevant. As such, this two-volume book represents the author's biased view of the subject of composite materials.

# 1
# Characterization of Surface Morphology

## 1.0   Introduction

Natural, as well as man-made, materials have enormous variations in their morphology, which consists of materials' geometry, topology and surface structure. The geometry refers to sizes of the micro- and mesoscale elements of the materials, as well as their shapes which range anywhere from completely ordered to complex and seemingly chaotic patterns. Generally speaking, regular Euclidean shapes are formed under close-to-equilibrium conditions, although even in such cases equilibrium thermodynamics is often incapable of describing the process that gives rise to such shapes. The topology of materials describes how the micro- and mesoscale elements are connected to one another. The structure of materials' surface, especially those that are produced under far-from-equilibrium conditions, is also very important because the surface is often very rough and possesses complex features. In recent years it has become clear that characterizing the surface roughness will go a long way toward giving us a much better understanding of materials' microstructure and hence many of their effective properties. However, when we speak of surface roughness, we must specify the length scales over which the roughness is measured. Even the most rugged mountains look perfectly smooth when viewed from the outer space! Therefore, surface roughness (and, more generally, all the morphological characteristics) *depends on the length scale of observations or measurements*. The effect of topology of disordered materials on their effective transport properties is quantified by percolation theory which, together with the effect of the geometry, was described in Chapters 2 and 3 of Volume I, and their significance was emphasized throughout Volume I where we analyzed effective linear properties of disordered materials. For other applications of percolation theory see Sahimi (1994a). Stauffer and Aharony (1992) present a simple introduction to the concepts of percolation theory. In this chapter, we consider the structure and characteristics of materials' surface, and describe various theoretical and experimental methods of studying rough surfaces, which are directly relevant to the nonlinear phenomena in heterogeneous materials considered in this Volume, particularly to their brittle fracture and dielectric breakdown.

An important example of a material with a rough surface are the thin films that produced by molecular beam epitaxy, and are utilized for manufacturing of semiconductors and computer chips. These films are made of silicon and other

elements, and are prepared by deposition of atoms on a very clean surface. Thin films with rough surfaces are also made by sputtering in which an energized beam of particles is sent toward the bulk of a material. Collision of the beam particles with the material causes ejection of some particles from the material's surface, which then deposit on another surface and start to grow a thin film of the original material.

Although the enormous variations in the morphology of natural, and even man-made, materials, particularly in their surface, are such that, up until a few decades ago, the problem of describing and quantifying such morphologies seemed hope-less, many experimental and theoretical developments of the past two decades have brightened the prospects for deeper understanding of materials' microstructures, and in particular the structure of their surface. Among them are the advent of pow-erful computers and novel experimental techniques that allow highly sophisticated computations of materials' properties and their measurement. In addition, the real-ization that the complex microstructure and behavior of a wide variety of materials can be quantitatively characterized by utilizing the ideas of fractal distributions, have advanced our understanding of materials' surface structure. As we discuss in this chapter, fractal concepts provide us with a powerful tool for characterizing the structure of materials' surface and its roughness, and the long-range correlations that often exist in their morphology.

The purpose of this chapter is to describe and discuss the essential features of surface morphology and its dynamics during the process in which it is formed, and how fractal concepts can be utilized for characterizing it. We already described in Chapter 2 of Volume I most of the main concepts of fractal geometry, and therefore in this chapter we restrict ourselves to a brief discussion of such concepts, after which we study and analyze rough surfaces.

## 1.1   Self-Similar Fractal Structures

An intuitive and informal definition of a self-similar fractal object is that, in such objects the part is reminiscent of the whole, implying that the object possesses scale-invariant properties, i.e., its morphology repeats itself at different length scales. This means that above a certain length scale—the lower cutoff scale for fractality—the structure of a piece of the object can be magnified to recover its structure at larger length scales up to another length scale—the upper cutoff for its fractality. Below the lower cutoff and above the upper cutoff scales the system loses its self-similarity. While there are disordered media that are self-similar at *any* length scale, natural materials and media that exhibit self-similarity typically lose their fractal characteristics at sufficiently small or large length scales.

One of the simplest characteristics of a self-similar fractal system is its fractal dimension $D_f$, which is defined as follows. We cover the fractal system by non-overlapping $d$-dimensional spheres of Euclidean radius $r$, or boxes of linear size $r$, and count the number $N(r)$ of such spheres that is required for complete coverage

of the system. The fractal dimension $D_f$ of the system is then defined by

$$D_f = \lim_{r \to 0} \frac{\ln N}{\ln(1/r)}.$$  (1)

Estimating the fractal dimension through the use of Eq. (1) is called the *box-counting method*. For non-fractal objects, $D_f = d$, where $d$ is the Euclidean dimensionality of the space in which they are embedded. Note that, in order to be able to write down Eq. (1), we have implicitly assumed the existence of a lower and an upper cutoff length scale for the fractality of the system which are, respectively, the radius $r$ of the spheres and the linear size $L$ of the system.

One can also define the fractal dimension $D_f$ through the relation between the system's mass $M$ and its characteristic length scale $L$. If the system is composed of particles of radius $r$ and mass $m$, then

$$M = cm(L/r)^{D_f},$$  (2)

where $c$ is a geometrical constant of order 1. Since we can fix the dependence of $M$ on $m$ and $r$, we can write

$$M(L) \sim L^{D_f}.$$  (3)

Often, measuring $M$ entails using an ensemble of samples with similar structures, rather than a single sample. In this case

$$\langle M \rangle = cm(L/r)^{D_f},$$  (4)

where $\langle \cdot \rangle$ implies an average over the mass of a large number of samples with linear sizes in the range $L \pm \delta L$, centered on $L$.

Most natural fractals are what we call *statistically self-similar* because their self-similarity is only in an average sense. One of the most important examples of such fractals is one which is generated by the *diffusion-limited aggregation* model (Witten and Sander, 1981). In this model the site at the center of a lattice is occupied by a stationary particle. A new particle is then injected into the lattice, far from the center, which diffuses on the lattice until it reaches a surface site, i.e., an empty site which is a nearest neighbor of the stationary particle, at which time the particle sticks to it and remains there permanently. Another diffusing particle is then injected into the lattice to reach another surface (empty) site and stick to it, and so on. If this process continues for a long time, a large aggregate is formed. The most important property of diffusion-limited aggregates is that they have a self-similar fractal structure (for a review see, for example, Meakin, 1998) with $D_f \simeq 1.7$ and 2.45 for 2D and 3D aggregates, respectively. A two-dimensional (2D) example of such aggregates is shown in Figure 1.1. Diffusion-limited aggregates have found wide applications, ranging from colloidal systems, to miscible displacement processes in porous media, to describing cellular patterns in human bone marrow (Naeim *et al.*, 1996). We will come back to this model in Chapters 5 and 8, where we describe models of dielectric breakdown and fracture of composite materials.

FIGURE 1.1. A two-dimensional diffusion-limited aggregate.

## 1.2   The Correlation Function

A powerful method for testing self-similarity of disordered media is to construct a correlation function $C_n(\mathbf{r}^n)$ defined by

$$C_n(\mathbf{r}^n) = \langle \rho(\mathbf{r}_0)\rho(\mathbf{r}_0 + \mathbf{r}_1) \cdots \rho(\mathbf{r}_0 + \mathbf{r}_n) \rangle, \tag{5}$$

where $\rho(\mathbf{r})$ is the density at position $\mathbf{r}$, and the average is taken over all possible values of $\mathbf{r}_0$. Here $\mathbf{r}^n$ denotes the set of points at $\mathbf{r}_1, \cdots, \mathbf{r}_n$. If an object is self-similar, then its correlation function defined by Eq. (5) should remain the same, up to a constant factor, if all the length scales of the system are rescaled by a constant factor $b$. Thus, one must have

$$C_n(b\mathbf{r}_1, b\mathbf{r}_2, \cdots, b\mathbf{r}_n) = b^{-nx} C_n(\mathbf{r}_1, \cdots, \mathbf{r}_n). \tag{6}$$

It is not difficult to see that only a *power-law* correlation function can satisfy Eq. (6). Moreover, it can be shown that one must have $x = d - D_f$, where the quantity $x$ is called the *co-dimensionality*. However, in most cases only the two-point, or the direct, correlation function can be computed or measured with high precisions, and therefore we focus on this quantity. In practice, to construct the direct correlation function for use in analyzing a self-similar fractal structure, one typically employs a digitized image of the system. The correlation function is then written as

$$C(\mathbf{r}) = \frac{1}{\Omega} \sum_{\mathbf{r}'} s(\mathbf{r}')s(\mathbf{r} + \mathbf{r}'), \tag{7}$$

where $s(\mathbf{r})$ is a function such that $s(\mathbf{r}) = 1$ if a point at $\mathbf{r}$ belongs to the system, $s(\mathbf{r}) = 0$ otherwise, and $r = |\mathbf{r}|$. Because of self-similarity of the system, the direct

correlation function $C(r)$ *decays* as

$$C(r) \sim r^{D_f - d}. \tag{8}$$

This power-law decay of $C(r)$ not only provides a test of self-similarity of a disordered medium or material, it also gives us a means of estimating its fractal dimension since, according to Eq. (8), if one prepares a logarithmic plot of $C(r)$ versus $r$, then for a fractal object one should obtain a straight line with a slope $D_f - d$. Estimating the fractal dimension based on the direct correlation function has proven to be a very robust and reliable method. Equation (8) has an important implication: There are long-range correlations in a self-similar fractal system, because $C(r) \to 0$ only when $r \to \infty$. The existence of such correlations has important implications for estimating the effective transport properties of disordered materials (see, for example, Sahimi, 1994b, 1995a, and references therein). Other experimental methods of estimating the fractal dimension were described in Chapter 2 of Volume I, and therefore are not repeated here.

## 1.3    Rough Surfaces: Self-affine Fractals

The self-similarity of a fractal structure implies that its microstructure is invariant under an isotropic rescaling of lengths, i.e., if all lengths in all directions are rescaled by the same scale factor. However, there are many fractals that preserve their scale-invariance only if lengths in different directions are rescaled by factors that are direction dependent. In other words, the scale-invariance of such systems is preserved only if lengths in $x$-, $y$-, and $z$-directions are scaled by scale factors $b_x$, $b_y$, and $b_z$, where in general these scale factors are not equal. This type of scale-invariance implies that the fractal system is, in some sense, *anisotropic*. Such fractal systems are called *self-affine*, a term that was first used by Mandelbrot (1985). If a fractal structure is self-affine, it can no longer be described by a single fractal dimension $D_f$, and in fact if one utilizes any of the methods of estimating a fractal dimension that were in Sections 1.1 and 1.2, then, the resulting fractal dimension would depend on the length scales over which the method is utilized.

A well-known example of a process that gives rise to a self-affine fractal is a marginally stable growth of an interface. For example, if water displaces oil in a porous medium, the interface between water and oil is a self-affine fractal. Well-known examples of man-made materials with rough and self-affine surfaces include thin films that are formed by molecular beam epitaxy. Among naturally-made surfaces that are rough and have self-affine properties are bacterial colonies, and pores and fractures of rock and other types of porous media. Many properties of such materials are described by a function $f(\mathbf{x})$ that also possesses a self-affine structure. For example, the surface height $h(x, y)$ at a lateral position $\mathbf{x}$ of a rough surface, e.g., the internal surface of a rock fracture, and the porosity distribution of rock along a well at depths $\mathbf{x}$, both have self-affine property. Self-affinity of many natural systems that are associated with Earth, such as various properties of natural rock, is quite understandable, since gravity plays a dominant role in one direction

but has very little effect in the other directions, hence generating anisotropy in the structure of rock. The interested reader is referred to Family and Vicsek (1991) for an excellent collection of articles which describe a wide variety of rough surfaces with self-affine properties.

Self-affine fractals that one encounters in practical situations are typically disordered, and thus their self-affinity is only in a statistical sense. For the problems that are of interest to us in this book, a disordered self-affine fractal can be thought of as the fluctuations about a straight line or a flat surface. Such fluctuations can generate rough self-affine curves or surfaces. If we consider the height difference between a pair of points $h(\mathbf{x}_1)$ and $h(\mathbf{x}_2)$ on a self-affine surface $h(\mathbf{x})$ that lie above or below points separated by a distance $x_1 - x_2 = x = |\mathbf{x}|$ on a flat reference surface (or line), then

$$\langle |h(\mathbf{x}_1) - h(\mathbf{x}_2)| \rangle \sim x^H, \tag{9}$$

where $H$ is called the Hurst exponent. One may generalize Eq. (9) to higher dimensions, and generate rough surfaces that are encountered in a variety of contexts, from surface of pores of a natural porous medium (see, for example, Sahimi, 1993b, 1995b, for comprehensive discussions) to fracture surface of heterogeneous materials (see Chapters 6 and 7), to thin films that are formed by a deposition process (see below).

## 1.4   Generation of Rough Surfaces: Fractional Brownian Motion

We now describe two fractal processes that are used for generating rough curves (1D profiles) and surfaces. The properties that these self-affine fractal processes possess may also be used as guides to better understanding of rough surfaces that one encounters in practical applications. In addition, because these stochastic processes generate fractal sets with long-range correlations, they have been widely used for modeling of a variety of phenomena in engineering and materials science in which the effect of long-range correlations is paramount. We first consider the 1D case, and define a stochastic process $B_H(t)$, called the fractional Brownian motion (fBm), by (Mandelbrot and Van Ness, 1968)

$$B_H(t) - B_H(0) = \frac{1}{\Gamma(H + 1/2)} \left[ \int_{-\infty}^{t} K(t - s) dB(s) \right], \tag{10}$$

where $t$ can be a spatial or temporal variable. Here $\Gamma(x)$ is the gamma function, $H$ is the Hurst exponent defined above, and the kernel $K(t - s)$ is given by

$$K(t - s) = \begin{cases} (t - s)^{H-1/2} & 0 \le s \le t \\ (t - s)^{H-1/2} - (-s)^{H-1/2} & s < 0. \end{cases} \tag{11}$$

It is not difficult to show that

$$B_H(bt) - B_H(0) \equiv b^H [B_H(t) - B_H(0)], \tag{12}$$

where "≡" means "statistically equivalent to." A remarkable property of fBm is that it generates correlations with *infinite* extent. To see this, consider the correlation function $C(t)$ of future increments $B_H(t)$ with past increments $-B_H(-t)$ which is defined by

$$C(t) = \frac{\langle -B_H(-t)B_H(t)\rangle}{\langle B_H(t)^2\rangle}. \tag{13}$$

It is straightforward to show that $C(t) = 2(2^{2H-1} - 1)$, *independent* of $t$. Moreover, the type of the correlations can be tuned by varying $H$. If $H > 1/2$, then fBm displays *persistence*, i.e., a trend (for example, a high or a low value) at $t$ is likely to be followed by a similar trend at $t + \Delta t$, whereas if $H < 1/2$, then fBm generates *antipersistence*, i.e., a trend at $t$ is not likely to be followed by a similar trend at $t + \Delta t$. For $H = 1/2$ the past and future are not correlated, and thus the increments in $B_H(t)$ are completely random and uncorrelated. Thus, varying $H$ allows us to generate infinitely long-range correlations or anticorrelations.

We can generalize the above 1D fBm to 2D or 3D. Hence, if we consider two arbitrary points $\mathbf{x}$ and $\mathbf{x}_0$ in 2D or 3D space, the fBm is defined by

$$\langle [B_H(\mathbf{x}) - B_H(\mathbf{x}_0)]^2\rangle \sim |\mathbf{x} - \mathbf{x}_0|^{2H}. \tag{14}$$

Figure 1.2 presents 1D and 2D rough profiles and surfaces generated by fBm. The increments in fBm are stationary but not ergodic. The variance of a fBm for a large enough array is divergent (i.e., the variance increases with the size of the array without bounds). Its trace in $d$ dimensions is a self-affine fractal with a *local* fractal dimension $D_f = d + 1 - H$. Fractional Brownian motion is not differentiable at any point, but by smoothing it over an interval one can obtain its approximate numerical derivative which is called *fractional Gaussian noise* (fGn), a 1D example of which is shown in Figure 1.3, which should be compared with its counterpart in Figure 1.2. We should point out that the correlation function $C(r)$ of a fBm is given by

$$C(r) - C(0) \sim r^{2H} \tag{15}$$

so that, as long as $H > 0$ (which are the only physically-acceptable values of $H$), the correlations *increase* as $r$ does.

Efficient and accurate generation of a $d$-dimensional array that follows the statistics of a fBm is not straightforward. Rambaldi and Pinazza (1994) describe a numerical algorithm based on Eqs. (10) and (11). In addition to their method, there are at least three other techniques for numerically generating a fBm array with a given Hurst exponent $H$ (Mehrabi *et al.*, 1997) which we now describe.

### 1.4.1  The Power-Spectrum Method

A convenient way of representing a stochastic function is through its power spectrum $S(\boldsymbol{\omega})$, the Fourier transform of its covariance. The power spectrum of a

FIGURE 1.2. Examples of one- and two-dimensional rough profiles and surfaces generated by the fractional Brownian motion with various Hurst exponents $H$.

$d$-dimensional fBm is given by

$$S(\boldsymbol{\omega}) = \frac{a_d}{\left(\sum_{i=1}^{d} \omega_i^2\right)^{H+d/2}}. \tag{16}$$

where $\boldsymbol{\omega} = (\omega_1, \cdots, \omega_d)$ is the Fourier-transform variable, and $a_d$ is a $d$-dependent constant. The spectral representation (16) also allows us to introduce a cutoff length scale $\ell_{co} = 1/\omega_{co}$ such that

$$S(\boldsymbol{\omega}) = \frac{a_d}{\left(\omega_{co}^2 + \sum_{i=1}^{d} \omega_i^2\right)^{H+d/2}}. \tag{17}$$

The cutoff $\ell_{co}$ allows us to control the length scale over which the spatial properties of a system are correlated (or anticorrelated). Thus, for length scales $L < \ell_{co}$ the properties preserve their correlations (anticorrelations), but for $L > \ell_{co}$ they become random and uncorrelated. Note that the power spectrum of fGn in, for

FIGURE 1.3. An example of one-dimensional fractional Gaussian noise.

example, 1D is given by

$$S(\omega) = \frac{b_d}{\omega^{2H-1}}, \tag{18}$$

where $b_d$ is another $d$-dependent constant. The spectral representation of fBm (and fGn) provides a convenient method of generating an array of numbers that follow the fBm statistics, using a fast Fourier transformation (FFT) technique. In this method, one first generates random numbers, distributed either uniformly in [0,1), or according to a Gaussian distribution with random phases, and assigns them to the sites of a $d$-dimensional lattice. In most cases the linear size $L$ of the lattice is a power of 2, but the only requirement is that $L$ can be partitioned into small prime numbers, so that a FFT algorithm can be used. One must also keep in mind that, since the variance $\sigma^2$ of a fBm increases with the size $L$ of the array, generating a fBm array with a given variance requires selecting an appropriate $L$. In any case, the Fourier transformation of the resulting $d$-dimensional array of the numbers is then calculated numerically, the resulting numbers are multiplied by $\sqrt{S(\omega)}$, and the results then inverse Fourier transformed back into the real space. The array so obtained follows the statistics of a fBm with the desired long-range correlations and the specified value of $H$. To avoid the problem associated with the periodicity of the numbers arising as a result of their Fourier transforming, one must generate

the array using a much larger lattice size than the actual size that is to be used in the analysis, and use the central part of the array (or lattice).

### 1.4.2   Successive Random Additions

In the successive random addition method (Voss, 1985) one begins with the two end points in the interval [0,1], and assigns a zero value to them. Then Gaussian random numbers $\Delta_0$ with a zero mean and unit variance are added to these values. In the next stage, new points are added at a fraction $r$ of the previous stage by interpolating between the old points (by either linear or spline interpolation), and Gaussian random numbers $\Delta_1$ with a zero mean and variance $r^{2H}$ are added to the new points. Thus, given a sample of $N_i$ points at stage $i$ with resolution $\lambda$, stage $i + 1$ with resolution $r\lambda$ is determined by first interpolating the $N_{i+1} = N_i/r$ new points from the old points, and then Gaussian random numbers $\Delta_i$ with a zero mean and variance $r^{2(i-1)H}$ are added to all of the new points. At stage $i$ with $r < 1$, the Gaussian random numbers have a variance

$$\sigma_i^2 \sim r^{2iH}. \tag{19}$$

This process is continued until the desired length of the data array is reached. Typically $r = 1/2$ is used to generate a fBm.

   The problem with this method is that the points that are generated in earlier generations are not statistically equivalent to those generated later. To remedy this, one can add, during the $n$th stage of the process, a random Gaussian displacement with a variance $r^{2(n-1)H}$ to *all* of the points. This of course increases the computation time (it roughly doubles it). Moreover, if one is interested in generating a fBm array with a very wide range, one may start the process by assigning a Gaussian random number with a variance $2^{2H}$ to one end of the [0, 1] interval. The generalization of this method to higher dimensions is straightforward.

### 1.4.3   The Weierstrass–Mandelbrot Algorithm

In the Weierstrass–Mandelbrot (WM) method (Voss, 1985) one first divides the interval [0,1] into $n - 1$ equally-spaced subintervals, where $n$ is the size of the data array that one wishes to generate, and assigns zero value to all the points in the interval. Then, to point $i$ at a distance $x_i$ from the origin one adds a random number generated by the Weierstrass function defined by

$$\mathcal{W}(x_i) = \sum_{j=-\infty}^{\infty} C_j r^{jH} \sin(2\pi r^{-j} x_i + \phi_j) \tag{20}$$

where $C_j$ and $\phi_j$ are random numbers distributed according to Gaussian and uniform distributions, respectively, and $r$ is a measure of the distance between the frequencies, which is usually chosen to be small, e.g., $r = 0.9$. The variance of $C_j$ is proportional to $r^{2jH}$, and the random phases $\phi_j$ are distributed uniformly on $[0, 2\pi]$. Usually, the infinite series in Eq. (20) is approximated by a finite number of terms, but the number of terms included in the series must be large to

ensure accuracy. For example, in our own work we have used up to 140 terms in $-70 \le j \le 70$ to obtain accurate results. The power spectrum of the data array generated by the WM method is discrete and does not contain all the frequencies. However, it is still proportional to $\omega^{-(2H+1)}$, in agreement with Eq. (16).

## 1.5   Scaling Properties of Rough Surfaces

How do we characterize a rough self-affine surface, either generated synthetically (numerically) or by a physical process, such as fracturing of a material? We define a height correlation function $C_n(\mathbf{x})$ by

$$C_n(\mathbf{x}) = \langle |h(\mathbf{x}_0 + \mathbf{x}) - h(\mathbf{x}_0)|^n \rangle^{1/n}, \tag{21}$$

where $h(\mathbf{x})$ is the height of the surface at a transverse position $\mathbf{x}$ above a reference surface that can be a smooth, coarse-grained approximation to the rough surface, and the averaging is over all the initial $\mathbf{x}_0$. The choice of the reference surface can be tricky. For example, if the rough surface has been grown from a planar substrate, then a plane parallel to the substrate and in a coordinate system that moves with the rough surface can be taken to be the reference plane. In any case, it has been found for many rough surfaces that

$$C_n(x) = \langle C_n(\mathbf{x}) \rangle_{|\mathbf{x}|=x} \sim x^{H(n)}, \tag{22}$$

where the averaging is taken with respect to all the origins $\mathbf{x}_0$ in the smooth reference plane. In most cases, the exponents $H(n)$ take on the same value $H$ for all $n$, but there are also some exceptions to this, as discussed by Barabási and Vicsek (1990). A surface with correlation function (22) is a self-affine fractal over the range of length scales in which $C_n(x)$ is computed. Typically, the height correlation function $C_2(x)$, denoted simply as $C(x)$, has been utilized for estimating $H$, and has proven to be a very robust and accurate method.

   In practical applications, such as analyzing rough fracture surfaces, the self-affinity of the surface is bounded by an upper correlation length $\xi^+$ and a lower correlation length $\xi^-$ in both the horizontal ($\parallel$) and vertical ($\perp$) directions. That is, self-affine behavior is restricted to the ranges, $\xi_\parallel^- < \delta x < \xi_\parallel^+$ and $\xi_\perp^- < \delta h < \xi_\perp^+$. Because of the self-affinity property we must have

$$\frac{\xi_\perp^+}{\xi_\perp^-} = \left( \frac{\xi_\parallel^+}{\xi_\parallel^-} \right)^H. \tag{23}$$

The correlation function $C_n(x)$ satisfies a general scaling equation given by

$$C_n(x) = x^H F_n(x/\xi_\parallel^+, x/\xi_\parallel^-). \tag{24}$$

For $x \gg \xi_\parallel^-$, such that $x/\xi_\parallel^- \to \infty$, scaling equation (24) simplifies to

$$C_n(x) = x^H f_n(x/\xi_\parallel^+), \tag{25}$$

where the scaling function $f(y)$ has the properties that $f(y) = c$ for $y \ll 1$ and $f(y) \sim y^{-H}$ for $y \gg 1$, where $c$ is a constant of order unity. Hereafter, we delete the superscripts and use $\xi_\parallel$ and $\xi_\perp$ for the upper cutoff length scales.

If the rough self-affine fractal surface is growing with the process time $t$ as in, for example, deposition on a flat surface, then one must define a more general correlation $C_n(x, t)$ in a manner similar to that used for $C_n(x)$, namely,

$$C_n(\mathbf{x}, t) = \left\langle [|h(\mathbf{x}_0 + \mathbf{x}, t + t') - h(\mathbf{x}_0, t')|]^n \right\rangle^{1/n}, \tag{26}$$

where the averaging is over all the initial position $\mathbf{x}_0$ and times $t'$. Then, due to self-affinity of the surface, the correlation function $C_n(x, t)$ has the property that

$$C_n(bx, b^z t) = b^\alpha C_n(x, t). \tag{27}$$

Similar to $C_2(\mathbf{x})$, one usually constructs $C_2(\mathbf{x}, t)$ and attempts to extract from it information about the surface. Under the dynamic conditions in which a rough and self-affine surface grows, there exists a time scale $t_c$ over which the time correlations are important. For rough surfaces that begin growing from a smooth surface, it has been found in most cases that $\xi_\perp$ and $\xi_\parallel$ satisfy the following power laws,

$$\xi_\perp \sim t^\beta, \quad t \ll t_c \tag{28}$$

$$\xi_\parallel \sim t^{1/z}, \quad t \ll t_c \tag{29}$$

where $t$ is either the time (for a growing rough surface) or the surface's mean thickness. For $t \gg t_c$ the magnitude of $\xi_\parallel$ saturates, $\xi_\parallel = L$. The quantity $z$ is called the *dynamical exponent* of the surface, while $\beta$ is called the *growth exponent*. The quantities $\xi_\perp$ and $\xi_\parallel$ are actually related to each other by

$$\xi_\perp \sim \xi_\parallel^\alpha. \tag{30}$$

$\alpha$ is called the *roughness exponent*. Although we are not aware of an experimental realization of a case for which $\alpha$ and the Hurst exponent $H$ are different, we keep both $\alpha$ and $H$ to make our discussion as general as possible.

The roughness of a dynamic, growing surface is characterized by the width $w(L)$ defined as,

$$w(L) = \left\langle [h(x) - \langle h \rangle_L]^2 \right\rangle^{1/2}, \tag{31}$$

where $h(x)$ is, as before, the height of the surface at position $x$, and $\langle h \rangle_L$ is its average over a horizontal segment of length $L$ (normalized by the "volume" $L^{d-1}$). According to the dynamic scaling theory of Family and Vicsek (1985) for growing rough surfaces, one has the following dynamic scaling equation

$$h(x) - \langle h \rangle_L \sim t^\beta f(x/t^{\beta/\alpha}), \tag{32}$$

where $\alpha$ and $\beta$, the two exponents defined above, satisfy the following scaling relation

$$\alpha + \frac{\alpha}{\beta} = 2, \tag{33}$$

and the scaling function $f(u)$ has the properties that $|f(u)| < c$ for $u \gg 1$, and $f(u) \sim L^\alpha f(Lu)$ for $u \ll 1$, where $c$ is a constant. Note that the ratio $\alpha/\beta$ can be replaced by the dynamical exponent $z$. It is then straightforward to show that

$$w(L, t) \sim L^\alpha g(t/L^{\alpha/\beta}), \tag{34}$$

where $g(u)$ is a universal scaling function. Note also that $w(L, t)$ is a measure of the correlation length $\xi_\perp$ along the direction of growth. As the rough surface grows, the wavelength of the spatial fluctuations and the length over which the fluctuations are correlated both grow with time. However, the length $L$ is the maximum spatial extent to which the correlations can grow in the $d - 1$ dimensions along the surface. When the correlations reach this scale, they cannot extend further, and therefore the rough surface reaches a steady-state which is characterized by a constant width. Then, the surface is scale invariant and the saturation value $w(L, \infty)$ is expected to have a power-law dependence on $L$:

$$w(L, \infty) \sim L^\alpha. \tag{35}$$

The correlation time $t_c$ also scales with $L$ as

$$t_c \sim L^{\alpha/\beta} \sim L^z, \tag{36}$$

Equation (34) indicates that, if one plots $w/L^\alpha$ versus $t/L^{\alpha/\beta}$, then, due to the universality of $g(u)$, all the results for various $t$ and $L$ should collapse onto a single universal curve [representing the scaling function $g(u)$]. Figure 1.4 presents such a data collapse for a rough surface grown by a ballistic deposition process (Vold, 1963). In the simplest version of ballistic deposition, one begins with a line of $L$, selects at random a horizontal line above the line of particles, and places a



FIGURE 1.4. Data collapse for a rough surface grown by ballistic deposition (courtesy of Ehsan Nedaaee Oskoee).

FIGURE 1.5. An example of a rough surface grown by ballistic deposition (courtesy of Ehsan Nedaaee Oskoee).

particle there. The particle is then allowed to fall along a straight line vertically downward. When the particle touches the original particles, it sticks to them and becomes part of the particle pile. A large deposit is then grown by repeating this procedure. Extensive numerical simulations indicate that the deposit is compact and non-fractal, but its surface is rough and self-affine. An example is shown in Figure 1.5.

## 1.6   Modeling of Growth of Thin Films with Rough Surface

How can we describe the growth of a rough surface? If the surface is characterized by a single-valued height function $h(\mathbf{x}, t)$, then in general we can describe the growth of the surface by the following equation

$$\frac{\partial h}{\partial t} = \mathcal{R}(\mathbf{x}, t) + \mathcal{N}(\mathbf{x}, t), \tag{37}$$

where $\mathcal{R}(\mathbf{x}, t)$ represents all the various (deterministic) physical factors that contribute to the rate of growth of $h(\mathbf{x}, t)$, and $\mathcal{N}$ represents the noise or randomness in

the growth of the rough surface. However, because of various constraints that are imposed by the physics of growing a rough surface, the set of acceptable functions $\mathcal{R}(\mathbf{x}, t)$ is limited. Some of these constraints are as follows (Barabási and Stanley, 1995).

(1) The growth of the surface should be independent of where $h = 0$ is defined, i.e., it should be invariant under the transformation $h \rightarrow h + \delta h$. Therefore, $\mathcal{R}$ cannot depend explicitly on $h$, but should be built from such terms as $\nabla^n h$ (with $n = 1, 2, \cdots$).

(2) The equation must have rotation and inversion symmetry with respect to the direction of the growth, implying that it cannot contain odd-order derivatives in the coordinates, such as $\nabla h$ and $\nabla(\nabla^2 h)$.

(3) The equation must be invariant under time translation $t \rightarrow t + \delta t$, which means that $\mathcal{R}$ cannot depend explicitly on $t$. It should also be translationally invariant in the direction perpendicular to the growth direction, and therefore $\mathcal{R}$ cannot contain terms that are explicit in $\mathbf{x}$.

(4) Since the fluctuations in the rough surface must be similar with respect to the mean position of the surface—the so-called up-down symmetry ($h \rightarrow -h$ invariance)—the equation cannot contain terms such as $(\nabla h)^n$ with $n$ being an even number. However, this symmetry can be broken if there exists a driving force $\mathcal{F}$, perpendicular to the rough surface, which selects a particular direction for the growth of the surface. The existence of this driving force is a necessary but not sufficient condition for breaking this symmetry.

Therefore, the most general form of the equation that describes the growth of a rough surface is given by

$$\frac{\partial h}{\partial t} = \nabla^2 h + \nabla^4 h + \cdots + (\nabla^2 h)(\nabla h)^2 + \cdots + (\nabla^{2k} h)(\nabla h)^{2j} + \mathcal{N}(\mathbf{x}, t).$$
(38)

To investigate the scaling properties of a growing surface, we consider the hydrodynamic limit, $t \rightarrow \infty$ and $\mathbf{x} \rightarrow \infty$. In this limit, the higher-order derivatives of $h$ are much smaller than the lowest-order one. Consider, as examples, $\nabla^2 h$ and $\nabla^4 h$. Writing $\mathbf{x} \rightarrow \mathbf{x}' \equiv b\mathbf{x}$, we must have $h \rightarrow h' \equiv b^\alpha h$, and thus, $\nabla^2 h \rightarrow \nabla'^2 h' \equiv b^{\alpha-2}\nabla^2 h$ and $\nabla^4 h \rightarrow \nabla'^4 h' \equiv b^{\alpha-4}\nabla^4 h$. In the limit $b \rightarrow \infty$, $\nabla^4 h$ decays much faster than $\nabla^2 h$ and can therefore be neglected.

Given such considerations, the simplest possible equation has the following form

$$\frac{\partial h}{\partial t} = \mathcal{D}\nabla^2 h + \mathcal{N}(\mathbf{x}, t),$$
(39)

which was proposed by Edwards and Wilkinson (1982). In most cases, the noise term has been assumed to be Gaussian:

$$\langle \mathcal{N}(\mathbf{x}, t)\mathcal{N}(\mathbf{x}', t)\rangle = 2A\delta(\mathbf{x} - \mathbf{x}')\delta(t - t'),$$
(40)

where $A$ is the amplitude of the noise. Equation (40) implies that there is no correlation in space or time, since the average $\langle \mathcal{N}(\mathbf{x}, t)\mathcal{N}(\mathbf{x}', t)\rangle$ vanishes (except,

of course, at $\mathbf{x} = \mathbf{x}'$ and $t = t'$). The Edwards–Wilkinson model, which satisfies the four constraints described above, can be solved exactly (this is made possible by the linearity of the equation). One obtains, $\alpha = \frac{1}{2}(2 - d)$, $\beta = \frac{1}{2}\alpha$, and hence $z = 2$. This model describes the growth of a surface by *random* deposition of particles on a growing surface, starting from a flat surface, in which, upon landing on the growing surface, the particles diffuse on the surface until they find a point with the lowest height at which they stop. Note that the Edwards–Wilkinson equation predicts that for $d = 2$ (growth on a 2D surface) $\alpha = 0$, which should be interpreted as implying a logarithmic dependence of the width $w$ on $L$, i.e., $w(L, \infty) \sim \ln L$.

The growth of a variety of thin films with rough, self-affine surfaces, such as those that are formed by ballistic deposition, and the dynamical scaling of the height and width of such surfaces, are described by the stochastic differential equation proposed by Kardar, Parisi, and Zhang (KPZ) (1986):

$$\frac{\partial h}{\partial t} = \mathcal{D}\nabla_T^2 h + \frac{1}{2}\mathrm{v}|\boldsymbol{\nabla}h|^2 + \mathcal{N}(\mathbf{x}, t), \tag{41}$$

where v is the growth velocity perpendicular to the surface, and $\mathcal{D}$ is a diffusivity. Equation (41) satisfies the first three constraints listed above, but violates the fourth constraint since, for example, in ballistic deposition there is lateral growth of the surface (i.e., the growth occurs in the direction of *local* normal to the growing surface), and this is equivalent to having a net driving force $\mathcal{F}$. The lateral growth is represented by the nonlinear term $\frac{1}{2}\mathrm{v}|\boldsymbol{\nabla}h|^2$. To see how this term arises, suppose that a new particle is added to the growing surface. If the surface grows in the direction of local normal to the surface, then its growth $\delta h$ is given by, $\delta h = [(\mathrm{v}\delta t)^2 + (\mathrm{v}\delta t \boldsymbol{\nabla}h)^2]^{1/2} = \mathrm{v}\delta t[1 + (\boldsymbol{\nabla}h)^2]^{1/2}$. Thus, if $|\boldsymbol{\nabla}h| \ll 1$, one must add a term $\frac{1}{2}\mathrm{v}(\boldsymbol{\nabla}h)^2$ to the Edwards–Wilkinson equation. In the literature one often finds that $\sigma$ is used instead of the diffusivity $\mathcal{D}$, and is referred to as a "surface tension," since $\nabla^2 h$ tends to smoothen the surface, as does a surface tension. However, we prefer to use $\mathcal{D}$ as the term $\mathcal{D}\nabla_T^2 h$ represents a diffusion process that arises when the depositing particles land on the growing surface, diffuse on the surface, and only stop when they find the point with the *lowest* height. This diffusion process also helps smoothen the growing surface (and counter the effect of lateral growth, represented by the nonlinear term $\frac{1}{2}\mathrm{v}|\boldsymbol{\nabla}h|^2$, which tends to roughen the surface). Kardar *et al.* (1986) considered the case in which the noise was assumed to be Gaussian with the correlation function (40). For their model, it has been proposed (Kim and Kosterlitz, 1989; Hentschel and Family, 1991) that for a $d$-dimensional surface,

$$\alpha = \frac{2}{d + 2}, \tag{42}$$

$$\beta = \frac{1}{d + 1}, \tag{43}$$

and therefore the dynamical exponent $z$ is given by, $z = 2(d + 1)/(d + 2)$. Equations (42) and (43) are not exact, but provide accurate estimates of $\alpha$ and $\beta$ (and hence $z$). Note that the KPZ equation predicts that $z = 2$ only when $d \to \infty$.

Another stochastic equation was proposed by Koplik and Levine (1985)

$$\frac{\partial h}{\partial t} = \mathcal{D}\nabla_T^2 h + v + A\mathcal{N}(\mathbf{r}, h), \tag{44}$$

a linear equation in which the term representing the noise is more complex than the corresponding term in the KPZ equation. For this model, the numerical simulations indicate that $\alpha(d = 2) \simeq 3/4$, which should be compared with that of the KPZ surfaces, $\alpha = 2/3$. The growth of a rough surface can sometimes stop because it is *pinned*. To see how the pinning occurs, consider Eq. (44) in zero transverse dimension:

$$\frac{\partial h}{\partial t} = v + A\mathcal{N}(h). \tag{45}$$

If $v > A\mathcal{N}_{max}$, where $\mathcal{N}_{max}$ is the maximum value of $\mathcal{N}$, then $\partial h/\partial t > 0$, and the surface always moves with a velocity that fluctuates around v. If, however, $v < A\mathcal{N}_{max}$, the surface will eventually arrive at a point where $v + A\mathcal{N} = 0$, and will be pinned down. Therefore, for a fixed v there must be a pinning transition at some finite value of $A$.

## 1.7    Measurement of Roughness Exponent

The numerical value of the Hurst exponent $H$ or the roughness exponent $\alpha$ is not enough for characterizing the roughness of a surface. It only indicates how the roughness (or the variance in the height) varies as the transverse length scale, over which it is measured, changes. A complete characterization of the rough surface would require not only $H$ or $\alpha$, but also the amplitudes of the height fluctuations as well as the transverse correlation lengths. One way of characterizing a rough surface is by measuring the width $w$ over a segment of size $\ell$ from the surface. Then for $\ell \ll \xi_\parallel$ we must have

$$w(\ell) \sim \ell^H. \tag{46}$$

For $\ell \gg \xi_\parallel$ we must of course have $w = \xi_\perp$.

Another method of characterizing a rough surface is by the so-called *slit island* method (Mandelbrot *et al.*, 1984). In this method, the rough surface is coated with another material and then polished carefully parallel to the flat reference surface (described above) to reveal a series of horizontal cuts. As the coating material is removed, islands of the surface material appear in a sea of the coating material. With further removal of the coating material, the islands will grow and merge. If we consider a region of linear size $\ell$ and height fluctuations $\Delta h = h(\mathbf{x}) - \langle h \rangle_\ell$, the distribution $P(\Delta h)$ can be described by the following scaling law

$$P(\Delta h) = w(\ell)^{-1} f[\Delta h/w(\ell)], \tag{47}$$

where $w(\ell)$ is the width of the region. Since $w(\ell)$ follows Eq. (46), the implication is that the density $\rho(\ell)$ in a cross-section of size $\ell$ is given by

$$\rho(\ell) \sim \ell^{-H}. \tag{48}$$

Equation (48) suggests that the interface between the two materials, i.e., between the rough surface and the coating material, in the cross-sections parallel to the reference plane is a *self-similar* fractal with a fractal dimension

$$D_f = d - H,$$

(49)

where $d$ is the Euclidean dimensionality of the reference surface. Therefore, if the fractal dimension $D_f$ can be estimated independently, then the Hurst exponent $H$ can also be evaluated. Typically, the islands that appear have a surface area distribution $n_S$ such that

$$n_S \sim S^{-\tau},$$

(50)

where $n_S$ is the number of islands with areas $S$ in the range $[S - \frac{1}{2}\Delta S, S + \frac{1}{2}\Delta S]$. The exponent $\tau$ is related to the fractal dimension $D_f$ through the following equation

$$\tau = \frac{1}{d}\left(D_f + d\right),$$

(51)

so that measurement of the islands' areas yields $D_f$, from which the Hurst exponent $H$ can be estimated.

The third method of analyzing a rough, self-affine surface is through its power spectrum which, in $d$ dimensions, is given by Eq. (16). However, as Hough (1989) pointed out, interpreting a power-law power spectrum is not without difficulties, and thus one must be careful in using such an analysis. In particular, a power-law power spectrum might also be the result of a non-stationary and non-fractal system. We will come back to this issue in Chapters 6 and 7, where we describe fracture surface of materials which are typically very rough.

## Summary

An important characteristics of morphology of disordered multiphase materials is the structure of their surface, and in particular their surface roughness. The concepts of modern statistical physics of disordered media can now quantify the roughness in terms of self-affine fractals, and the roughness or Hurst exponent. The dynamics of growth of such surfaces can also be described by dynamical scaling, discrete models of material growth, and suitable continuum differential equations. Moreover, fractal geometry, and the associated power-law correlation functions, point to the fundamental role of length scale and long-range correlations in the macroscopic homogeneity of a heterogeneous material. If the largest relevant length scale of the material, e.g., its linear size, is less than the length scale at which it can be considered homogeneous, then the classical equations that describe transport processes in the material must be fundamentally modified.

# Part I

# Effective Properties of Heterogeneous Materials with Constitutive Nonlinearities

# 2
# Nonlinear Conductivity and Dielectric Constant: The Continuum Approach

## 2.0   Introduction

The main focus of Volume II is on nonlinear properties of heterogeneous materials. In general, there are two fundamental classes of nonlinearity that one may encounter in disordered materials:

(1)   One class of nonlinear materials is described by what we call *constitutive nonlinearity*, which is one in which the basic *local* constitutive law that expresses the relation between the flux (of current, force, etc.) and the potential (voltage, stress, etc.) gradient is nonlinear. As a result, the macroscopic behavior of such materials must also be described by nonlinear transport equations. In particular, the effective transport properties of such materials are nonlinear in the sense of being functions of the external potential gradient. Such materials are of great practical importance, since, for example, one may be able to design new nonlinear optical materials by tuning their nonlinear response which can be achieved by, for example, changing the volume fraction of their constituents. For example, it has been suggested that strong local field effects, such as the large local field at the surface plasmon resonance frequency of a metallic inclusion, may lead to enhanced nonlinear response in a heterogeneous material. Constitutive nonlinearity is the subject of this and the next two chapters. Even within this restricted class of nonlinear materials, one may imagine a very large number of nonlinear constitutive equations (similar to those that have been proposed, for example, for polymeric fluids). Therefore, while we describe in this chapter results for general constitutive nonlinearity, their application is restricted mostly to strongly nonlinear materials, i.e., those that are described by a *power-law* relation between the flux and the current. In the next two chapters we will also describe the macroscopic behavior of nonlinear materials that can be described by a few other types of nonlinear constitutive equations, for which considerable progress has been made, and a comparison between the theoretical predictions and the experimental data is possible.

(2)   In the second class of nonlinearities, a material is characterized by *thresholds* in the (local as well as macroscopic) potential gradient. Then, depending on the physics of the phenomenon under study, one of the following two scenarios may arise.

(i) The transport properties of the material vanish below the threshold, but above the threshold the material behaves linearly (or, possibly, exhibits constitutive nonlinearity) and possesses non-zero effective transport properties. For example, consider a resistor network in which each bond is insulating if the voltage drop between its two ends is less than a threshold value, but becomes conducting (either linearly or nonlinearly) if the voltage drop exceeds a threshold. An example of a material to which such a model is directly relevant is foam. As described in Chapter 9 of Volume I, foams behave both as solid materials (in the sense of exhibiting an elastic response when exposed to an external stress or strain), and as a fluid when the applied stress that they are exposed to reaches a threshold value. Therefore, foams do not flow if the stress applied to them is less than the threshold. As a result, if we consider, for example, flow of foams in a porous medium (which is usually modeled as a network of tubes), there would be no macroscopic flux of foams unless the pressure gradient applied to the porous medium exceeds a threshold. We must, however, point out that this type of threshold behavior is *not* the same as that of a percolation system below and above the percolation threshold, i.e., this threshold behavior is *not* a geometrical effect, although, as we will show in Chapter 3, there are certain similarities between the two types of phenomena.

(ii) The second scenario arises when the material behaves linearly (or, possibly, exhibits constitutive nonlinearity) if the applied potential gradient is *less than* a threshold, but exhibits highly nonlinear properties when the threshold is exceeded. Well-known examples of this type of phenomenon are brittle fracture and dielectric breakdown of solid materials, phenomena that will be studied beginning with Chapter 5.

Compared to linear systems, the number of studies in which an attempt has been made to obtain estimates of the effective nonlinear properties is small. This is particularly true in the context of continuum models of disordered materials. Discrete models have received much more attention, and will be described and discussed in Chapter 3. To our knowledge, Marcellini (1978) was perhaps the first to undertake a systematic study of effective transport properties of nonlinear materials, and attempted to estimate their effective dielectric constant. He considered a two-phase composite in which one phase had a constant dielectric constant, while the dielectric constant of the second phase, that consisted of spherical inclusions, was a function of the local electric field. The particles were arranged either at random or in a periodic manner, similar to the periodic models that were described and analyzed in detail in Chapter 4 of Volume I. Miksis (1983) obtained slightly more general results for the effective properties of periodic arrays, and random distributions of nonlinear spherical inclusions in a linear matrix. The methods of Marcellini and Miskis were more or less straightforward generalization of those described in Chapter 4 of Volume I, and hence need not be described again. Willis (1986) applied the approach of Hill (1963) (see below; see also Chapter 7 of Volume I for more details) to nonlinear dielectrics. In terms of deriving rigorous bounds for

the effective nonlinear electrical conductivity and dielectric constant, Talbot and Willis (1985, 1987, 1994) and Willis (1986) proposed extensions of the Hashin–Shtrikman variational principles (Hashin and Shtrikman, 1962a,b, 1963) (see also Chapters 4 and 7 of Volume I) to nonlinear heterogeneous materials. In a series of papers, Ponte Castañeda and co-workers (Ponte Castañeda, 1992b, 1998; Ponte Castañeda and Kailasam, 1997) analyzed the effective nonlinear conductivity and dielectric constant of two-phase heterogeneous materials using two different techniques. One of the methods is exact to first-order in contrast between the properties of the two phases, and is capable of delivering rigorous lower bounds and approximate estimates for the upper bounds (*not* the upper bounds themselves), while the second method is exact to second order in the contrast between the phases' properties. To our knowledge, their work is the most advanced attempt in the area of continuum description of the effective nonlinear conductivity and dielectric constant of disordered materials, and is described in detail in this chapter.

## 2.1    Variational Principles

Volume I of this book should have made it abundantly clear that the effective linear properties of heterogeneous materials are not characterized by simple averages of the properties of the constituent phases, weighted, for example, by their respective volume fractions. In fact, in addition to the volume fractions, the effective properties depend in general on certain microstructural parameters which are themselves functions of the volume fractions. The same is true about nonlinear effective properties of disordered materials. Moreover, due to the nonlinearity, a proper definition of the effective properties is even more important than the linear case because, for example, nonlinear effective properties may exhibit sensitive dependence on the boundary conditions.

Consider a heterogeneous dielectric material that occupies a region $\Omega$ in space. The nonlinear constitutive behavior of the material may be characterized in terms of an electric energy-density function $w(\mathbf{x}, \mathbf{E})$ which depends on the position $\mathbf{x}$ and the electric field $\mathbf{E}(\mathbf{x})$, such that the electric displacement field $\mathbf{D}(\mathbf{x})$ is given by

$$\mathbf{D}(\mathbf{x}) = \frac{\partial}{\partial \mathbf{E}} w(\mathbf{x}, \mathbf{E}). \tag{1}$$

Furthermore, if one assumes local isotropy, then, $w(\mathbf{x}, \mathbf{E}) = e(\mathbf{x}, E)$, where $e : \Omega \times R \to R$ is continuous, convex and coercive (in the sense that, $w \to \infty$ as $E \to \infty$) which satisfies the conditions, $e(\mathbf{x}, \mathbf{E}) \geq 0$ and $e(\mathbf{x}, 0) = 0$. Here, $R$ is the set of the extended real numbers.

The effective constitutive behavior of the heterogeneous material is then defined by

$$\langle \mathbf{D} \rangle = \frac{\partial}{\partial \mathbf{E}} \mathcal{H}_e(\mathbf{E}), \tag{2}$$

where $\langle \cdots \rangle$ denotes an spatial average. We should keep in mind that the effective behavior of the heterogeneous material, as characterized by the energy functional $\mathcal{H}_e(\langle \mathbf{E} \rangle)$, may in general be anisotropic, even if the material's phases themselves are isotropic. In principle, $\mathcal{H}_e(\langle \mathbf{E} \rangle)$ is determined by solving the usual electrostatic problem on $\Omega$, defined by, $\nabla \times \mathbf{E} = \mathbf{0}$, and $\nabla \cdot \mathbf{D} = 0$, subject to a uniform boundary condition, $\varphi = -\langle \mathbf{E} \rangle \cdot \mathbf{x}$ on the external surface of $\Omega$, where $\varphi$ is the electrostatic potential defined by, $\mathbf{E} = -\nabla \varphi(\mathbf{x})$ in $\Omega$. This boundary condition ensures that the average of the electric field is in fact $\langle \mathbf{E} \rangle$, in the sense that

$$\langle \mathbf{E} \rangle = \int_\Omega \mathbf{E}(\mathbf{x}) \, d\mathbf{x}. \tag{3}$$

Moreover, the average displacement field is defined by a similar relation:

$$\langle \mathbf{D} \rangle = \int_\Omega \mathbf{D}(\mathbf{x}) \, d\mathbf{x}, \tag{4}$$

so that one obtains the effective energy $\mathcal{H}_e$ that evaluates the pertinent energy functional for the heterogeneous material,

$$\mathcal{H}_e(\mathbf{E}) = \int_\Omega w(\mathbf{x}, \mathbf{E}) \, d\mathbf{x}, \tag{5}$$

at the actual electric field solving the electrostatic problem for a given microstructure. Due to the complexity of the morphology of real materials, it is not practical to solve the electrostatic problem. For this reason, variational formulations of the problem based on the minimum energy and minimum complementary-energy principles provide useful alternative routes for analyzing the problem. Thus, let us state these principles here (which were also utilized in Volume I for obtaining estimates of effective linear properties).

According to the minimum energy principle, expressed in terms of the energy functional $\mathcal{H}$, one can obtain the following expression for the effective energy $\mathcal{H}_e$ of a heterogeneous material,

$$\mathcal{H}_e(\langle \mathbf{E} \rangle) = \min_{\mathbf{E} \in S_1} \mathcal{H}(\mathbf{E}), \tag{6}$$

where

$$S_1 = \{\mathbf{E} | \mathbf{E} = -\nabla \varphi(\mathbf{x}) \text{ in } \Omega, \text{ and } \varphi = -\langle \mathbf{E} \rangle \cdot \mathbf{x} \text{ on } \partial\Omega\}. \tag{7}$$

Note that, to guarantee the existence of the minimizer (6), certain conditions on the behavior of $w$ (or $e$) as $E \to \infty$ are required, which is why one assumes that $w$ is coercive. Moreover, strict convexity of $\mathcal{H}_e$ guarantees uniqueness of the solution, convexity of $w$ ensures that of $\mathcal{H}_e$, and if the fields are smooth enough, Eq. (6) will be equivalent to the original electrostatic problem defined above.

The second characteristic of the heterogeneous material is obtained from its complementary-energy function $\mathcal{H}_e^c$, defined in terms of the principle of minimum complementary energy:

$$\mathcal{H}_e^c(\langle \mathbf{D} \rangle) = \min_{\mathbf{D} \in S_2} \mathcal{H}^c(\mathbf{D}), \tag{8}$$

where

$$\mathcal{H}^c(\mathbf{D}) = \int_\Omega w^*(\mathbf{x}, \mathbf{D}) \, d\mathbf{x} \tag{9}$$

is the complementary energy functional, expressed in terms of

$$w^*(\mathbf{x}, \mathbf{D}) = \sup_{\mathbf{E}}\{\mathbf{E} \cdot \mathbf{D} - w(\mathbf{x}, \mathbf{E})\}, \tag{10}$$

with

$$S_2 = \{\mathbf{D}|\nabla \cdot \mathbf{D} = 0 \text{ in } \Omega, \text{ and } \mathbf{D} \cdot \mathbf{n} = \langle\mathbf{D}\rangle \cdot \mathbf{n} \text{ on } \partial\Omega\} \tag{11}$$

being the set of admissible electric displacement fields. Note that, if Eq. (3) is reinterpreted as a definition for the average electric field, then, one has

$$\langle\mathbf{E}\rangle = \frac{\partial}{\partial\langle\mathbf{D}\rangle}\mathcal{H}^c_e(\langle\mathbf{D}\rangle). \tag{12}$$

In general, it can be shown that

$$\mathcal{H}^c_e(\langle\mathbf{E}\rangle) \geq \mathcal{H}^{c*}_e(\langle\mathbf{E}\rangle). \tag{13}$$

The reason for the inequality (13) is related to the fact that definitions of $\mathcal{H}$ and $\mathcal{H}^c$ correspond to different boundary conditions on the heterogeneous material (Dirichlet versus Neumann conditions), hence leading to generally distinct effective energies. However, the strict equality holds in (13) if the composite can be homogenized, in the sense that it can be considered as homogeneous on a large enough scale. Finally, note that

$$w^*(\mathbf{x}, \mathbf{D}) = e^*(\mathbf{x}, D), \tag{14}$$

where $e^*$ is the convex polar function (Legendre transform) of $e$ and $D$ is the magnitude of $\mathbf{D}$.

Ponte Castañeda (1992b) proposed new variational principles in order to obtain upper and lower bounds and estimates for the effective energy functions of nonlinear materials. These variational principles are equivalent to the standard ones described above, under appropriate hypothesis on the energy-density function. The new variational principle is based on a change of variables $r = h(E)$, with $h : R^+ \rightarrow R^+$ ($R^+$ is the set of non-negative reals) given by $h(E) = E^2$. One than obtains a function $f : \Omega \times R^+ \rightarrow R^+$, such that

$$f(\mathbf{x}, r) = e(\mathbf{x}, E) = w(\mathbf{x}, \mathbf{E}), \tag{15}$$

has the same dependence on $\mathbf{x}$ as $e$ and $w$, and that it is continuous and coercive (but not necessarily convex) in $r$. Moreover, $f$ is a non-negative function satisfying, $f(\mathbf{x}, 0) = 0$. Then, if we define the Legendre transform (convex polar) of $f$ by

$$f^*(\mathbf{x}, p) = \sup_{r \geq 0}\{rp - f(\mathbf{x}, r)\}, \tag{16}$$

it follows that

$$f(\mathbf{x}, r) \geq \sup_{p \geq 0}\{rp - f^*(\mathbf{x}, p)\}. \tag{17}$$

Note that $\mathbf{x}$ is fixed in (16) and (17), and that the suprema are evaluated over the sets of non-negative $r$ and $p$, respectively. In addition, the right-hand side of (17) is the bipolar of $f$, which has the geometric interpretation of the convex envelope of $f$, and hence the inequality. The equality in (17) is achieved if $f$ is convex and continuous in $r$. Therefore, assuming that the energy function $w$ in (16) is such that $f$ is convex (note that convexity of $f$ implies that of $w$), one obtains from (17) the following representation for the local energy density function of the nonlinear heterogeneous material,

$$w(\mathbf{x}, \mathbf{E}) = \sup_{\epsilon^0 \geq 0} \{w^0(\mathbf{x}, \mathbf{E}) - v(\mathbf{x}, \epsilon^0)\}, \tag{18}$$

where [from Eq. (16)]

$$v(\mathbf{x}, \epsilon^0) = \sup_{\mathbf{E}} \{w^0(\mathbf{x}, \mathbf{E}) - w(\mathbf{x}, \mathbf{E})\}, \tag{19}$$

where $p$ has been identified with $\frac{1}{2}\epsilon^0$ and $r$ with $E^2$, in such a fashion that $w^0(\mathbf{x}, \mathbf{E}) = \frac{1}{2}\epsilon^0(\mathbf{x})E^2$ and $v(\mathbf{x}, \frac{1}{2}\epsilon^0) = f^*(\mathbf{x}, \epsilon^0)$. Thus, $w^0$ corresponds to the local energy-density function of a linear, heterogeneous comparison material with arbitrary (but not necessarily constant) non-negative dielectric constant $\epsilon^0(\mathbf{x})$. The minimum energy formulation of the variational principle follows by making use of the representation (18) in the classical minimum energy principle, and interchanging the order of the infimum in (16) and the supremum in (18). The result (Ponte Castañeda, 1992b) is the following theorem.

**Theorem 1:**  *Suppose that the local energy-density function $w$ of a given nonlinear heterogeneous material with isotropic phases satisfies condition (15) with $f$ a non-negative, continuous, coercive and convex function of $r = E^2$, with $f(\mathbf{x}, 0) = 0$. Then, the effective energy function of the nonlinear heterogeneous material $\mathcal{H}_e$ is determined by the variational principle,*

$$\mathcal{H}_e = \sup_{\epsilon^0(\mathbf{x}) \geq 0} \{\mathcal{H}_e^0(\mathbf{E}) - V(\epsilon^0)\}, \tag{20}$$

*where*

$$V(\epsilon^0) = \int_\Omega v[\mathbf{x}, \epsilon^0(\mathbf{x})] \, d\mathbf{x}, \tag{21}$$

*and $\mathcal{H}_e^0$ is the effective energy function of a linear heterogeneous comparison material with local energy function $w^0$, such that*

$$\mathcal{H}_e^0 = \min_{\mathbf{E} \in S_1} \int_\Omega w^0(\mathbf{x}, \mathbf{E}) \, d\mathbf{x}. \tag{22}$$

The complementary-energy formulation of the new variational principle follows in a similar fashion from the change of variables $s = h(D)$, where $h$ is the same function as before which induces a function $g : \Omega \times R^+ \to R^+$, such that

$$g(\mathbf{x}, s) = e^*(\mathbf{x}, D) = w^*(\mathbf{x}, \mathbf{D}), \tag{23}$$

where $g$ is continuous and coercive in $s$, and is a non-negative function such that,

$g(\mathbf{x}, 0) = e^*(\mathbf{x}, 0) = 0$. Then, if one defines the concave polar of $g$ by

$$g_*(\mathbf{x}, q) = \inf_{s \geq 0} \{sq - g(\mathbf{x}, s)\}, \tag{24}$$

it follows that

$$g(\mathbf{x}, s) \leq \inf_{q \geq 0} \{sq - g_*(\mathbf{x}, q)\}, \tag{25}$$

with the equality holding true if $g$ is concave. Assuming then that the complementary energy density function $w^*$ of the nonlinear heterogeneous material is such that $g$ is concave, it follows from (25) that

$$w^*(\mathbf{x}, \mathbf{D}) = \inf_{\epsilon^0 \geq 0} \{w^{0*}(\mathbf{x}, \mathbf{D}) + v(\mathbf{x}, \epsilon^0)\}, \tag{26}$$

where $q$ has been identified with $(2\epsilon^0)^{-1}$ and $s$ with $D^2$, such that $w^{0*}(\mathbf{x}, \mathbf{D}) = [\frac{1}{2}\epsilon^0(\mathbf{x})]D^2$ is the complementary-energy function of the linear, heterogeneous comparison material with arbitrary non-negative dielectric coefficient $\epsilon^0(\mathbf{x})$, and $v(\mathbf{x}, \epsilon^0) = g^*(\mathbf{x}, \frac{1}{2}\epsilon^0)$. Given these, one can state the following theorem (Ponte Castañeda, 1992b)

**Theorem 2:** *Suppose that the (convex) local complementary-energy function $w^*$ of a given nonlinear heterogeneous material with isotropic phases satisfies condition (23) with $g$ being a non-negative, continuous, coercive and concave function of $s = D^2$, and $g(\mathbf{x}, 0) = 0$. Then, the effective complementary-energy function $\mathcal{H}_e^c$ of the nonlinear heterogeneous material is given by*

$$\mathcal{H}_e^c(\langle \mathbf{D} \rangle) = \inf_{\epsilon^0(\mathbf{x}) \geq 0} \{\mathcal{H}_e^{0c}(\langle \mathbf{D} \rangle) + V(\epsilon^0)\}, \tag{27}$$

*where*

$$\mathcal{H}_e^{0c}(\langle \mathbf{D} \rangle) = \min_{\mathbf{D} \in S_2} \int_\Omega w^{0*}(\mathbf{x}, \mathbf{D}) \, d\mathbf{x} \tag{28}$$

*is the effective complementary-energy function of the linear comparison material.*

Note that without the hypotheses of convexity of $f$ and concavity of $g$ the equivalence between the classical minimum energy and the new variational principles would not hold. It can be shown that concavity of $g$ implies convexity of $f$. Moreover, recall that, so far, it has only been assumed explicitly that $w$ is convex and coercive. Since concavity of $g$ implies convexity of $f$, it implies in turn that $w \geq \alpha E^2$ $(\alpha > 0)$ as $E \to \infty$. Thus, a sensible condition may be that, $w(\mathbf{x}, E) \sim E^{1+n} (n \geq 1)$ as $E \to \infty$. Then, $f$ is stronger than, or at least as strong as, affine at infinity, consistent with its convexity. On the other hand, the above assumption for $w$ implies that $w^*(\mathbf{x}, \mathbf{D}) \sim D^{1+1/n}$ as $D \to \infty$, and therefore $g$ is weaker than, or at least as weak as, affine at infinity, consistent with its concavity. Other conditions are possible, but the bounds and estimates that are derived below may require reinterpretation, if the conditions are different. For example, if one lets $n$ in the above conditions be such that $0 < n \leq 1$, then, the suprema and infima in the above relations would have to be replaced by infima and suprema, respectively.

## 2.2   Bounds on the Effective Energy Function

One can now determine bounds and estimates for the effective energy functions of nonlinear heterogeneous materials that are characterized by some appropriate statistical data on their morphology. The main idea of Ponte Castañeda (1992b) is to make use of corresponding bounds and estimates for linear heterogeneous comparison materials, which were described in detail in Chapters 4 and 7 of Volume I, such as the Wiener one-point bounds, the Hashin–Shtrikman two-point bounds, and the Beran three-point bounds, in order to derive the corresponding results for the nonlinear materials. The linear comparison material has the same morphology as the nonlinear composite. In particular, consider heterogeneous materials with $N$ homogeneous phases, characterized by the isotropic energy functions $e_i$ ($i = 1, \cdots, N$), such that the local energy function $w$ of the heterogeneous material is given by

$$w(\mathbf{x}, \mathbf{E}) = \sum_{i=1}^{N} m_i(\mathbf{x}) e_i(E), \tag{29}$$

where $m_i(\mathbf{x})$ is the exclusion indicator function of phase $i$ defined by, $m_i(\mathbf{x}) = 1$ if $\mathbf{x}$ is in phase $i$, and $m_i(\mathbf{x}) = 0$ otherwise. The volume fractions $\phi_i$ of the constituent phases are assumed fixed and given by

$$\phi_i = \int_{\Omega} m_i(\mathbf{x}) \, d\mathbf{x}. \tag{30}$$

Before proceeding with the determination of the bounds and estimates, the following useful corollaries to Theorems 1 and 2 must be stated. Their proofs (which are simple) are given by Ponte Castañeda (1992b).

**Corollary 1:**   *Suppose that Eq. (29) characterizes the local energy-density function of a $N$-phase nonlinear composite, satisfying the hypotheses of Theorem 1. Then, the effective energy function $\mathcal{H}_e$ of the composite satisfies the inequality*

$$\mathcal{H}_e(\langle \mathbf{E} \rangle) \geq \sup_{\epsilon_i^0 > 0} \left\{ \mathcal{H}_e^0(\langle \mathbf{E} \rangle) - \sum_{i=1}^{N} \phi_i v_i(\epsilon_i^0) \right\}, \tag{31}$$

*where $\mathcal{H}_e^0$ is the effective energy function of a linear comparison material with $N$ phases of dielectric constants $\epsilon_i^0$ with volume fractions $\phi_i$, such that the effective dielectric constant $\epsilon_e^0$ of the comparison composite is given by*

$$\epsilon_e^0(\mathbf{x}) = \sum_{i=1}^{N} m_i(\mathbf{x}) \epsilon_i^0. \tag{32}$$

The function $v_i$ is given by Eq. (19), written for the $i$th phase, and the supremum in (31) is evaluated over the set of constants $\epsilon_i^0$ ($i = 1, \cdots, N$).

**Corollary 2:**   *Suppose that the appropriate complementary version of (29) characterizes the local complementary energy function $w^*$ of a $N$-phase non-*

*linear composite, satisfying the hypotheses of Theorem 2. Then, the effective complementary-energy function $\mathcal{H}_e^c$ satisfies*

$$\mathcal{H}_e^c(\langle \mathbf{D} \rangle) \leq \inf_{\epsilon_i^0} \left\{ \mathcal{H}_e^{0c}(\langle \mathbf{D} \rangle) + \sum_{i=1}^N \phi_i v_i(\epsilon_i^0) \right\}, \tag{33}$$

*where $\mathcal{H}_e^{0c}$ is the effective complementary-energy function of a linear comparison composite with N phases of dielectric constants $\epsilon_i^0$ and volume fractions $\phi_i$, such that the effective dielectric constant of the comparison composite is given by*

$$\epsilon_e^0(\mathbf{x}) = \sum_{i=1}^N m_i(\mathbf{x})\epsilon_i^0. \tag{34}$$

## 2.2.1  Lower Bounds

Similar to the effective linear conductivity and dielectric constant of disordered materials described in Chapter 4 of Volume I, we can now derive one-, two- and three-point bounds for the effective nonlinear conductivity and dielectric constant of materials. What follows is a description of derivation of such bounds.

### 2.2.1.1  One-Point Bounds

Consider utilizing the one-point lower bound of Wiener (1912) for linear, anisotropic materials described in Chapter 4 of Volume I for generating a corresponding bound for nonlinear, anisotropic composites. Recall that the bounds are given by, $\langle g(\mathbf{r})^{-1} \rangle \leq g_e \leq \langle g(\mathbf{r}) \rangle$. Although these bounds are not very sharp, their derivation for nonlinear materials provides a useful demonstration of utility of the variational principles of Ponte Castañeda (1992b), described above, for deriving rigorous bounds which will then be used in order to derive the Hashin–Shtrikman and Beran bounds. The Wiener lower bound may be specified as a bound on the effective energy functions of linear composites with dielectric constants $\epsilon_i^0$ and volume fractions $\phi_i$ (with $i = 1, \cdots, N$) via the relation

$$\mathcal{H}_e^0 \geq \frac{1}{2} \left( \sum_{i=1}^N \frac{\phi_i}{\epsilon_i^0} \right)^{-1} \langle \mathbf{E} \rangle^2, \tag{35}$$

where $\mathcal{H}_e^0 = \frac{1}{2}(\epsilon_e^0 \langle \mathbf{E} \rangle) \cdot \langle \mathbf{E} \rangle$ is the effective energy function of the linear material with effective dielectric tensor $\epsilon_e^0$. The nonlinear Wiener lower bound for the effective energy functions $\mathcal{H}_e$ of the nonlinear materials is obtained by applying Eq. (31) to the set of nonlinear composites with given phase volume fractions, and combining the result with the lower bound (35) for the corresponding linear comparison materials. The result is

$$\mathcal{H}_e \geq \sup_{\epsilon_i^0 > 0} \left\{ \frac{1}{2} \left( \sum_{i=1}^N \frac{\phi_i}{\epsilon_i^0} \right)^{-1} \langle E \rangle^2 - \sum_{i=1}^N \phi_i v_i(\epsilon_i^0) \right\}, \tag{36}$$

with

$$v_i(\epsilon_i^0) = \sup_{s>0} \left\{ \frac{1}{2}\epsilon_i^0 s^2 - e_i(s) \right\}. \tag{37}$$

Clearly, the number of optimizations implicit in (36) and (37) is $2N$, but this number may be significantly reduced by using the identity,

$$\left( \sum_{i=1}^{N} \frac{\phi_i}{\epsilon_i^0} \right)^{-1} = \inf_{\omega_i} \left\{ \sum_{i=1}^{N} \phi_i \epsilon_i (1 - \omega_i)^2 \right\}, \tag{38}$$

where the infimum is over the set of variables $\omega_i$ ($i = 1, \cdots, N$) which are subject to a zero-average constraint, i.e., $\langle \omega \rangle = \sum_{i=1}^{N} \phi_i \omega_i = 0$. This identity, when applied to the nonlinear lower bound for $\mathcal{H}_e$ in (36), yields

$$\mathcal{H}_e \geq \sup_{\epsilon_i^0 > 0} \left\{ \inf_{\omega_i} \left\{ \sum_{i=1}^{N} \phi_i \left[ \frac{1}{2}\epsilon_i (1 - \omega_i)^2 \langle E \rangle^2 - v_i(\epsilon_i^0) \right] \right\} \right\}, \tag{39}$$

which in turn leads to

$$\mathcal{H}_e \geq \inf_{\omega_i} \left\{ \sum_{i=1}^{N} \phi_i \sup_{\epsilon_i^0 > 0} \left\{ \frac{1}{2}\epsilon_i (1 - \omega_i)^2 \langle E \rangle^2 - v_i(\epsilon_i^0) \right\} \right\}. \tag{40}$$

In (40), the saddle point theorem and the fact that the argument of the nested supremum and infimum is concave in $\epsilon_i^0$ (since the functions $v_i$ are convex in $\epsilon_i^0$) and convex in $\omega_i$ have been used in order to justify the interchange of the supremum and infimum operations. Finally, application of Eq. (18), specialized to each phase in the form

$$e_i(s) = \sup_{\epsilon_i^0 > 0} \left\{ \frac{1}{2}\epsilon_i^0 s^2 - v_i(\epsilon_i^0) \right\}, \tag{41}$$

leads to

$$\mathcal{H}_e \geq \inf_{\omega_i} \left\{ \phi_i e_i(|1 - \omega_i| \langle E \rangle) \right\}, \tag{42}$$

which is much simpler than the bounds (36) and (37), as it involves only a $N$-dimensional optimization, with one linear constraint, which can easily be embedded in the optimization operation by suitable relabelling of the optimization variables. For example, for a two-phase material, bound (42) becomes

$$\mathcal{H}_e \geq \inf_{\omega} \left\{ \phi_1 e_1(|1 - \phi_2 \omega| \langle E \rangle) + \phi_2 e_2(|1 + \phi_1 \omega| \langle E \rangle) \right\}, \tag{43}$$

where the optimization variable $\omega$ is now unconstrained.

## 2.2.1.2  Two-Point Bounds

One can now use the same technique to derive the Hashin–Shtrikman-type bound for nonlinear isotropic materials. To do this, one should first note that the effective

dielectric tensor of a linear isotropic heterogeneous material is isotropic (i.e., $\epsilon_e^0 = \epsilon_e^0 \mathbf{U}$, where $\mathbf{U}$ is the identity tensor). Then, the Hashin–Shtrikman lower bound $\epsilon_e^{(l)}$ for the effective dielectric constant $\epsilon_e^0$, satisfying $\epsilon_e^0 \geq \epsilon_e^{(l)}$ is given by the expression (see also Chapters 4 and 7 of Volume I)

$$\epsilon_e^{(l)} = \inf_{\omega_i} \left[ \sum_{i=1}^{N} \frac{\phi_i}{\epsilon_i^0 + (d-1)\epsilon^{(l)}} \right]^{-1} - (d-1)\epsilon^{(l)}, \tag{44}$$

where $\epsilon^{(l)} = \inf_s \{\epsilon_s^0\}$. Equation (44), which is subject to the constraint that, $\langle \omega \rangle = 0$, may be rewritten as

$$\epsilon_e^{(l)} = \inf_{\omega_i} \left\{ \sum_{i=1}^{N} \phi_i \left[ \epsilon_i (1-\omega_i)^2 + (d-1)\epsilon^{(l)}\omega_i^2 \right] \right\}. \tag{45}$$

Observe that the effective energy functions $\mathcal{H}_e$ of the macroscopically-isotropic, nonlinear materials can be estimated from relation (31), where $\mathcal{H}_e^0$ now represents the effective energy function of the linear comparison materials with phases of dielectric constants $\epsilon_i^0$ and volume fractions $\phi_i$. Note that, while not all microstructures that are isotropic for linear materials are also isotropic in the nonlinear context, nonlinear isotropic microstructures must also be isotropic in the linear context. Therefore, a lower bound for the effective energy function of linear, isotropic comparison materials is also a lower bound for the subclass of linear comparison composites with "nonlinearly isotropic" microstructure. Hence, replacing $\mathcal{H}_e^0$ in (31) by the lower bound given by (45) generates a lower bound for the nonlinear isotropic composites, with the result being

$$\mathcal{H}_e \geq \sup_{\epsilon_i^0 > 0} \left\{ \inf_{\omega_i} \left\{ \sum_{i=1}^{N} \phi_i \left[ \frac{1}{2}(\epsilon_i(1-\omega_i)^2 + (d-1)\epsilon^{(l)}\omega_i^2)\langle E \rangle^2 - v_i(\epsilon_i^0) \right] \right\} \right\}$$

$$= \inf_{\omega_i} \left\{ \sup_{\epsilon_i^0 > 0} \left\{ \sum_{i=1}^{N} \phi_i \left[ \frac{1}{2}\left(\epsilon_i(1-\omega_i)^2 + (d-1)\epsilon^{(l)}\omega_i^2\right)\langle E \rangle^2 - v_i(\epsilon_i^0) \right] \right\} \right\}, \tag{46}$$

where the saddle point theorem has been used to justify interchanging the supremum and infimum operations. Then, using (41), one obtains

$$\mathcal{H}_e \geq \min_s \left\{ \inf_{\omega_i} \left\{ \sum_{i=1, i \neq s}^{N} \phi_i e_i (|1-\omega_i|\langle E \rangle) \right. \right.$$

$$\left. \left. + \phi_s e_s \sqrt{(1-\omega_s)^2 + (d-1)\frac{1}{\phi_s}\sum_{j=1}^{N}\phi_j \omega_j^2} \langle E \rangle \right\} \right\}, \tag{47}$$

which represents the Hashin–Shtrikman lower bound for nonlinear isotropic materials with isotropic phases of given volume fractions, and is denoted by $\mathcal{H}_{\mathrm{HS}}^{(l)}$.

For a two-phase material, the nonlinear lower bound reduces to

$$\mathcal{H}_{\text{HS}}^{(l)}(\langle E \rangle) = \min \begin{Bmatrix} \inf_{\omega}\{\phi_1 e_1(|1 - \phi_2\omega|\langle E \rangle) + \phi_2 e_2\sqrt{(1 + \phi_1\omega)^2 + (d - 1)\phi_1\omega^2} \langle E \rangle \} \\ \inf_{\omega}\{\phi_1 e_1\sqrt{(1 - \phi_2\omega)^2 + (d - 1)\phi_2\omega^2} \langle E \rangle + \phi_2 e_2(|1 + \phi_1\omega|\langle E \rangle) \} \end{Bmatrix}.$$

$$(48)$$

Note that for a two-phase nonlinear material, the bounds given above involve only one optimization. Moreover, the method described here has a distinct advantage in that, it utilizes the linear heterogeneous comparison material in conjunction with its linear bounds and estimates (other than, for example, the Hashin–Shtrikman bounds) to yield the corresponding nonlinear bounds and estimates.

### 2.2.1.3   Three-Point Bounds

As another illustration of this feature of the method, the lower bounds of the Beran-type for two-phase, nonlinear isotropic materials are derived which, as discussed in Chapters 4 and 7 of Volume I, are generally tighter than the Hashin–Shtrikman bounds except, of course, for those microstructures for which the Hashin–Shtrikman bounds become exact estimates, such as the coated-spheres model (see Sections 4.4 and 7.2.3 of Volume I). As discussed in Chapters 4 and 7 of Volume I, the Beran bound (Beran, 1965), simplified by Milton (1981a,b), depends on the volume fraction of the phases and on one additional microstructural parameter $\zeta_i$, and is given by

$$\epsilon_e^{(l)} = \left[ \sum_{i=1}^{2} \frac{\phi_i}{\epsilon_i^0 + (d - 1)\epsilon^{(l)}} \right]^{-1} - (d - 1)\epsilon^{(l)}, \qquad (49)$$

which is identical in form to (44), except that $\epsilon^{(l)}$ is now given by

$$\epsilon^{(l)} = \left( \sum_{i=1}^{2} \frac{\zeta_i}{\epsilon_i^0} \right)^{-1}, \qquad (50)$$

where the third-order microstructural parameters $\zeta_1$ and $\zeta_2 = 1 - \zeta_1$ are both in the interval [0,1], and were described in detail in Chapters 4 (see Sections 4.5.2 and 4.5.3) and 7 of Volume I (see section 7.4.3). Substituting (49) into the lower-bound approximation (31) and following a procedure very similar to that used for the Hashin–Shtrikman bound, one arrives at the following lower bound for the nonlinear energy function,

$$\mathcal{H}_{\text{B}}^{(l)}(\langle E \rangle) = \inf_{\omega, \gamma} \Bigg\{ \phi_1 e_1 \sqrt{(1 - \phi_2\omega)^2 + (d - 1)\phi_2\zeta_1\omega^2(1 - \zeta_2\gamma)^2} \langle E \rangle$$

$$+ \phi_2 e_2 \sqrt{(1 + \phi_1\omega)^2 + (d - 1)\phi_1\zeta_2\omega^2(1 + \zeta_1\gamma)^2} \langle E \rangle \Bigg\}. \qquad (51)$$

Note that the corresponding nonlinear Hashin–Shtrikman lower bound follows immediately from (51) by choosing either $\zeta_1 = 0$ or $\zeta_2 = 1$, whichever yields the lowest value (note also that the infimum problem over $\gamma$ becomes trivial in

either case), which is completely analogous to the corresponding result for linear two-phase materials.

## 2.2.2 Approximate Estimates of the Effective Energy

Although the above developments were for the effective dielectric constant of nonlinear materials, they are equally applicable to the problem of estimating their nonlinear conductivity. We will discuss this problem in detail later in this chapter, but it is useful to note here the work of Gibiansky and Torquato (1998a). They wrote Eq. (51) in a more general form

$$\mathcal{H}_e(\langle E \rangle) = \inf_{\omega,\gamma} \left\{ \phi_1 e_1 \sqrt{(1 - \phi_2\omega)^2 + (d - 1)\phi_2\zeta_1\omega^2(1 - \zeta_2\gamma)^2 + B\phi_2\zeta_1\zeta_2\omega^2\gamma^2} \; \langle E \rangle \right.$$

$$\left. + \phi_2 e_2 \sqrt{(1 + \phi_1\omega)^2 + (d - 1)\phi_1\zeta_2\omega^2(1 + \zeta_1\gamma)^2} \; \langle E \rangle \right\} \tag{52}$$

which must be optimized over the two scalar variables $\omega \in (-\infty, \infty)$ and $\gamma \in (-\infty, \infty)$. The optimization can be carried out either analytically, if the energy functions of the nonlinear phases are sufficiently simple, or numerically. Here, $B$ is a parameter which is given by (Torquato, 1985a,b)

$$B = (d - 1)\frac{(d - 1) - \zeta_2}{1 - (d - 1)\zeta_2}. \tag{53}$$

We can now consider two important limiting cases.

### 2.2.2.1 Conductor–Superconductor Composites

If we assume that the inclusion phase 2 is a superconducting material, i.e., if

$$e_2(E) = \begin{cases} 0, & \text{if } E = 0, \\ \infty, & \text{if } E \neq 0, \end{cases} \tag{54}$$

then, for such a composite, the right-hand side of Eq. (52) will be divergent unless the argument of the function $e_2$ is equal to zero, i.e., unless

$$\sqrt{(1 + \phi_1\omega)^2 + (d - 1)\phi_1\zeta_2\omega^2(1 + \zeta_1\gamma)^2} \; \langle E \rangle = 0,$$

which is possible (for $d \neq 1$) only if,

$$\omega = -(\phi_1)^{-1}, \quad \gamma = -(\zeta_2)^{-1}, \tag{55}$$

which represent the optimal values of these parameters. An approximate expression for the effective energy of the nonlinear material is then obtained:

$$\mathcal{H}_e(\langle E \rangle) = \phi_1 e_1 \left[ \sqrt{\frac{\zeta_1 + (d - 1)\phi_2 + B\phi_2\zeta_2}{\zeta_1\phi_1^2}} \; \langle E \rangle \right]. \tag{56}$$

Therefore, if, for example, the matrix is a strongly nonlinear material with the energy function, $e_1 = g_1^{(n)} E^n / n$, and if the effective nonlinear conductivity $g_e^{(n)}$

is defined by, $\mathcal{H}_e(\langle E \rangle) = g_e^{(n)} \langle E \rangle^n / n$, one obtains

$$\frac{g_e^{(n)}}{g_1^{(1)}} = \phi_1 \left[ \frac{\zeta_1 + (d-1)\phi_2 + B\phi_2\zeta_2}{\zeta_1\phi_1^2} \right]^{n/2}. \tag{57}$$

Equations (56) and (57) can now be utilized for estimating the effective energy of nonlinear composites with superconducting inclusions, provided that the appropriate expressions for the microstructural parameters $\zeta_1$ and $\zeta_2 = 1 - \zeta_1$ are available, a matter that was discussed in detail in Section 4.5.3 of Volume I. In particular, it can be shown that Eq. (56) provides an estimate of the effective energy which is always larger than the estimates provided by Eqs. (48) and (51), hence satisfying these rigorous bounds.

### 2.2.2.2   Conductor–Insulator Composites

Consider now the opposite limit in which the inclusion phase is insulating, so that $e_2(E) = 0$ for all $E$. Then, the optimal values of $\omega$ and $\gamma$ are obtained by minimizing

$$\sqrt{(1 - \phi_2\omega)^2 + (d-1)\phi_2\zeta_1\omega^2(1 - \zeta_2\gamma)^2 + B\phi_2\zeta_1\zeta_2\omega^2\gamma^2} \, \langle E \rangle,$$

with respect to these parameters. It is straightforward to show that the optimal values are given by

$$\omega = \frac{B + (d-1)\zeta_2}{B\phi_2 + (d-1)(B\zeta_1 + \phi_2\zeta_2)}, \quad \gamma = \frac{d-1}{B + (d-1)\zeta_2}, \tag{58}$$

which then lead to

$$\mathcal{H}_e(\langle E \rangle) = \phi_1 e_1 \left[ \sqrt{\frac{(d-1)B\zeta_1}{B\phi_2 + (d-1)(B\zeta_1 + \zeta_2\phi_2)}} \, \langle E \rangle \right]. \tag{59}$$

For a strongly-nonlinear (power-law) matrix, the effective conductivity of the composite is then given by

$$\frac{g_e^{(n)}}{g_1^{(n)}} = \phi_1 \left[ \frac{(d-1)B\zeta_1}{B\phi_2 + (d-1)(B\zeta_1 + \zeta_2\phi_2)} \right]^{n/2}. \tag{60}$$

Let us mention that Eqs. (56), (57), (59) and (60) are accurate only if the inclusion phase does not form large clusters.

## 2.2.3   Upper Bounds and Estimates

The derivation of upper bounds for the effective energy functions of nonlinear materials is intrinsically more difficult than the corresponding lower bounds. This is because approximations such as (31) do not work in this case. While it is possible to derive the Wiener upper bound, derivation of upper bounds of the Hashin–Shtrikman- and Beran-type bounds has proven to be very difficult. Instead, one may obtain *upper estimates* or, more precisely, *lower estimates for the upper bound*, of the Hashin–Shtrikman- and Beran-types.

The derivation of the Wiener upper bound is accomplished by the corresponding upper bound for linear materials with an arbitrary dielectric constant $\epsilon^0(\mathbf{x})$, and is given by

$$\mathcal{H}_e^0(\langle \mathbf{E} \rangle) \leq \frac{1}{2} \left[ \int_\Omega \epsilon^0(\mathbf{x}) d\mathbf{x} \right] \langle E \rangle^2. \tag{61}$$

Then, application of (61) to (20) leads to

$$\mathcal{H}_e(\langle \mathbf{E} \rangle) \leq \sup_{\epsilon^0(\mathbf{x}) \geq 0} \left\{ \frac{1}{2} \left[ \int_\Omega \epsilon^0(\mathbf{x}) d\mathbf{x} \right] \langle E \rangle^2 - \int_\Omega v[\mathbf{x}, \epsilon^0(\mathbf{x})] d\mathbf{x} \right\}$$

$$= \int_\Omega \sup_{\epsilon^0 \geq 0} \left\{ \frac{1}{2} \epsilon^0 \langle E \rangle^2 - v(\mathbf{x}, \epsilon^0) \right\} d\mathbf{x} = \int_\Omega e(\mathbf{x}, E) d\mathbf{x}, \tag{62}$$

which, via (29), leads to the nonlinear Wiener upper bound

$$\mathcal{H}_W^{(u)} \leq \sum_{i=1}^N \phi_i e_i(\langle E \rangle). \tag{63}$$

The determination of an estimate for the Hashin–Shtrikman upper bound, or the upper estimate, is accomplished by application of approximation (31) to the Hashin–Shtrikman upper bounds for the linear comparison material. The upper bound for the effective energy function of the linear comparison material may be given in terms of the upper bound for its effective dielectric constant:

$$\epsilon_e^+ = \left[ \sum_{i=1}^N \frac{\phi_i}{\epsilon_i^0 + (d-1)\epsilon^+} \right]^{-1} - (d-1)\epsilon^+, \tag{64}$$

where $\epsilon^+ = \sup_i\{\epsilon_i^0\}$. The procedure that utilizes the lower bound (44) for the linear comparison material to obtain a lower bound for the nonlinear material may now be repeated. To derive the upper estimates, one utilizes (64) instead of (44), in which case the result would be the same as (47) and (48) for the $N$-phase and two-phase nonlinear materials, respectively, with the difference that the outermost minimum operations must now be replaced by maximum operations. The result, denoted by $\mathcal{H}_{HS}^+$, is referred to as the Hashin–Shtrikman upper estimate. However, as shown below, $\mathcal{H}_{HS}^+$ is not, in general, an upper bound for $\mathcal{H}_e$.

The same arguments and analyses also apply to the Beran upper bounds. Hence, one can obtain upper estimate for the Beran-type bounds. If

$$\epsilon^+ = \sum_{i=1}^2 \zeta_i \epsilon_i^0, \tag{65}$$

then, the corresponding result for the upper estimate (which, in general, *is not* an upper bound) for nonlinear isotropic materials is given by

$$\mathcal{H}_B^+(\langle \mathbf{E} \rangle) = \inf_\omega \left\{ \phi_1 e_1 \sqrt{(1 - \phi_2 \omega)^2 + (d-1)\phi_2 \zeta_1 \omega^2} \, \langle E \rangle \right.$$

$$\left. + \phi_2 e_2 \sqrt{(1 + \phi_1 \omega)^2 + (d-1)\phi_1 \zeta_2 \omega^2} \, \langle E \rangle \right\}. \tag{66}$$

## 2.3    Exact Results for Laminates

Having derived the rigorous lower bounds and also the lower estimates for the upper bounds, two important issues must be now addressed.

(1) How accurate are the lower bounds given above for any type of materials' morphology?
(2) Do the upper estimates represent rigorous bounds?

To address these issues, one can, for example, analyze the effective properties of sequentially-laminated materials (Ponte Castañeda, 1992b) which have provided useful insights into the properties of linear materials, even though they represent highly ideal models. A sequentially-laminated material (or laminate, for short) is an iterative construction obtained by layering one type of laminated material with other types of laminated materials, or directly with the homogeneous phases that make up the composite, in such a way as to produce hierarchical microstructures of increasing complexity. The *rank* of the laminate is the number of layering operations required to reach the final iterated morphology. Figure 2.1 presents a first-rank laminate, constructed by mixing layers of two homogeneous phases to obtain a simple laminate with layering direction $\mathbf{n}_1$. A second-rank laminate, also shown in Figure 2.1, is obtained by layering the first-rank laminate with a third phase or, alternatively, with one of the original phases (say 2), in a different layering direction $\mathbf{n}_2$. In general, $\mathbf{n}_1$ and $\mathbf{n}_2$ can take on any orientation. It is assumed that the length scale of the embedded laminates is small compared with the length scale of the embedding laminates. Under this assumption, the fields will be essentially constant within each elemental layer, provided that the boundary conditions applied to the laminate are uniform. This feature greatly simplifies the computation of effective properties, thereby making sequentially-laminated materials very useful constructions. With such a microstructure, the effective energy function of a simple



FIGURE 2.1. Examples of first-rank (left) and second-rank laminates (right).

linear laminate lies within and attains (for specific orientations of the applied fields) the Wiener bounds. Thus, at least in this case, the Wiener bounds on the effective energy function of arbitrarily anisotropic-linear materials are sharp.

In the context of two-phase linear materials, it is known that only iterated laminates of rank greater than or equal to the dimension of the underlying physical space ($d = 2$ or 3) can have isotropic properties. The isotropy is obtained by choosing the relative volume fractions and the layering directions of each of the embedded laminates in such a way that the tensor representing the effective property of interest is isotropic, while the absolute volume fractions of the constituent phases remain fixed. One might criticize sequentially-laminated materials by noting that the inclusions are flat, whereas in practice the inclusions are often equi-axed. However, one must note that iterated laminates can be used to model arbitrarily close the properties of *any* two-phase microstructure (Milton, 1986). For example, the coated-spheres model of Hashin and Shtrikman (1962a,b, 1963) possesses exactly the same effective properties as an isotropic iterated laminate with the same volume fractions. In the coated-spheres model (see also Chapters 3, 4 and 7 of Volume I) the material consists of composite spheres that are composed of a spherical core of conductivity $g_2$ and radius $a$, surrounded by a concentric shell of conductivity $g_1$ with an outer radius $b > a$. The ratio $a/b$ is fixed, and the volume fraction $\phi_2$ of inclusions in $d$ dimensions is given by $\phi_2 = (a/b)^d$. The composite spheres fill the space, implying that there is a sphere size distribution that extends to infinitesimally-small spheres. In Chapters 4 and 7 of Volume I we derived exact expression for the effective conductivity and elastic moduli of the coated-spheres model and low-rank laminates.

As shown in Chapters 4 and 7 of Volume I, the Hashin–Shtrikman bounds for the coated-spheres model, which represent isotropic microstructures, are exact estimates. Thus, it may seem that the coated-spheres model may be more realistic than the iterated laminates. However, the laminates have a distinct advantage over the coated-spheres model in that, they contain a finite number of length scales, in contrast with the coated-spheres microstructure which involves an infinite number of length scales because, as described above, the composite spheres must cover all sizes to fill the space. Another advantage of sequentially-laminated materials is that, when subjected to uniform boundary conditions, the fields are piecewise constant within the material (regardless of whether the composite's phases are linear or nonlinear), except in small boundary layer regions at the interfaces separating laminates of different ranks, the effect of which is made negligible by the hypothesis of separated length scales. This fact was used for deriving the exact results for the effective linear properties of the laminates presented in Chapters 4 and 7 of Volume I.

To compute the effective energy function of nonlinear rank$-d$ laminates ($d = 2$ and 3) with layering directions $\mathbf{n}_1, \cdots, \mathbf{n}_d$, we denote by $\phi_I$ the volume fraction of phase 1 with energy function $e_1$ in the first-rank laminate, and note that $1 - \phi_I$ is the corresponding volume fraction of phase 2 with energy $e_2$. Ponte Castañeda (1992b) showed that the effective energy function of the nonlinear, first-rank laminate is

given by

$$\mathcal{H}_e^I(\langle \mathbf{E} \rangle) = \inf_{\omega_I^{(1)}, \omega_I^{(2)}} \{\phi_I e_1(s_1) + (1 - \phi_I)e_2(s_2)\}, \tag{67}$$

subject to the constraints that $\langle \omega_I \rangle = \phi_I \omega_I^{(1)} + (1 - \phi_I)\omega_I^{(2)} = 0$, and where

$$s_1 = \sqrt{\langle E \rangle^2 - E_1^2 + [1 - \omega_I^{(1)}]^2 E_1^2},$$

$$s_2 = \sqrt{\langle E \rangle^2 - E_1^2 + [1 - \omega_I^{(2)}]^2 E_1^2}, \tag{68}$$

where $E_1 = \langle \mathbf{E} \rangle \cdot \mathbf{n}_1$.

Consider now the second-rank laminate obtained by mixing layers of the first-rank laminate with layers of a third phase characterized by an energy function $e_3$ and relative (to the second-rank laminate) volume fractions $\phi_{II}$ and $1 - \phi_{II}$, respectively. The new lamination direction $\mathbf{n}_2$ is orthogonal to $\mathbf{n}_1$. Then, the following energy function for the nonlinear second-rank laminate in dimension $d \geq 2$ is obtained (Ponte Castañeda, 1992b):

$$\mathcal{H}_e^{II}(\langle \mathbf{E} \rangle) = \inf_{\omega_I^{(1)}, \omega_I^{(2)}, \omega_{II}^{(1)}, \omega_{II}^{(2)}} \{\phi_{II}\phi_I e_1(s_1) + \phi_{II}(1 - \phi_I)e_2(s_2) + (1 - \phi_{II})e_3(s_3)\}, \tag{69}$$

subject to the constraints that $\langle \omega_I \rangle = \langle \omega_{II} \rangle = 0$ (where $\langle \omega_{II} \rangle$ is defined in a manner analogous to $\langle \omega_I \rangle$), and where

$$s_1 = \sqrt{\langle E \rangle^2 - E_1^2 - E_2^2 + [1 - \omega_I^{(1)}]^2 E_1^2 + [1 - \omega_{II}^{(1)}]^2 E_2^2},$$

$$s_2 = \sqrt{\langle E \rangle^2 - E_1^2 - E_2^2 + [1 - \omega_I^{(2)}]^2 E_1^2 + [1 - \omega_{II}^{(1)}]^2 E_2^2}, \tag{70}$$

$$s_3 = \sqrt{\langle E \rangle^2 - E_2^2 + [1 - \omega_{II}^{(2)}]^2 E_2^2},$$

where $E_i = \langle \mathbf{E} \rangle \cdot \mathbf{n}_i$.

A similar result can be obtained for a two-phase, nonlinear second-rank laminate. In this case, the result for $\mathcal{H}_e^{II}$ is generally anisotropic and direction-dependent, but it may be used in two dimensions (2D) for deriving an isotropic result, for each value of $\langle E \rangle$, by an appropriate choice of $\phi_{II}$ (but not the choice that makes the corresponding linear second-rank laminate isotropic), which is obtained by requiring that $\phi_{II}$ $(0 \leq \phi_{II} \leq 1)$ and $E_1$ satisfy the following relations

$$\frac{\partial \mathcal{H}_e^{II}}{\partial E_1} = 0 \quad \text{and} \quad \frac{\partial \mathcal{H}_e^{II}}{\partial \phi_{II}} = 0, \tag{71}$$

where the first relation is subject to the constraint that, $E_1^2 + E_2^2 = \langle E \rangle^2$, while in the second relation one assumes that $\langle E \rangle$ is fixed. These conditions follow by performing a Taylor series expansion of (69) in $\phi_{II}$ and $E_1$ and requiring that the expansion yield the same result for any choice of $\phi_{II}$ and $E_1$. Physically, this corresponds to selecting a microstructure (by choosing $\phi_{II}$)—with fixed overall volume fractions of the phases—for each value of $\langle E \rangle$, ensuring that $\mathcal{H}_e^{II}$ is

independent of the direction of $\langle \mathbf{E} \rangle$, thus guaranteeing that the resulting energy function is isotropic. However, the resulting energy function does not correspond to a fixed microstructure, rather to a family of (anisotropic) microstructures, each one of which is obtained from one value of the applied electric field.

The effective energy function of a nonlinear third-rank laminate is obtained by analyzing the effective behavior of a simple laminate made up of layers of the second-rank laminate and of layers of a fourth phase with energy function $e_4$ and volume fractions $\phi_{III}$ and $1 - \phi_{III}$, respectively. The new layering direction $\mathbf{n}_3$ is selected to be orthogonal to both $\mathbf{n}_2$ and $\mathbf{n}_1$. Then, the effective energy function of the nonlinear third-rank laminate is given by (Ponte Castañeda, 1992b)

$$\mathcal{H}_e^{III}(\langle \mathbf{E} \rangle) = \inf_{\omega_I^{(1)}, \omega_I^{(2)}, \omega_{II}^{(1)}, \omega_{II}^{(2)}, \omega_{III}^{(1)}, \omega_{III}^{(2)}} \{ \phi_{III} \phi_{II} \phi_I e_1(s_1)$$

$$+ \phi_{III} \phi_{II}(1 - \phi_I) e_2(s_2) + \cdots + \phi_{III}(1 - \phi_{II}) e_3(s_3) + (1 - \phi_{III}) e_4(s_4) \},$$

(72)

subject to the constraints that $\langle \omega_I \rangle = \langle \omega_{II} \rangle = \langle \omega_{III} \rangle = 0$, and where

$$s_1 = \sqrt{\langle E \rangle^2 - E_1^2 - E_2^2 - E_3^2 + [1 - \omega_I^{(1)}]E_1^2 + [1 - \omega_{II}^{(1)}]^2 E_2^2 + [1 - \omega_{III}^{(1)}]^2 E_3^2},$$

$$s_2 = \sqrt{\langle E \rangle^2 - E_1^2 - E_2^2 - E_3^2 + [1 - \omega_I^{(2)}]E_1^2 + [1 - \omega_{II}^{(1)}]^2 E_2^2 + [1 - \omega_{III}^{(1)}]^2 E_3^2},$$

$$s_3 = \sqrt{\langle E \rangle^2 - E_2^2 - E_3^2 + [1 - \omega_{II}^{(2)}]E_2^2 + [1 - \omega_{III}^{(1)}]^2 E_3^2},$$

$$s_4 = \sqrt{\langle E \rangle^2 - E_3^2 + [1 - \omega_{III}^{(2)}]^2 E_3^2},$$

(73)

where, as before, $E_i = \langle \mathbf{E} \rangle \cdot \mathbf{n}_i$.

The effective energy function of a two-phase, nonlinear third-rank laminate may be obtained by letting $e_4 = e_3 = e_2$ in (72). Then, for 3D third-rank laminates Eq. (72) may be used to obtain an isotropic energy by choosing $\phi_{II}$ and $\phi_{III}$, and $E_1$, $E_2$, and $E_3$ with $E_1^2 + E_2^2 + E_3^2 = \langle E \rangle^2$, such that

$$\frac{\partial \mathcal{H}_e^{III}}{\partial E_1} = \frac{\partial \mathcal{H}_e^{III}}{\partial E_2} = 0, \quad \frac{\partial \mathcal{H}_e^{III}}{\partial \phi_{II}} = \frac{\partial \mathcal{H}_e^{III}}{\partial \phi_{III}} = 0.$$

(74)

## 2.4    Effective Dielectric Constant of Strongly Nonlinear Materials

To illustrate the application of the methods described above, we consider two important examples that we have been considering throughout this book, both in Volume I and the present Volume. Both limits involve a nonlinear matrix with isotropic potential $e_2 = e$ (subject to the restrictions of Theorem 1), and an inclusion phase that, similar to the case of nonlinear conductivity discussed above, has either an infinite dielectric constant or, alternatively, a zero dielectric constant. In the first case, $e_1 = 0$ if $E = 0$, or $e_1 = \infty$ otherwise, while in the second, $e_1 = 0$

regardless of the value of $E$. Moreover, we specialize the results to the case in which the nonlinearity of the matrix is of power-law type, which is usually referred to as strong nonlinearity. This type of nonlinearity is characterized by the energy-density function

$$e(E) = (n + 1)^{-1} \epsilon^{(n)} E^{n+1}, \tag{75}$$

where $n \geq 1$, and $\epsilon^{(n)}$ is the nonlinear dielectric constant. Equation (75) has the advantage that it yields the same type of behavior for the isotropic composite materials with perfectly conducting or insulating inclusions. Thus, for both types of isotropic composites, we have

$$\mathcal{H}_e(\langle E \rangle) = (n - 1)^{-1} \epsilon_e^{(n)} \langle E \rangle^{n+1}, \tag{76}$$

where $\epsilon_e^{(n)}$ is the effective nonlinear dielectric constant of the material. For the anisotropic materials, the form of the effective energy will, in general, be different, but the Wiener bounds will be of the same form. We can then characterize the behavior of the Wiener, Hashin–Shtrikman and isotropic (in the sense defined earlier) laminates for this class of materials in terms of the effective nonlinear dielectric constant.

## 2.4.1 Inclusions with Infinite Dielectric Constant

The results for the bounds and estimates of 2D materials are not essentially different from those for 3D composites, and therefore only the results for the 3D materials are presented. Consider first the Wiener and the Hashin–Shtrikman lower bounds and the isotropic laminate estimate for nonlinear materials with perfectly conducting inclusions. These results, expressed in terms of the effective nonlinear dielectric constant, are given by (Ponte Castañeda, 1992b)

$$\frac{\epsilon_{\mathrm{W}}^{(l)}}{\epsilon^{(n)}} = (1 - \phi)^{-n}, \tag{77}$$

$$\frac{\epsilon_{\mathrm{HS}}^{(l)}}{\epsilon^{(n)}} = \frac{(1 + 2\phi)^{(n+1)/2}}{(1 - \phi)^n}, \tag{78}$$

$$\frac{\epsilon_e^{III}}{\epsilon} = \sup_{x,y} \left\{ (1 - y)^{2n/(n-1)} + \left( \frac{2 - y}{y} \right)^{(n+1)/(n-1)} \left[ (xy - \phi)^{2n/(n-1)} \right. \right.$$
$$\left. \left. \times \left( \frac{2 - x}{x} \right)^{(n+1)/(n-1)} + [(1 - x)y]^{2n/(n-1)} \right]^{(1-n)/2} \right\}, \tag{79}$$

where $\phi = \phi_2$ is the volume fraction of the inclusions, and the optimization variables $x$ and $y$ are subject to the constraints, $0 \leq x, y \leq 1$ and $xy \geq \phi$. Note that as $n \to \infty$,

$$\frac{\epsilon_e^{III}}{\epsilon^{(n)}} \to (1 - 8\phi + 12\phi^{4/3} - 6\phi^{5/3} + \phi^2)^{-n/2}, \tag{80}$$

which is different from, but close to, $\epsilon_{\mathrm{HS}}^{(l)}/\epsilon^{1/n}$ in the same limit. In general, the Hashin–Shtrikman bound provides estimates that are very close to those for

the laminates, and both differ strongly from the Wiener bound, with the latter yielding estimates that are larger than the former two.

### 2.4.2 Inclusions with Zero Dielectric Constant

Consider now the corresponding results for 3D nonlinear materials with perfectly insulating inclusions. The results for the Wiener upper bound, the Hashin–Shtrikman upper estimate, and the exact estimate for the isotropic laminate are given by (Ponte Castañeda, 1992b)

$$\frac{\epsilon_{\mathrm{W}}^{(u)}}{\epsilon^{(n)}} = 1 - \phi, \tag{81}$$

$$\frac{\epsilon_{\mathrm{HS}}^{+}}{\epsilon^{(n)}} = \frac{1 - \phi}{1 + \frac{1}{2}\phi^{(n+1)/2}}, \tag{82}$$

$$\frac{\epsilon_{e}^{III}}{\epsilon^{(n)}} = \sup_{x,y} \left\{ \frac{[(1 - x)y + (1 - y)p]^{(n+1)/2}}{[(xy - \phi)q^{(n+1)/(n-1)} + (1 - x)y + (1 - y)p^{(n+1)/(n-1)}]^{(n-1)/2}} \right\}, \tag{83}$$

where $p$ is the root of the quadratic equation,

$$\frac{1}{2}\frac{1 - y}{1 - x}(2 - x)p^2 - \left[\frac{1}{2}xy + 2(1 - y)\right]p + (1 - x)(1 - y) = 0,$$

and

$$q = \frac{xy}{xy - \phi}\frac{1 - x}{2 - x}.$$

In this case, the Hashin–Shtrikman upper estimates for the isotropic composite lie well below the Wiener bounds for arbitrarily anisotropic composites. On the other hand, the exact estimates for the nonlinear isotropic laminates lie above the Hashin–Shtrikman upper estimates, hence verifying that the Hashin–Shtrikman upper estimates are not in general upper bounds. This is due to the fact that the isotropic laminates correspond to specific microstructures within the class of isotropic composite materials, and if the Hashin–Shtrikman upper estimates were rigorous bounds for such materials, they would have to lie above all possible isotropic microstructures, and, in particular, they must lie above the isotropic laminates. Nevertheless, the effective dielectric constants of the isotropic laminates are not far from the Hashin–Shtrikman upper estimates.

## 2.5   Effective Conductivity of Nonlinear Materials

The above methods of deriving bounds and estimates for the effective dielectric constant of heterogeneous nonlinear materials can also be used for estimating their effective conductivity (Ponte Castañeda, 1998). Equation (75) is now written as

$$e_i(E) = (n + 1)^{-1}g_i^{(n)}|E|^{n+1}, \tag{84}$$

where $g_i^{(n)}$ is the generalized nonlinear conductivity of phase $i$. The linear comparison materials are now defined by the quadratic energy-density function,

$$w^0(\mathbf{x}, \mathbf{E}) = \frac{1}{2} g^0(\mathbf{x}) E^2, \tag{85}$$

where $g^0(\mathbf{x})$ is the conductivity of the fictitious linear material. Then, under the hypothesis that the functions $e_i$ of the nonlinear material are convex on $E^2$, the analogues of Eqs. (37) and (41) for the conductivity problem are given by

$$e_i(E) = \max_{g^0 \geq 0} \left\{ \frac{1}{2} g^0(\mathbf{x}) E^2 - v_i(g^0) \right\}, \tag{86}$$

$$v_i(g^0) = \max_{E} \left\{ \frac{1}{2} g^0 E^2 - e_i(E) \right\}. \tag{87}$$

Note that if the functions $e_i$ are smooth, the maxima in Eqs. (86) and (87) are attained at

$$\frac{1}{2} E^2 = \frac{\partial v_i}{\partial g^0}, \quad g^0 = \frac{1}{E} \frac{\partial e_i}{\partial E}, \tag{88}$$

respectively, which are inverse of each other. Then, the analogue of Eq. (31) for the conductivity problem is given by

$$\mathcal{H}_e(\langle \mathbf{E} \rangle) = \max_{g^0(\mathbf{x}) \geq 0} \left\{ \mathcal{H}_e^0(\langle \mathbf{E} \rangle) - \sum_{i=1}^{N} \phi_i \langle v_i [g^0(\mathbf{x})] \rangle_i \right\}, \tag{89}$$

where $\mathcal{H}_e^0$ is the effective energy function of the linear comparison material, with local energy function (85), such that

$$\mathcal{H}_e^0(\langle \mathbf{E} \rangle) = \min_{\mathbf{E} \in S_2} \langle w^0(\mathbf{x}, \mathbf{E}) \rangle, \tag{90}$$

as before, where $S_2$ is the set defined by (11).

   Equation (89), together with Eqs. (87) and (90), provide variational representation of the effective energy function of the nonlinear material in terms of the effective energy function of a fictitious linear composite, the choice of which is determined by Eq. (89). It should be emphasized that the conductivity $g^0(\mathbf{x})$ of the comparison material is an *arbitrary* non-negative function of $\mathbf{x}$, and that the minimum principle (89) is valid only under the hypothesis that the functions $e_i$ are convex on $E^2$. If these functions are concave on $E^2$ (as when, for example, $0 \leq n < 1$), an analogous result would hold, but with the maximum in Eq. (89) replaced by a minimum, and with the function $v_i$ redefined such that the maximum in Eq. (87) is replaced by a minimum.

## 2.5.1  Materials with Nonlinear Isotropic Phases

Even if each of the nonlinear phases is homogeneous, the solutions for the comparison conductivities $g^0(\mathbf{x})$ in the variational principle (89) is not, in general,

constant over the individual phases, unless the actual fields are constant throughout the phases. However, as discussed earlier in this chapter, a lower bound for $\mathcal{H}_e^0$ can be obtained by restricting the class of trial comparison conductivity fields to be constant within each phase such that

$$g^0(\mathbf{x}) = \sum_{i=1}^{N} m_i(\mathbf{x})g_i^0, \tag{91}$$

where $g_i^0$ is constant, and $m_i(\mathbf{x})$ is the indicator function of phase $i$, as before. Equation (91) follows from the fact that the maximum over a set is, in general, larger than the maximum over any subset of the original set. Therefore, from Eqs. (89) and (91), it follows that

$$\mathcal{H}_e(\langle \mathbf{E} \rangle) \geq \max_{g_i^0 > 0} \left\{ \mathcal{H}_e^0(\langle \mathbf{E} \rangle) - \sum_{i=1}^{N} \phi_i v_i(g_i^0) \right\}, \tag{92}$$

where

$$\mathcal{H}_e^0(\langle \mathbf{E} \rangle) = \frac{1}{2} \langle \mathbf{E} \rangle \cdot [\mathbf{g}_e^{(l)} \langle \mathbf{E} \rangle] = \min_{E \in S_2} \left\{ \frac{1}{2} \sum_{i=1}^{N} \phi_i g_i^0 \langle E^2 \rangle_i \right\}. \tag{93}$$

Here, $\mathbf{g}_e^{(l)}$ is the effective conductivity tensor of a linear comparison material with precisely the same morphology as the original nonlinear composite which, in general, is anisotropic.

As discussed earlier in this chapter, the above estimates for $N$-phase nonlinear materials represent lower bounds for $\mathcal{H}_e$. Thus, lower bounds for $\mathbf{g}_e^{(l)}$ may be used to generate lower bounds for $\mathcal{H}_e$, but upper bounds for $\mathbf{g}_e^{(l)}$ cannot be used for deriving upper bounds for $\mathcal{H}_e$. In this case, one can ignore the inequality in (92) and reinterpret it as an *approximate* equality in order to obtain estimates for specific types of materials. Denoting by $\hat{g}_i^0$ the optimal values of $g_i^0$ from Eq. (92), it follows that the average current field $\langle \mathbf{I} \rangle$ is given by

$$\langle \mathbf{I} \rangle = \mathbf{g}_e^{(l)}(\hat{g}_1^0, \cdots, \hat{g}_N^0)\langle \mathbf{E} \rangle + \sum_{i=1}^{N} \left\{ \frac{1}{2}\langle \mathbf{E} \rangle \cdot \left[ \frac{\partial \mathbf{g}_e^{(l)}}{\partial g_i^0}(\hat{g}_1^0, \cdots, \hat{g}_N^0)\langle \mathbf{E} \rangle \right] - \phi_i \frac{\partial v_i}{\partial g_i^0}(\hat{g}_i^0) \right\} \frac{\partial \hat{g}_i^0}{\partial \langle \mathbf{E} \rangle}, \tag{94}$$

so that, the maximum in (92) for the general bound is attained at

$$\frac{1}{2}\langle \mathbf{E} \rangle \cdot \left[ \frac{\partial \mathbf{g}_e^{(l)}}{\partial g_i^0}(\hat{g}_1^0, \cdots, \hat{g}_N^0)\langle \mathbf{E} \rangle \right] = \phi_i \frac{\partial v_i}{\partial g_i^0}(\hat{g}_i^0) \quad (i = 1, \cdots, N). \tag{95}$$

The constitutive relation that defines the effective conductivity of the nonlinear material reduces to

$$\langle \mathbf{I} \rangle = \mathbf{g}_e^{(l)}(\hat{g}_1^0, \cdots, \hat{g}_N^0)\langle \mathbf{E} \rangle. \tag{96}$$

Note that, Eq. (96) is fully nonlinear because the variables $\hat{g}_i^0$ depend nonlinearly on $\langle \mathbf{E} \rangle$. Since the linear conductivity $\mathbf{g}_e^{(l)}$ is a homogeneous function of degree

one in the conductivity constants $g_i^0$ of the linear comparison material (see also Chapters 2, 4 and 6 of Volume I), then

$$\sum_{i=1}^{N} g_i^0 \frac{\partial \mathbf{g}_e^{(l)}}{\partial g_i^0} = \mathbf{g}_e^{(l)}. \tag{97}$$

Therefore, Eq. (83) is rewritten as

$$\mathcal{H}_e(\langle \mathbf{E} \rangle) \geq \max_{g_i^0 > 0} \sum_{i=1}^{N} \phi_i \left\{ \frac{1}{2} \frac{g_i^0}{\phi_i} \langle \mathbf{E} \rangle \cdot \left[ \frac{\partial \mathbf{g}_e^{(l)}}{\partial g_i^0} \langle \mathbf{E} \rangle \right] - v_i(g_i^0) \right\}, \tag{98}$$

and Eq. (86) implies that

$$\mathcal{H}_e(\langle \mathbf{E} \rangle) \geq \sum_{i=1}^{N} \phi_i e_i(\hat{E}_i), \tag{99}$$

where

$$\hat{E}_i = \left\{ \frac{1}{\phi_i} \langle \mathbf{E} \rangle \cdot \left[ \frac{\partial \mathbf{g}_e^{(l)}}{\partial g_i^0} (\hat{g}_1^0, \cdots, \hat{g}_i^0) \langle \mathbf{E} \rangle \right] \right\}^{1/2} \quad (i = 1, \cdots, N). \tag{100}$$

Finally, the constitutive relation describing the effective behavior of the nonlinear material is written in the following form,

$$\langle \mathbf{I} \rangle = \mathbf{g}_e^{(l)} \left[ \frac{1}{\hat{E}_1} \frac{\partial e_1}{\partial E} (\hat{E}_1), \cdots, \frac{1}{\hat{E}_n} \frac{\partial e_N}{\partial E} (\hat{E}_N) \right] \langle \mathbf{E} \rangle, \tag{101}$$

where $\hat{E}_i$ are functions of the (average) applied field $\langle \mathbf{E} \rangle$, the nonlinear properties of the constituent phases of the material, and the material's microstructure.

## 2.5.2  Strongly Nonlinear Materials with Isotropic Phases

Consider now the class of materials that is defined by Eq. (84) for the phase potentials $e_i$, for which it is possible to simplify further the two equivalent forms (92) and (99). Thus, since

$$\sum_{i=1}^{N} \phi_i g_i^0 (\hat{E}_i)^2 = \langle \mathbf{E} \rangle \cdot \left[ \mathbf{g}_e^{(l)} \langle \mathbf{E} \rangle \right], \tag{102}$$

then, for a power-law material,

$$\mathcal{H}_e(\langle \mathbf{E} \rangle) \geq \sum_{i=1}^{N} \phi_i e_i(\hat{E}_i) = \frac{1}{n+1} \sum_{i=1}^{N} \phi_i g_i^0 (\hat{E}_i)^2, \tag{103}$$

and therefore

$$\mathcal{H}_e(\langle \mathbf{E} \rangle) \geq \frac{1}{n+1} \langle \mathbf{E} \rangle \cdot \left\{ \mathbf{g}_e^{(l)} \left[ \frac{1}{\hat{E}_1} \frac{\partial e_1}{\partial E} (\hat{E}_1), \cdots, \frac{1}{\hat{E}_N} \frac{\partial e_N}{\partial E} (\hat{E}_N) \right] \langle \mathbf{E} \rangle \right\}. \tag{104}$$

If the material's microstructure is statistically isotropic, then writing

$$\mathcal{H}_e(\langle \mathbf{E} \rangle) = \frac{1}{n+1} g_e^{(n)} \langle E \rangle^{n+1},\tag{105}$$

and using (95), we obtain an equation for $g_e^{(n)}$, the effective nonlinear conductivity of the material (Wan *et al.*, 1996; Ponte Castañeda, 1998).

For statistically-isotropic, nonlinear materials, we need to consider only isotropic linear comparison composites with $\mathcal{H}_e^0(\langle \mathbf{E} \rangle) = \frac{1}{2} g_e^{(l)} \langle E \rangle^2$, where $g_e^{(l)}$ is now a scalar function of the nonlinear conductivities $g_i^0$, the volume fractions $\phi_i$, and the material's microstructure. In particular, as discussed in Section 4.6.1.1 of Volume I, for a two-phase material, there are several closely related bounds and estimates for linear materials which can be all characterized in terms of the $d$-dimensional quantity,

$$g_e^{(l)} = \phi_1 g_1^0 + \phi_2 g_2^0 - \frac{\phi_1 \phi_2 (g_1^0 - g_2^0)^2}{\phi_2 g_1^0 + \phi_1 g_2^0 + (d-1)g^0},\tag{106}$$

where $g^0$ takes on different values for different types of estimates. For example, assuming that $g_1^0 > g_2^0$, then,

(1) the limits $g^0 \to \infty$ and 0 correspond, respectively, to the Wiener (one-point) upper and lower bounds.
(2) The choices $g^0 = g_1^0$ and $g_2^0$ yield the two Maxwell–Garnett approximations for particulate microstructures with phases 1 and 2, respectively, in the matrix phase (recall from Section 9.4.9 of Volume I that the Maxwell–Garnett approximation is *not* symmetric with respect to the two phases).
(3) The same choices as in (2) also lead to the Hashin–Shtrikman upper and lower bounds.
(4) If we choose, $g^0 = \zeta_1 g_1^0 + \zeta_2 g_2^0$ and $(\zeta_1/g_1^0 + \zeta_2/g_2^0)^{-1}$, we obtain, respectively, the upper and lower bounds of Beran, in terms of the microstructural parameters $\zeta_1$ and $\zeta_2 = 1 - \zeta_1$.
(5) Finally, the choice $g^0 = g_e^{(l)}$ yields the effective-medium approximation (EMA).

For a two-phase material, one can obtain an expression for $\mathcal{H}_e$ in terms of only one nonlinear equation for the ratio $\hat{g}_1^0/\hat{g}_2^0$. Computing the variables $\hat{E}_1$ and $\hat{E}_2$ in terms of this ratio, the resulting effective nonlinear conductivity $g_e^{(n)}$ is presented in Figure 2.2 for 2D, statistically isotropic, two-phase, power-law conductors with $n = 3$ and $g_2^{(n)}/g_1^{(n)} = 1000$, where $W^{(l)}$ and $W^{(u)}$ correspond to the rigorous upper and lower Wiener bounds for heterogeneous materials with arbitrary microstructures, $MG^{(l)}$ and $MG^{(u)}$ represent the Maxwell–Garnett estimates for particulate microstructures with the less and more conducting materials occupying the matrix phase, respectively. Because the $MG^{(l)}$ estimate for the conductivity of a linear material coincides with the Hashin–Shtrikman lower bound for the set of all statistically isotropic composites, the $MG^{(l)}$ results are identical to the rigorous nonlinear Hashin–Shtrikman lower bound. The results of numerical

FIGURE 2.2. Comparison of various bounds and estimates for the effective nonlinear conductivity $g_e^{(n)}$ of 2D, isotropic, two-phase, power-law conductors with power exponent $n = 3$ and $g_2^{(n)}/g_1^{(n)} = 1000$. Symbols show the results of numerical simulations with random resistor networks (RRN). The lower bound $MG^{(l)}$ obtained from the Maxwell–Garnett approximation is identical with that obtained from the Hashin–Shtrikman lower bound. $MG^{(u)}$ and $B^{(u)}$ denote, respectively, the estimates for the upper bound using the Maxwell–Garnett approximation and the Beran upper bound. (after Ponte Castañeda, 1998).

simulation using a random resistor network (RRN) model, obtained by Wan *et al.* (1996) using a square network, are also shown. The resistor network models will be described and discussed in detail in Chapter 3. The label $B^{(l)}$ represents the rigorous lower bound of Beran for statistically isotropic microstructures with $\zeta_1 = \phi_1$, which is presumably appropriate for symmetric cell microstructures that are similar to the RRN models. The label $B^{(u)}$ denotes the estimate (not a rigorous bound; see above) which is obtained by using the Beran upper bound for the linear comparison composite. As expected, the EMA estimates are in good agreement with the RRN simulations. The Wiener, Hashin–Shtrikman, and Beran bounds progressively narrow the range of possible behavior by introducing, as discussed in detail in Chapters 3, 4, and 7 of Volume I, first-, second- and third-order statistical information about the microstructure of the material, respectively. Although the Maxwell–Garnett approximation and the EMA are generally accurate for particulate- and granular-type microstructures, respectively (see Chapters 4 and

7 of Volume I), the Beran bounds provide a way of estimating the effective properties of more general types of microstructures for which the Maxwell–Garnett and EMA estimates may not be accurate.

## 2.6    Second-Order Exact Results

The method described and utilized so far is most suited for deriving lower bounds and estimates for the upper bounds. These estimates are exact to first order in the contrast between the properties of various phases of a multiphase material. In this section, we describe and discuss another method, developed by Ponte Castañeda and Kailasam (1997), which yields estimates that are exact to second order in the contrasts. As such, they are more accurate than the predictions that are provided by the method described above.

We should mention that Blumenfeld and Bergman (1991b) developed a general method for reducing the solution of the scalar-potential field problems to the solution of a set of linear Poisson-type equations in suitably rescaled coordinates. In particular, for power-law type nonlinearities, they solved explicitly for the effective dielectric constant of a two-phase material to second order in the contrast between the phases' properties. Despite its elegance, their solution yields unphysical results for strong nonlinearity, even when the contrast is not very large, whereas the method described below does not suffer from this shortcoming. We will come back to this point at the end of this section.

The key idea of Ponte Castañeda and Kailasam (1997) is developing a Taylor expansion for the phase energy functions $w_i$, around appropriately defined phase-average electric fields $\langle \mathbf{E}_i \rangle$, so that

$$w_i(\mathbf{E}) = w_i(\langle \mathbf{E}_i \rangle) + \mathbf{I}^{(i)} \cdot (\mathbf{E} - \langle \mathbf{E}_i \rangle) + \frac{1}{2}(\mathbf{E} - \langle \mathbf{E}_i \rangle) \cdot [\hat{\mathbf{g}}^{(i)}(\mathbf{E} - \langle \mathbf{E}_i \rangle)], \quad (107)$$

where $\mathbf{I}^{(i)}$ and $\hat{\mathbf{g}}^{(i)}$ are reference current densities and conductivity tensors with components

$$I^{(i)} = \frac{\partial w_i}{\partial E_i}(\langle \mathbf{E}_i \rangle), \quad g_{jk}^{(i)} = \frac{\partial^2 w_i}{\partial E_i \partial E_j}(\mathbf{E}^{(i)}), \quad (108)$$

where $\mathbf{E}^{(i)}$ is a reference electric field given by, $\mathbf{E}^{(i)} = \lambda^{(i)} \langle \mathbf{E}_i \rangle + [1 - \lambda^{(i)}]\mathbf{E}$, with $0 < \lambda^{(i)} < 1$. We now rewrite Eq. (107) in terms of the average $\langle \mathbf{E} \rangle$ and fluctuating $\mathbf{E}'$ components, $\mathbf{E} = \langle \mathbf{E} \rangle + \mathbf{E}'$:

$$w_i(\langle \mathbf{E} \rangle + \mathbf{E}') = v_i(\langle \mathbf{E} \rangle) + \mathbf{P}_i \cdot \mathbf{E}' + \frac{1}{2}\mathbf{E}' \cdot [\hat{\mathbf{g}}^{(i)}\mathbf{E}'], \quad (109)$$

where

$$v_i(\langle \mathbf{E} \rangle) = w_i(\langle \mathbf{E}_i \rangle) + \mathbf{P}_i \cdot (\langle \mathbf{E} \rangle - \langle \mathbf{E}_i \rangle) - \frac{1}{2}(\langle \mathbf{E} \rangle - \langle \mathbf{E}_i \rangle) \cdot [\hat{\mathbf{g}}^{(i)}(\langle \mathbf{E} \rangle - \langle \mathbf{E}_i \rangle)],$$

$$\mathbf{P}_i = \mathbf{I}^{(i)} + \hat{\mathbf{g}}^{(i)}(\langle \mathbf{E} \rangle - \langle \mathbf{E}_i \rangle). \quad (110)$$

Then, the effective energy $\mathcal{H}_e$ of the nonlinear composite material is given by

$$\mathcal{H}_e(\langle\mathbf{E}\rangle) = \min_{\mathbf{E}'\in S'} \langle v + \mathbf{P}\cdot\mathbf{E}' + \frac{1}{2}\mathbf{E}'\cdot(\hat{\mathbf{g}}\mathbf{E}')\rangle, \tag{111}$$

where $S'$ denotes the set of admissible fields $\mathbf{E}'$, such that, $\mathbf{E}' = \nabla\varphi'$ in the subspace $\Omega$ and $\varphi' = 0$ on $\partial\Omega$, and

$$v(\mathbf{x}) = \sum_{i=1}^{N} m_i(\mathbf{x})v_i, \quad \mathbf{P}(\mathbf{x}) = \sum_{i=1}^{N} m_i(\mathbf{x})\mathbf{P}_i, \quad \hat{\mathbf{g}}(\mathbf{x}) = \sum_{i=1}^{N} m_i(\mathbf{x})\hat{\mathbf{g}}^{(i)}. \tag{112}$$

Equation (111) assumes that the reference fields $\mathbf{E}^{(i)}$ are known in terms of the $\lambda^{(i)}$, which, in general, are functions of the actual electric field $\mathbf{E}$, as well as of the (unknown) $\langle\mathbf{E}_i\rangle$, and therefore the problem posed by (111) for $\mathcal{H}_e$ is nonlinear. However, provided that the second derivatives of $w_i$ in Eq. (107) vary slowly with the $\mathbf{E}^{(i)}$, Eq. (111) suggests, as an approximation, replacing $\mathbf{E}^{(i)}$ by a (as-yet unknown) constant, in which case $\hat{\mathbf{g}}^{(i)}$, $\mathbf{P}_i$ and $v_i$ will also be constant within each phase, hence leading to the following expression for $\mathcal{H}_e$,

$$\mathcal{H}_e(\langle\mathbf{E}\rangle) = \sum_{i=1}^{N} \phi_i v_i(\langle\mathbf{E}\rangle) + \tilde{P}(\langle\mathbf{E}\rangle), \tag{113}$$

where

$$\tilde{P}(\langle\mathbf{E}\rangle) = \min_{\mathbf{E}'\in S'} \langle\mathbf{P}\cdot\mathbf{E}' + \frac{1}{2}\mathbf{E}'\cdot(\hat{\mathbf{g}}\mathbf{E}')\rangle. \tag{114}$$

The interesting feature of Eq. (113) is that it requires only the solution of the linear problem (114) for $\tilde{P}$ which is, physically, equivalent to a problem for a linear conductor with $N$ anisotropic constituents with conductivity tensors $\hat{\mathbf{g}}^{(i)}$ and prescribed polarizations $\mathbf{P}_i$, a problem much simpler to analyze than the original nonlinear problem for $\mathcal{H}_e$. The question then arises as to what the best choices are for these constants.

The optimal choice for each $\langle\mathbf{E}_i\rangle$ is $\langle\mathbf{E}\rangle_i$, the average of the actual field $\mathbf{E}$ over phase $i$:

$$\langle\mathbf{E}_i\rangle = \langle\mathbf{E}\rangle_i, \tag{115}$$

where $\langle\cdots\rangle_i$ denotes a volume average over phase $i$. Although $\langle\mathbf{E}\rangle_i$ cannot be obtained exactly, a consistent estimate for it may be obtained by noting that, $\langle\mathbf{E}\rangle_i = \langle\mathbf{E}\rangle + \langle\mathbf{E}'\rangle_i$, where

$$\langle\mathbf{E}'\rangle_i = \frac{1}{\phi_i}\frac{\partial\tilde{P}}{\partial\mathbf{P}_i} \tag{116}$$

with $\hat{\mathbf{g}}^{(i)}$ held fixed. On the other hand, although the best choice for the $\mathbf{E}^{(i)}$ is not *a priori* clear, given the approximation that was made in deriving (113), the choice

$$\mathbf{E}^{(i)} = \langle\mathbf{E}\rangle_i, \tag{117}$$

is simple and plausible. In particular, Eqs. (115) and (117) are exact for laminated materials, where the fields are constant within each phase. Thus, any solution for the problem posed by (114), together with the associated estimates (116), may be utilized for obtaining corresponding estimates for $\mathcal{H}_e$ via Eq. (113), together with the (self-consistent) equations (115) and (117). Note that $\mathbf{E}^{(i)} = \langle \mathbf{E}_i \rangle$, and, for this reason, $\hat{\mathbf{g}}^{(i)}$ is henceforth denoted by $\mathbf{g}_i$, the phase conductivity tensor. In particular, for two-phase composite materials one can show that

$$\tilde{P}(\langle \mathbf{E} \rangle) = \frac{1}{2} \left[ (\mathbf{g}_e^{(l)} - \langle \mathbf{g} \rangle)(\Delta \mathbf{g})^{-1} \Delta \mathbf{P} \right] \cdot (\Delta \mathbf{g})^{-1} \Delta \mathbf{P}, \qquad (118)$$

from which it follows, using (116), that

$$\langle \mathbf{E} \rangle_1 = \langle \mathbf{E} \rangle + \frac{1}{\phi_1} (\Delta \mathbf{g})^{-1} (\mathbf{g}_e^{(l)} - \langle \mathbf{g} \rangle)(\Delta \mathbf{g})^{-1} \Delta \mathbf{P}, \qquad (119)$$

$$\langle \mathbf{E} \rangle_2 = \langle \mathbf{E} \rangle + \frac{1}{\phi_2} (\Delta \mathbf{g}^{-1})(\mathbf{g}_e^{(l)} - \langle \mathbf{g} \rangle)(\Delta \mathbf{g})^{-1} \Delta \mathbf{P}, \qquad (120)$$

where $\Delta \mathbf{g} = \mathbf{g}_1 - \mathbf{g}_2$, $\Delta \mathbf{P} = \mathbf{P}_1 - \mathbf{P}_2$, $\langle \mathbf{g} \rangle = \phi_1 \mathbf{g}_1 + \phi_2 \mathbf{g}_2$, and $\mathbf{g}_e^{(l)}$ is the effective conductivity tensor of a two-phase *linear* material with phase conductivity tensors $\mathbf{g}_1$ and $\mathbf{g}_2$, volume fractions $\phi_1$ and $\phi_2$, and precisely the same microstructure as the nonlinear composite. This means that *any* estimate that is available for the effective conductivity tensor $\mathbf{g}_e^{(l)}$ of a two-phase linear material, including, for example, the Maxwell–Garnett and EMA estimates, can be used for generating the corresponding estimates for $\mathcal{H}_e$ of a two-phase nonlinear material. Note that the *approximate* estimate of $\mathcal{H}_e$ given by Eq. (113) is a convex function. Since the *exact* expression for $\mathcal{H}_e$ is also known to be convex, it follows that derivatives of the approximate expressions for $\mathcal{H}_e$ should provide a reasonably accurate approximation to the exact constitutive relation.

## 2.6.1 Strongly Nonlinear Isotropic Materials

Consider now a class of two-phase materials for which Eq. (84) describes the constitutive relation. We already saw that for statistically isotropic materials, the macroscopic behavior is described by Eq. (105). For such materials, it is reasonable to assume that the reference fields $\mathbf{E}^{(i)}$ and $\mathbf{I}^{(i)}$ are aligned with the corresponding applied fields. If so, one can define scalar variables $\omega_i$ and $\nu_i$, such that

$$\mathbf{E}^{(i)} = (1 + \omega_i)\langle \mathbf{E} \rangle, \quad \mathbf{I}^{(i)} = (1 + \nu_i)\langle \mathbf{I} \rangle, \qquad (121)$$

from which it follows that $e_i = \langle \mathbf{E} \rangle / \langle E \rangle$ and $\mathbf{I}_i = \langle \mathbf{I} \rangle / \langle I \rangle$, for all the phases $i$. This implies that the conductivity tensor $\mathbf{g}_i$ of all the phases in the linear comparison material has exactly the same symmetry. Then, using (121), we find from Eq. (110) that

$$\mathbf{P}_i = g_i^{(p)} \langle E \rangle^{n-1} \mathbf{E}, \quad \nu_i(\langle \mathbf{E} \rangle) = (1 + n)^{-1} g_i^{(v)} \langle E \rangle^{1+n}, \qquad (122)$$

with

$$g_i^{(p)} = g_i[1 + (1 - n)\omega_i], \tag{123}$$

$$g_i^{(v)} = g_i \left[ 1 + (1 - n)\omega_i + \frac{1}{2}n(n - 1)\omega_i^2 \right]. \tag{124}$$

Then, from Eqs. (84), (113) and (118) we obtain

$$g_e^{(n)} = \langle g^{(v)} \rangle + \frac{n + 1}{2n^2} \left[ mg_e^{(l)} - n\langle g \rangle \right] \left[ \frac{g_1^{(p)} - g_2^{(p)}}{g_1 - g_2} \right]^2, \tag{125}$$

where here $g_i = g_i^{(n)}|1 + \omega_i|^{n-1}$, with $g_i^{(n)}$ being the nonlinear conductivity of phase $i$, and $\langle g^{(v)} \rangle$ defined in a manner analogous to $\langle g \rangle$. The variables $\omega_i$ are determined by Eqs. (115), (119), (120) and (121); they yield, $\omega_1 = \phi_2\omega, \omega_2 = \phi_1\omega$, where

$$\omega = \frac{1}{\phi_1\phi_2} \frac{1}{n^2} \frac{mg_e^{(l)} - n\langle g \rangle}{g_1 - g_2} \frac{g_1^{(p)} - g_2^{(p)}}{g_1 - g_2}. \tag{126}$$

Estimates of $g_e^{(n)}$ based on the Maxwell–Garnett approximation and the EMA can now be obtained by using their corresponding estimates for the effective conductivity tensor $\mathbf{g}_e^{(l)}$, which are in terms of the phase conductivity tensor $\mathbf{g}_i$. Since the Maxwell–Garnett approximation is not symmetric in material's phases, one obtains two classes of Maxwell–Garnett estimates, corresponding to particulate microstructures with the less and more conducting material designated as the matrix phase. On the other hand, due to its symmetry, the estimates provided by the EMA are unique. Moreover, it should be pointed out that the Maxwell–Garnett and EMA estimates for the effective nonlinear conductivity $g_e^{(n)}$ are *not* exactly equivalent. In fact, it can be shown that while for sufficiently weak nonlinearity (i.e., for $n \simeq 1$) these estimates are in close agreement with each other, they can be significantly different for stronger nonlinearities (i.e., as $n \to 0$ or $\infty$). The reason for the differences is associated with the nature of the approximations made in going from the exact estimate (111) for $\mathcal{H}_e$ to the approximation (113), assuming that the reference conductivity tensors $\hat{\mathbf{g}}^{(i)}$ vary slowly with $\mathbf{E}^{(i)}$, so that the replacement of the $\mathbf{E}^{(i)}$ by $\langle \mathbf{E} \rangle_i$ does not introduce significant errors. In what follows, we summarize the results obtained with the Maxwell–Garnett and EMA estimates. The details of derivation of these results, which is straightforward, are given by Ponte Castañeda and Kailasam (1997).

### 2.6.1.1   The Maxwell–Garnett Estimates

The Maxwell–Garnett estimates that correspond to designating the matrix as phase 2 are obtained from the following equations (see Section 4.9.4 of Volume I for the corresponding Maxwell–Garnett equations for *linear* materials),

$$mg_e^{(l)} - n\langle g \rangle = -n\phi_1\phi_2(g_1 - g_2)^2 \left[ \frac{1}{n}\alpha(n)g_2 + \phi_2(g_1 - g_2) \right]^{-1}, \tag{127}$$

where $\alpha(n)$ is a function of $n$, given by

$$\alpha(n) = n\sqrt{n}, \quad \text{2D}, \tag{128}$$

$$\alpha(n) = (n-1) \left( 1 - \frac{1}{\sqrt{n-1}} \arcsin \sqrt{\frac{n-1}{n}} \right)^{-1}, \quad \text{3D}. \tag{129}$$

The 3D expression for $\alpha$ is valid for $n \geq 1$, but the corresponding expressions for $n \leq 1$ may be easily obtained by analytic continuation. Then, Eq. (127), together with Eqs. (125) and (126), provide one of the Maxwell–Garnett estimates for $g_e^{(n)}$. The other Maxwell–Garnett estimate, with the matrix designated as phase 1, is obtained by simply interchanging the roles of 1 and 2.

### 2.6.1.2   Effective-Medium Approximation Estimates

In this case,

$$g_e^{(l)} = \frac{n}{m} \left\{ \frac{\langle g \rangle - m(g_1 + g_2)/\alpha(m)}{2[1 - m/\alpha(m)]} \right.$$
$$\left. + \sqrt{\left\{ \frac{\langle g \rangle - m(g_1 + g_2)/\alpha(m)}{2[1 - m/\alpha(m)]} \right\}^2 + \frac{g_1 g_2}{\alpha(m)/m - 1}} \right\} \tag{130}$$

and

$$g_e^{(l)} = \frac{\langle g \rangle - (g_1 + g_2)/\beta(m)}{2[1 - 1/\beta(m)]}$$
$$+ \sqrt{\left\{ \frac{\langle g \rangle - m(g_1 + g_2)/\beta(m)}{2[1 - m/\beta(m)]} \right\}^2 + \frac{g_1 g_2}{\beta(m)/m - 1}} \tag{131}$$

where

$$\beta(n) = \begin{cases} 1 + \sqrt{n}, & \text{2D}, \\ 2(1-n) \left( 1 - \frac{n}{\sqrt{n-1}} \arcsin \sqrt{\frac{n-1}{n}} \right)^{-1}, & \text{3D}, \end{cases} \tag{132}$$

and $\langle g \rangle = \phi_1 g_1 + \phi_2 g_2$, as before. Equations (130) and (131), obtained from the two independent components of the anisotropic tensor $\mathbf{g}_e^{(l)}$, depend on the functions $\alpha$ and $\beta$ which, in turn, are known functions of the unknown parameter $m$. Therefore, $m$ is obtained by equating (130) and (131). Once $m$ is obtained, $g_e^{(l)}$ and hence $g_e^{(n)}$ are computed.

We now consider the application of these results to estimating the effective nonlinear conductivity of two important classes of heterogeneous materials that we have been studying throughout this book, namely, those with superconducting or insulating inclusions. As we emphasized in Volume I, because these two composites represent two extreme limits of contrast between the properties of the two phases, they provide stringent tests of any theory. In other words, if a theory is reasonably accurate in these limits, it will be even more accurate in less extreme cases.

## 2.6.2  Conductor–Superconductor Composites

It is straightforward to show that, in this limit, Eq. (125) yields the following Maxwell–Garnett estimates for $g_e^{(n)}$:

$$\frac{g_e^{(n)}}{g_2^{(n)}} = \frac{1}{\phi_2^n}\left\{1 + \frac{1}{2}n(n+1)\phi_1\left[\frac{1}{n}\alpha(n) - 1\right]\right\}, \qquad (133)$$

where $0 < n < 1$. The corresponding EMA estimates for $g_e^{(n)}$ are given by Eq. (133) with the factor $\alpha(n)/n - 1$ replaced by $[\alpha(m) - m]/[m - \alpha(m)\phi_1]$, where $m$ is the solution of the equation,

$$m - \alpha(m)\phi_1 = n[1 - \beta(m)\phi_1]. \qquad (134)$$

The EMA estimates are valid for $\phi_1 < 1/\beta(m)$. The limit $\phi_1 = 1/\beta(m)$ defines the percolation thresholds for $g_e^{(n)}$ at which $g_e^{(n)} \to \infty$. The Maxwell–Garnett estimates, on the other hand, do not exhibit any percolation behavior, which is an undesirable aspect of these approximations, as already pointed out in Chapter 4 of Volume I.

Figure 2.3 presents the 3D Maxwell–Garnett, EMA, Hashin–Shtrikman and Wiener estimates for the effective resistivity $R_e^{(n)}/R_2^{(n)}$ of the composite material, as functions of the volume fraction $\phi_1$ of the inclusions, for the power-law exponent $n = 3$. As one might expect, both the Maxwell–Garnett and EMA estimates lie below the Wiener and Hashin–Shtrikman upper bounds. Moreover, it can be shown that the differences between the new Maxwell–Garnett estimates and the old Hashin–Shtrikman bounds (derived earlier in this chapter) increase as $n$ increases, whereas they agree for $n = 1$. As usual, the EMA estimates exhibit sharply the percolation limit at a finite value of $\phi_1$, a distinct advantage of this method.

## 2.6.3  Conductor–Insulator Composites

In this limit [when $g_1^{(n)} \to 0$], Eq. (125) yields the following Maxwell–Garnett estimates for $g_e^{(n)}$:

$$\frac{g_e^{(n)}}{g_2^{(n)}} = \phi_2|1 - \phi_1\omega|^n\left[1 + \frac{1}{2}\phi_1\omega(n-1)\right], \qquad (135)$$

with $\omega = [\phi_1 + \alpha(n) - n]^{-1}$. The corresponding EMA estimates are obtained from Eq. (135) with $\omega = [\phi_1 + n\alpha(m)\phi_2/m - 1]^{-1}$, where $m$ is the root of the following equation

$$m\left[1 + \frac{\phi_1\beta(m)}{1 - \beta(m)}\right] = n\left[1 + \frac{\phi_1\alpha(m)}{m - \alpha(m)}\right]. \qquad (136)$$

The EMA estimates are valid for $\phi_1 \leq 1 - \beta(m)^{-1}$. The limit $\phi_1 = 1 - \beta(m)^{-1}$ defines the percolation threshold at which $g_e^{(n)}$ vanishes.

FIGURE 2.3. The effective resistivity $R_e^{(n)}$ of 3D, isotropic, two-phase, power-law materials, as predicted by the various approximations, versus the volume fraction $\phi_1$ of the supercon-ducting inclusions, with $n = 3$. Note that only the effective-medium approximation indi-cates the existence of a percolation threshold (after Ponte Castañeda and Kailasam, 1997).

Figure 2.4 presents the 3D Maxwell–Garnett, EMA, Hashin–Shtrikman and Wiener estimates for $g_e^{(n)}/g_2^{(n)} = [R_2^{(n)}/R_e^{(n)}]^n$, where $R_e^{(n)}$ is the effective resis-tivity of the material, as functions of the volume fraction $\phi_1$ of the inclusions, for $n = 3$. Both the Maxwell–Garnett and EMA estimates lie below the Wiener up-per bound for $g_e^{(n)}$ [Eq. (81)], while the Maxwell–Garnett estimates lie above the Hashin–Shtrikman lower bound for $R_e^{(n)}$ for particulate microstructures. The EMA estimates that correspond to granular microstructures (which are different from particulate microstructures) are not constrained to satisfy the Hashin–Shtrikman bound and vanish at a finite value of $\phi_1$, the percolation threshold. The difference between the Maxwell–Garnett estimates and the Hashin–Shtrikman lower bound increases with increasing $n$; recall that they are identical in the limit $n = 1$.

It can also be shown that all the nonlinear Maxwell–Garnett and EMA estimates for the effective conductivity or resistivity agree to first order in the volume fraction $\phi_1$, with the result being

$$\frac{R_e^{(n)}}{R_2^{(n)}} = 1 + \gamma(n)\phi_1 + O(\phi_1^2), \tag{137}$$

FIGURE 2.4. Same as in Figure 2.3, but with insulating inclusions (after Ponte Castañeda and Kailasam, 1997).

with

$$\gamma(n) = \frac{1}{n} \left\{ \frac{n+1}{2[\alpha(n) - n]} - 1 \right\}. \tag{138}$$

Analogous expressions can also be derived for the Wiener and Hashin–Shtrikman bounds. Figure 2.5 presents a comparison of $\gamma(n)$ for the MG/EMA estimates versus the Wiener and Hashin–Shtrikman bounds for the 2D materials, along with the numerical results of Lee and Mear (1992), who reported their results for transverse shear of fiber-reinforced power-law ductile composite conductors. As one might expect, the MG/EMA estimates lie above the rigorous Wiener lower bound (for the resistivity) for particulate microstructures. Moreover, the new estimates are in excellent agreement with the numerical estimates of Lee and Mear (1992).

### 2.6.4  General Two-Phase Materials

In addition to the above limits, one may also consider general two-phase power-law materials at arbitrary contrast between the phases. Since the behavior of the general estimate (125) for $g_e^{(n)}$ for $n < 1$ is similar to that of the effective resistivity $R_e^{(n)}$ for $n > 1$, we discuss the various types of estimates for $R_e^{(n)}$ for $n > 1$. In addition, recall that since two types of Maxwell–Garnett estimates are possible

FIGURE 2.5. Dependence of the coefficient $\gamma$ on the power-law exponent $n$, for a 2D, isotropic, two-phase materials with insulating inclusions. Symbols represent the numerical results of Lee and Mear (LM) (1992) (after Ponte Castañeda and Kailasam, 1997).

for a given value of the ratio $R_1^{(n)}/R_2^{(n)}$, depending on whether phase 1 or 2 is designated as the matrix phase (and vice versa for the inclusion phase), we restrict our attention to $R_1^{(n)}/R_2^{(n)} > 1$, and denote by MG1 and MG2 the two estimates corresponding to designating phases 1 and 2, respectively, as the matrix phase.

Ponte Castañeda and Kailasam (1997) showed that as the volume fraction $\phi_1$ of the inclusions increases, the MG2 estimates for $R_e^{(n)}$ also increase. However, the rate of the increase decreases with increasing $n$. In particular, for sufficiently large $n$, there is hardly any increase in $R_e^{(n)}$ over the matrix resistivity $R_2^{(n)}$. The reason for this effect is the fact the current density becomes concentrated in the more conducting matrix phase as $n$ increases, and therefore the effect of the inclusions becomes insignificant. Moreover, as the volume fraction $\phi_1$ of the inclusions increases, the MG1 estimates for $R_e^{(n)}$ decrease, with the rate of the decrease increasing with increasing $n$. In addition, as is the case for the estimates of the linear EMA (see Chapter 4 of Volume I), the nonlinear EMA estimates for $R_e^{(n)}$ agree with the MG1 and MG2 estimates in the limits of small volume fractions of phases 2 and 1, respectively.

The 3D Maxwell–Garnett and EMA estimates for $R_e^{(n)}$ also agree with the corresponding small-contrast asymptotic results of Blumenfeld and Bergman (1991b), which are known to be exact to second order in the contrast, and with the Wiener upper and lower bounds (see above). In fact, the Maxwell–Garnett and EMA es-

timates for $g_e^{(n)}$ and $R_e^{(n)}$ reproduce the asymptotic estimates of Blumenfeld and Bergman:

$$g_e^{(n)} \sim \langle g \rangle - \frac{n+1}{2\alpha(n)} \frac{\langle g^2 \rangle - \langle g \rangle^2}{\langle g \rangle}, \tag{139}$$

$$R_e^{(n)} \sim \langle R \rangle - \frac{n+1}{\beta(n)} \frac{\langle R^2 \rangle - \langle R \rangle^2}{\langle R \rangle}, \tag{140}$$

where $\langle R \rangle = \phi_1 R_1 + \phi_2 R_2$, with $R_1$ and $R_2$ being the resistivities of phases 1 and 2, respectively. The agreement for small enough contrast (to second-order) is a consequence of the fact that the effective behavior of weakly heterogeneous, nonlinear materials with statistically isotropic microstructures is dependent only upon the phase volume fractions (Blumenfeld and Bergman 1991). However, while the small-contrast expansions of Blumenfeld and Bergman (1991) for $g_e^{(n)}$ and $R_e^{(n)}$ diverge as $n \to 0$ and $\infty$, respectively, and can therefore yield unphysical results even at relatively small contrasts, the estimates provided by Eq. (125) do not diverge and always yield physically meaningful results. The new Maxwell–Garnett estimates presented in this section satisfy all the known rigorous bounds, including the Wiener bounds and the Hashin–Shtrikman upper bounds of Ponte Castañeda (1992b) derived earlier in Sections 2.4 and 2.5, and the lower bounds of Talbot and Willis (1994, 1996) for nonlinear composites with statistically isotropic particulate microstructures (with $n \geq 1$).

Finally, let us point out that Gibiansky and Torquato (1998b) derived cross-property bounds that link the effective conductivity of nonlinear disordered materials to their effective elastic moduli. Such cross-property bounds were already described in Section 7.9 of Volume I for linear materials, and will be presented in Chapter 4 for nonlinear composites.

## Summary

In this chapter, we described and discussed general procedures for estimating the effective conductivity and dielectric constant of nonlinear materials. These procedures, which represent the generalization of those described in Chapters 4 and 7 of Volume I for linear materials, provide bounds and estimates for the effective conductivity and dielectric constant. One procedure leads to rigorous bounds and estimates that are exact to first order in the phase property contrast, while the second technique yields estimates that are exact to second order in the contrast.

One important difference between the results obtained by the two procedures must be emphasized. By design, the results presented in Section 10.6 are exact to second-order in the phase contrast, and thus are consistent with the asymptotic results of Blumenfeld and Bergman (1991b), whereas the results presented in Sections 10.4 and 10.5 are nonlinear estimates that are exact only to first order in the phase contrast. On the other hand, while the first-order results provide rigorous bounds for the effective energy function of nonlinear materials (and hence their

generalized effective conductivity and dielectric constant), the second-order esti-mates do not lead to any bound, either the lower or upper bound. Nevertheless, comparison of these estimates with the numerical results and the known bounds suggests that the second-order results provide accurate estimates for the effec-tive conductivity and dielectric constant of general nonlinear materials, and in particular, strongly nonlinear, power-law type composites.

# 3
# Nonlinear Conductivity, Dielectric Constant, and Optical Properties: The Discrete Approach

## 3.0   Introduction

In this chapter we study nonlinear transport and optical properties of heterogeneous materials, representing their morphology by a discrete model. In particular, we consider two-phase materials with percolation disorder which represents a strong type of heterogeneity, although all the theoretical developments that are described in this chapter (and throughout this book) are equally applicable to other types of disorder. As we emphasized in Volume I, we believe that if a theory can provide accurate predictions for transport and optical properties of materials with percolation disorder, i.e., materials in which the contrast between the properties of its two phases is strong, it should also be able to do so for almost any other type of disorder.

There are many transport processes in which the current density is not related to the applied field through a linear relation. Such nonlinearities, in the limit of zero frequency, play an important role in many phenomena, including dielectric breakdown, field dependence of hopping conductivity in heavily-doped semiconductors, and many others. They are, at finite frequencies, the basis of nonlinear optical phenomena in many disordered materials. By suitably tuning of the material's parameters, such as the volume fraction of the conducting material and its nonlinear susceptibility, one can design a wide variety of composite materials with specific properties that have important industrial applications. Chapter 2 described the theoretical methods for estimating the effective conductivity and dielectric constant of nonlinear disordered materials, based on the continuum models. In the present chapter we consider several classes of nonlinear transport processes and describe and discuss, based on the discrete models of heterogeneous materials, the progress that has been made in understanding such phenomena.

## 3.1   Strongly Nonlinear Composites

In most of our discussions in this chapter we use a resistor network model for describing transport in heterogeneous materials. Strongly nonlinear composites are those in which the relation between the current $i$ and the voltage $v$ for any

bond of the network is of power-law type and is given by

$$i = g v^{1/n}, \tag{1}$$

where, as in Chapter 2, we interpret $g$ as a *generalized conductance* of the bond. Equation (1) defines a power-law resistor. If we replace $i$ and $v$ with $q$ and $\Delta P$, the flow rate and pressure drop in a tube or pore of a porous material, then Eq. (1) also defines a power-law fluid, widely used for modeling flow of polymers (Bird *et al.*, 1987). Experimentally, Eq. (1) has been observed in certain classes of conductors, such as ZnO ceramics. More generally, Eq. (1) describes the response of a material when the magnitude of the applied field is very large, so that a linear relation between $i$ and $v$ breaks down completely.

In theoretical analyses of a nonlinear materials, certain subtleties, that are not encountered in linear systems, arise that must be addressed. For example, the nature of boundary conditions that are imposed on the system is very important to the solution of the transport problem. In our discussions in this chapter, we consider only *two-terminal* networks, i.e., those into which one injects a constant current at *one* node and extracts it at another node. Little is known, at least in the context of the problems that we discuss here, about multiterminal networks (i.e., those with more than one injection and one extraction node). It can be shown (see, for example, Straley and Kenkel, 1984) that an equation similar to (1) is also valid for two-terminal networks made of such nonlinear resistors. That is, if $I$, $g_e$ and $V$ are the macroscopic current, effective generalized conductivity and voltage drop in the network, then $I = g_e V^{1/n}$. To prove this, one proceeds as follows (Straley and Kenkel, 1984). One defines a function

$$G_{kj}(v) = \int_0^v i_{kj}(v) dv \tag{2}$$

for each bond $kj$ and constructs a function

$$F = \sum_{k,j} G_{kj}(v_k - v_j). \tag{3}$$

Because $G_{kj}$ has a lower bound, so does the function $F$, and therefore it has a minimum. The existence of this minimum is equivalent to the existence of a solution to Kirchhoff's equations for the resistor network. This can be easily shown by calculating $\partial F / \partial v_k$ and showing that it vanishes at node $k$, hence demonstrating that the net current reaching node $k$ is zero.

However, because this is a nonlinear system, the proof is complete only one also proves that, in addition to existing, the solution to Kirchhoff's equations is also *unique*. This can also be proven (Straley and Kenkel, 1984) by assuming that the function $F$ has two minima for the voltage distributions $\{v_k^{(1)}\}$ and $\{v_k^{(2)}\}$. If so, then $F$ must also have a saddle point at $\{v_k^{(s)}\}$, because along any path in the voltage space that connects the two distributions $\{v_k^{(1)}\}$ and $\{v_k^{(2)}\}$, $F$ must have a maximum. If this saddle point exists, it must be a solution to Kirchhoff's equations, $\partial F / \partial v_k = 0$. However, if the function $i(v)$ [e.g., one that is defined by Eq. (1)] is differentiable with a positive derivative, then it is not difficult to show that the

saddle point cannot exist, since $F$ can be expanded in a series,

$$
\begin{aligned}
F = F^{(s)} + \sum_{k,j} \partial G_{kj}/\partial(v_k - v_j)|_s (v_k - v_j^{(s)} - v_j \\
+ v_j^{(s)}) + \sum_{k,j} \partial^2 G_{kj}/\partial(v_k - v_j)^2|_s (v_k - v_j^{(s)} - v_j + v_j^{(s)})^2 + \cdots
\end{aligned}
\tag{4}
$$

In this equation, the linear term must vanish as $\partial G_{kj}/\partial v_k = 0$, and the quadratic terms are all positive since we assumed that, $di/dv > 0$, and therefore the saddle point does not exist, implying that the function $F$ has a *unique* minimum, i.e., the solution to Kirchhoff's equations is unique. We note that Larson (1981) showed that for slow flow of a power-law fluid in a porous medium *with one injection point and one producing point* (which is the analogue of a two-terminal network) in which flow in each pore is governed by Eq. (1), an equation similar to (1) is also valid at the macroscopic scale, i.e., one has, at the macroscopic scale, $Q = G\Delta P^{1/n}$, or, $Q = G(\Delta P/L)^{1/n}$, where $L$ is the length of the porous medium. The reason that the general form of power-law (1) survives at the macroscopic scale is that, such power-laws are self-similar and therefore they preserve their identity under a microscopic-to-macroscopic transformation (that is, power laws propagate self-similarly).

Calculating the voltage distribution in a nonlinear resistor network is a difficult task, since the nonlinear Kirchhoff's equations may have multiple solutions (all but one of which would be unphysical), and thus one must be careful with the numerical technique used in the simulation (see Uenoyama *et al.*, 1992, for a discussion of this point).

### 3.1.1  Exact Solution for Bethe Lattices

The simplest non-trivial discrete model of strongly nonlinear composites that can be analyzed exactly is a Bethe lattice of coordination number $Z$, which is an endlessly branching network without any closed loops, an example of which is shown in Figure 3.1. We assume that each bond of the Bethe lattice is a power-law conductor. If the lattice contains percolation-type disorder, then the solution of the problem corresponds to the mean-field limit of percolation, i.e., the limit in which the dimensionality of the system is $d \geq 6$. To derive the solution we need the appropriate rules for determining the equivalent conductance of power-law resistors that are in series or parallel. It is not difficult to show that for $N$ power-law resistors in series or parallel, the equivalent conductivity $g_N$ is given by

$$
g_N = \begin{cases} \sum_{i=1}^{N} g_i, & \text{parallel,} \\ \left( \sum_{i=1}^{N} g_i^{-n} \right)^{-1/n}, & \text{series.} \end{cases}
\tag{5}
$$

Suppose now that the bonds' conductances are distributed according to a probability density distribution $f(g)$. We derive an integral equation, from the solution of which all the properties of interest can be computed (Sahimi, 1993a). Consider

FIGURE 3.1. A Bethe lattice of coordination number $Z = 3$.



a branch of a Bethe lattice of coordination number $Z$ which starts at the origin $O$. The conductance of the branch can be computed by simply realizing that, it is the equivalent conductance of one bond, say $OA$, of the branch that starts at $O$ with conductance $g_i$ in series with the branch that starts at $A$ and has a conductance $G_i$. Suppose now that the lattice is grounded at infinity and that a unit voltage has been imposed at $O$. Then, the total conductance $G$ of the network between $O$ and infinity is that of $(Z - 1)$ branches that are in parallel. Therefore,

$$G = \left[ \sum_{i=1}^{Z-1} \left( \frac{1}{g_i^n} + \frac{1}{G_i^n} \right) \right]^{-1/n}. \tag{6}$$

For an infinitely large Bethe lattice, $G$ and $G_i$ are statistically equivalent. Thus, if $H(G)$ represents the statistical distribution of $G$, we must have

$$H(G) = \int \cdots \int \delta \left\{ G - \left[ \sum_{i=1}^{Z-1} \left( \frac{1}{g_i^n} + \frac{1}{G_i^n} \right) \right]^{-1/n} \right\} \prod_{i=1}^{Z-1} f(g_i) H(G_i) dg_i \, dG_i. \tag{7}$$

If we now take the Laplace transform of both sides of Eq. (7), we obtain (Sahimi, 1993a)

$$\tilde{H}(s) = \int_0^\infty \exp(-sG) H(G) dG$$

$$= \left\{ \int \int \exp \left[ -s \left( \frac{1}{g^n} + \frac{1}{G^n} \right)^{-1/n} \right] f(g) H(G) dg dG \right\}^{Z-1}. \tag{8}$$

From the numerical solution of integral equation (8) we obtain all the properties of interest. Note that, in the limit $n = 1$, Eq. (8) reduces to the corresponding integral equation for Bethe lattices with linear resistors which was analyzed in Chapter

5 of Volume I. To our knowledge, no exact solution of Eqs. (8) and (9), for any distribution $f(g)$ and any value of $n$, has been derived.

### 3.1.1.1  Microscopic Versus Macroscopic Conductivity

In general, one may calculate two different effective conductivities for a Bethe lattice. One is $g_m$, the *microscopic* conductivity of the lattice, obtained by grounding the lattice at infinity, imposing a unit voltage at site $O$ of the lattice, and calculating $g_m$ as the current that flows out along one of the outgoing bonds connected to $O$. It is not difficult to see that, aside from a constant factor, $g_m$ is the average $\langle G \rangle$ of the distribution $H(G)$, $g_m = Z\langle G \rangle/(Z - 1)$. Using the properties of the Laplace transform, one can then show (Sahimi, 1993a) that for a Bethe lattice of coordination number $Z$,

$$g_m = Z \left[ \int \int f(g) H(G) \left( \frac{1}{g^n} + \frac{1}{G^n} \right)^{-1/n} dg \, dG \right]^{Z-2}, \qquad (9)$$

which reduces to the equation given by Stinchcombe (1974) and Heinrichs and Kumar (1975) for the $n = 1$ limit, derived in Chapter 5 of Volume I.

Equation (9) is valid for any $f(g)$, the statistical distribution of the bond conductances. Consider then percolation-type disorder, i.e., one for which

$$f(g) = (1 - p)\delta(g) + ph(g), \qquad (10)$$

where $p$ is the fraction of the conducting bonds with conductances that are selected from $h(g)$, which can be any normalized probability density function. It is then not difficult to show that near the percolation threshold $p_c$

$$g_m = \frac{2c(Z-1)^{2+1/n}}{h_n^{1/n}(Z-1)} \left[ \frac{1}{(Z-1)^n - 1} \right]^{1/n} [n\Gamma(J - n - 1)]^{1/n} (p - p_c)^{1+1/n}. \qquad (11)$$

In Eq. (11), $c$ is a constant of order unity, $\Gamma$ is the gamma function, $J = 2 + [n]$, where $[n]$ denotes the integer part of $n$, $p_c = 1/(Z - 1)$ is the percolation threshold of the Bethe lattice, and

$$h_n = \int_0^\infty \frac{h(g)}{g^n} \, dg.$$

The power law implied by Eq. (11) for the dependence on $p$ of $g_m$ near $p_c$ was first derived by Straley and Kenkel (1984), except that they did not provide the exact form of the numerical factor given by Eq. (11). Equation (11) predicts that, for the linear ($n = 1$) limit, one has

$$g_m \propto (p - p_c)^2. \qquad (12)$$

On the other hand, the *macroscopic* or effective conductivity $g_e$, which is what one usually calculates for 2D or 3D networks, is the average current density per unit applied field. The difference between the two cases is due to the geometry of the Bethe lattice, which has a peculiar structure (lacking any closed loops while keeping the length of the bonds constant), and the boundary conditions at infinity

(Straley, 1977). To estimate $g_e$ one may proceed as follows (Straley and Kenkel, 1984). The average power $\mathcal{P}$ dissipated per unit volume is given by

$$\mathcal{P} = \frac{IV}{AL} = g_e \left(\frac{V}{L}\right)^{1+1/n}, \qquad (13)$$

where $A$ and $L$ are, respectively, the surface area and linear size of the sample. The voltage difference across a chain of the lattice is controlled by the geometrical distance $\xi_p$ between the ends of the chain, where $\xi_p$ is the correlation length of percolation. In general, $\xi_p$ is less than $\mathcal{L}$, the length of the chain, since the chain is twisted. However, in a Bethe lattice, the chain performs a random walk in space, implying that, $\xi_p^2 \sim \mathcal{L}$, and therefore the current $I_c$ that is carried by a chain is given by

$$I_c = \left(\frac{\xi_p V}{L\mathcal{L}}\right)^{1/n} = (p - p_c)^{1/2n} \left(\frac{V}{L}\right)^{1/n}. \qquad (14)$$

However, the chain will carry no current at all unless its ends are connected to the sample-spanning percolation cluster. To calculate the probability of this connection, we note that the probability that a given site is connected by a particular bond to the sample-spanning cluster is $P(p)$, the percolation probability, and therefore, near $p_c$, the two ends of the chain are connected to the cluster with a probability $P^2(p) \sim (p - p_c)^{2\beta}$. As $\beta = 1$ for the Bethe lattice, we find that the probability that the chain is connected to the sample-spanning cluster is proportional to $(p - p_c)^2$. Therefore, the dissipated power is given by

$$\mathcal{P} = P^2(p)[(p - p_c)^{1/2n}(V/L)^{1/n}]^{1+n} = (p - p_c)^{(5+1/n)}(V/L)^{1+1/n}, \quad (15)$$

which, when compared with Eq. (13), implies that, near $p_c$,

$$g_e \sim (p - p_c)^{(5+1/n)/2}. \qquad (16)$$

Observe that the critical exponents that characterize the near threshold behavior of both $g_m$ and $g_e$ depend on $n$. In particular, Eq. (16) indicates that if, in general, near $p_c$ one has

$$g_e \sim (p - p_c)^{\mu(n)}, \qquad (17)$$

where $\mu(n)$ is the analogue of the conductivity critical exponent $\mu$ for the linear case; that is, for linear resistor networks near $p_c$ one has

$$g_e \sim (p - p_c)^{\mu}, \qquad (18)$$

then in the mean-field approximation (the solution of which is obtained by solving the problem on a Bethe lattice) $\mu(n) = \mu_n$ is given by

$$\mu_n = \frac{1}{2}(5 + n^{-1}), \qquad (19)$$

which implies that, in the linear ($n = 1$) limit, one has

$$g_e \sim (p - p_c)^3, \qquad (20)$$

in agreement with the result derived in Chapter 5 of Volume I.

The effective conductivity of 3D linear resistor networks near $p_c$ follows the power law (18) with $\mu \simeq 2.0$. Therefore, Eq. (12) is similar to the power-law behavior of the effective linear conductivity of 3D networks. Because by varying the coordination number $Z$ of the Bethe lattice, its percolation threshold, $p_c = 1/(Z-1)$, can be adjusted to closely match that of a 3D network (for example, the percolation threshold of a Bethe lattice with $Z = 5$ is $p_c = 1/4$, which is essentially the same as the bond percolation threshold of a simple-cubic lattice, $p_c \simeq 0.249$), it is clear that for linear transport (the limit $n = 1$) $g_m$ should provide an excellent approximation to the conductivity of 3D networks (Heiba $et\ al.$, 1982, 1992) and this has been shown to be indeed the case (Sahimi, 1993b). For power-law transport considered here one may also use $g_m$ as an approximation to the effective nonlinear conductivity of 3D networks. Figure 3.2 compares the conductivity $g_m$ obtained from the numerical solution of Eqs. (8) and (9) with $Z = 5$ with that of a simple-cubic network obtained by Monte Carlo calculations, and it is clear that the agreement between the two is very good.



FIGURE 3.2. Comparison of the microscopic conductivity of a Bethe lattice of coordination number $Z = 5$ with the effective conductivity of a simple-cubic network obtained by Monte Carlo simulations (dashed curve). The bonds of the two lattices are power-law resistors with a power-law exponent $n = 0.4$. The other two curves are, from top to bottom, the predictions of Eqs. (32) and (31) (after Sahimi, 1993a).

### 3.1.1.2    Effective-Medium Approximation for Bethe Lattices

Using Eq. (8), one can also construct an effective-medium approximation (EMA) for power-law electrical transport in a Bethe lattice. As pointed out in Section 5.3.2 of Volume I, in the effective-medium approach, the probability distribution $H(G)$ is expected to achieve its maximum around a mean value $G^*$, and thus we may approximate $H(G)$ by $H(G) \simeq \delta(G - G^*)$, so that $\tilde{H}(s) = \exp(-sG^*)$. Then, Eq. (8) becomes

$$\exp(-sG^*) = \left\{ \int \int \exp\left[ -s \left( \frac{1}{g^n} + \frac{1}{(G^*)^n} \right)^{-1/n} \right] f(g) dg \right\}^{Z-1}. \quad (21)$$

To determine $G^*$, we take the derivative of Eq. (21) with respect to $s$ and evaluate the result at $s = 0$; we find that

$$\int_0^\infty \left\{ \left[ \frac{1}{g^n} + \frac{1}{(G^*)^n} \right]^{-1/n} - \frac{G^*}{Z-1} \right\} f(g) dg = 0. \quad (22)$$

The effective conductivity $g_e$ of the network is obtained if we set in Eq. (22), $f(g) = \delta(g - g_e)$ (because in the EMA approach, all bonds of the network have the same conductance $g_e$), in which case Eq. (22) yields, $(G^*)^n = g_e^n[(Z-1)^n - 1]$. Substituting this result in Eq. (22) yields the desired EMA (Sahimi, 1993a):

$$\int_0^\infty \left[ \frac{(Z-1)g}{\{g^n + [(Z-1)^n - 1]g_e^n\}^{1/n}} - 1 \right] f(g) dg = 0. \quad (23)$$

Typical of all the EMAs, and similar to the EMAs derived in Volume I for the effective linear properties, Eq. (23) provides accurate estimates of $g_e$ if the disorder is not too strong, implying that the EMA cannot be very accurate near $p_c$.

## 3.1.2    Effective-Medium Approximation for Three-Dimensional Materials

Unlike the EMA for linear electrical transport which was derived and discussed in Chapters 5 and 6 of Volume I, derivation of an EMA for the nonlinear transport is not unambiguous. In particular, several of such approximations have been proposed in the past in order to estimate the effective conductivity of random resistor networks with power-law conductances, all of which are purported to represent some sort of an EMA. We should point out, however, that any reasonable EMA (and similar approximations) should possess two important properties.

(1) It should reduce, in the limit $n = 1$, to the well-known EMA for linear random resistor networks derived and analyzed in Chapters 5 and 6 of Volume I:

$$\int_0^\infty \frac{g - g_e}{g + (Z/2 - 1)g_e} f(g) dg = 0. \quad (24)$$

(2) It should predict the same bond percolation threshold, $p_c = 2/Z$, that the linear EMA predicts, as the location of the percolation thresholds is independent of $n$.

One of the first EMAs for resistor networks with power-law conductances was proposed by Sahimi (1993a), and is given by

$$\int_0^\infty \left\{ \frac{gZ/2}{[g^n + ((Z/2)^n - 1)G^n]^{1/n}} - 1 \right\} f(g)dg = 0, \qquad (25)$$

which reduces to Eq. (24) in the limit $n = 1$, as it should. Another EMA was derived by Tua and Bernasconi (1988) for a 2D isotropic continuum (with circular inclusions), which was extended (Sahimi, 1993a) to networks of random conductances with coordination number $Z$. In this approach one first defines a *tangent* or *differential* conductance $\sigma$ by

$$\sigma = \frac{di}{dv}, \qquad (26)$$

which, in the limit $n = 1$, yields the usual $\sigma = g$. Equation (26), when combined with (1), yields

$$\sigma = \frac{g}{n} v^{(1-n)/n}. \qquad (27)$$

Consider now a two-phase material with its phase tangent conductances being $\sigma_1$ and $\sigma_2$, both of which depend on the voltage $v$. Recall from Chapter 5 of Volume I that in the EMA approach one inserts in the effective medium a bond with its true conductance and determines the voltage fluctuations along this bond, i.e., the extra voltage in the effective medium generated by the replacement of the conductance of the bond in the effective medium by its true value. Carrying out this replacement for component $j$ ($j = 1, 2$) yields

$$v_j = \frac{\sigma_e(v_1, v_2)Z/2}{\sigma_j(v_j) + \sigma_e(Z/2 - 1)} v_e, \qquad (28)$$

where $v_e$ is the voltage along the bond in the effective medium, and $\sigma_e$ is the effective value of $\sigma$. If we now apply the usual idea of an EMA, namely, that the average of $v_j$ must be equal to $v_e$ (or that the average of the voltage fluctuations must be zero), we obtain

$$\int_0^\infty \frac{\sigma_j(v_j) - \sigma_e}{\sigma_j + \sigma_e(Z/2 - 1)} f(\sigma_j)d\sigma_j = 0, \qquad (29)$$

which is the same as Eq. (24) except that the conductances $\sigma_j$ and $\sigma_e$ are functions of the voltage. If the composite consists of two phases with (volume) fractions $p$ and $(1 - p)$, then

$$pv_1 + (1 - p)v_2 = v_e. \qquad (30)$$

The generalization of Eq. (30) to an $N$-component system is obvious. Equations (29) and (30) are then used for determining $\sigma_e$. Having determined this quantity, we calculate $g_e$ using Eq. (27).

To test the accuracy of these two approximations, let us consider a simple case, namely, a resistor network with a percolation-type conductance distribu-

tion, $f(g) = (1 - p)\delta(g) + p\delta(g - 1)$. In this limit, Eq. (25) reduces to (Sahimi, 1993a)

$$g_e = \left[ \frac{(pZ/2)^n - 1}{(Z/2)^n - 1} \right]^{1/n}, \tag{31}$$

while Eqs. (29) and (30) predict that (Sahimi, 1993a)

$$g_e = p^{(n^2-1)/n} \left( \frac{p - 2/Z}{1 - 2/Z} \right)^{1/n}. \tag{32}$$

Equations (31) and (32) do meet the two criteria that we set above, namely, that they both reduce to the linear EMA for $n = 1$, and their predictions for the percolation threshold are the same as in the case of linear transport: Both equations predict that $g_e$ vanishes at $p = p_c = 2/Z$, the same as that predicted by Eq. (24) for linear transport. We can also compare the predictions of these EMAs with those for the effective microscopic conductivity of the Bethe lattice. For example, for $n = 1/2$ Eq. (7) predicts that $\mu_n = 3$, whereas the numerical estimate for 3D systems (see below) for $n = 1/2$ is $\mu_n \simeq 2.35$. However, unlike the two EMAs described above, the region near $p_c$ in which the conductivity of a Bethe lattice is different from that of a 3D network is so narrow that it can hardly be detected (see Figure 3.2).

Consider now the case in which the nonlinear composite material obeys a current-field response of the following form

$$\mathbf{I} = g|\mathbf{E}|^{1/n}\mathbf{E}, \tag{33}$$

which is a slight generalization of Eq. (1). Bergman (1989) and Lee and Yu (1995) developed an EMA for computing the effective conductivity of this type of composite materials. Bergman developed an EMA for any value of $n$, while Lee and Yu considered only the $n = 1/2$ limit. In both cases a 2D continuum model (but with percolation disorder) in which inclusions, consisting of long cylinders (or circles) of nonlinear conductance $g_\alpha$ ($\alpha = i, h$), representing the inclusion and the host matrix, were embedded in an effective medium with a nonlinear conductance $g_e$. As usual, one applies a uniform far field $\mathbf{E}_0$, calculates the local field $\mathbf{E}_\alpha$, and insists that $\langle \mathbf{E}_\alpha \rangle = \mathbf{E}_0$. We supplement Eq. (33) by the usual electrostatic equations, namely,

$$\nabla \cdot \mathbf{I} = 0, \quad \nabla \times \mathbf{E} = \mathbf{0}. \tag{34}$$

Then, there exists a potential $\varphi$ such that

$$\mathbf{E} = -\nabla\varphi. \tag{35}$$

If the potential $\varphi$ is known, then, one can calculate $\mathbf{E}_\alpha$. Trial functions of the following form,

$$\varphi_\alpha(r, \theta) = -E_0(1 - b_\alpha)r \cos\theta, \quad r < R, \tag{36}$$

$$\varphi_e(r, \theta) = -E_0(r - b_\alpha R^2/r) \cos\theta, \quad r > R, \tag{37}$$

are now selected, where $b_\alpha$ is a variational parameter, and $R$ is the radius of the cylinder. With these choices, the energy functional of the composite is given by

$$\mathcal{H}_\alpha = \left[ g_e + p_\alpha g_e \left( -1 + 4b_\alpha + 4b_\alpha^2 + \frac{1}{3}b_\alpha^4 \right) + p_\alpha g_e (1 - b_\alpha)^4 \right] V_0^4, \quad (38)$$

where $p_\alpha$ is the volume fraction of material of type $\alpha$. If we now define $y_\alpha = g_\alpha/g_e$ and minimize the energy functional, we obtain

$$(1 + y_\alpha)b_\alpha^3 - 9y_\alpha b_\alpha^2 + 3(2 + 3y_\alpha)b_\alpha + 3(1 - y_\alpha) = 0, \quad (39)$$

which provides an equation for $b_\alpha$ and $\varphi_\alpha$, and hence $\mathbf{E}_\alpha$. If the system is such that inclusions of nonlinear conductivity $g_i$ and volume fraction $p_i$ are randomly distributed in a host of conductivity $g_h$ with volume fraction $p_h$ ($p_i + p_h = 1.0$), then the EMA equation is simply given by

$$p_i b_i(y_i) + p_h b_h(y_h) = 0. \quad (40)$$

Figure 3.3 compares the predictions of this EMA with the results of numerical simulation, demonstrating the accuracy of the predictions.



FIGURE 3.3. Effective nonlinear conductivity $g_e^{(n)}$, normalized by the effective conductivity of the system in the linear regime, versus the fraction $p$ of the good conducting bonds. Solid curves are the predictions of the EMA, Eqs. (39) and (40), while symbols show the results of numerical simulations. The results are, from top to bottom, for conductivity ratios $y = 0.5, 0.1, 0.01$ and $0.001$ (after Lee and Yu, 1995).

### 3.1.3   The Decoupling Approximation

Equation (29) is quite general and can be used with a variety of composites. For example, Wan *et al.* (1996) analyzed a general two-phase composite consisting of materials $a$ and $b$ with volume fractions $p$ and $(1 - p)$, respectively, such that the constitutive equation that related the current density $\mathbf{I}$ to the electric field $\mathbf{E}$ was given by Eq. (33). The effective generalized conductivity $g_e$ is then defined by the usual equation, $\langle\mathbf{I}(\mathbf{x})\rangle = g_e|E_0|^{1/n}\mathbf{E}_0$, where $\langle\cdot\rangle$ denotes an average over the volume of the system. For each region $i$ of the composite ($i = a$ or $b$), the **I-E** relation is approximated by, $\mathbf{I}(\mathbf{x}) = g_i\langle|E(\mathbf{x})|^{1/n}\rangle_i\mathbf{E}(\mathbf{x}) \equiv \sigma_i\mathbf{E}(\mathbf{x})$, where $\langle\cdot\rangle_i$ denotes an average over volume of region $i$. Similarly, for the composite as a whole, one can define, $\mathbf{I} = g_e\langle|E(\mathbf{x})|^{1/n}\rangle\mathbf{E}(\mathbf{x}) \equiv \sigma_e\mathbf{E}(\mathbf{x})$. Therefore, similar to our discussion presented above, the composite is treated as a *linear* material, but with field-dependent conductivities $\sigma_a$ and $\sigma_b$. It is not difficult to show that

$$\langle E^2\rangle_i = \frac{1}{p_i}\frac{\partial\sigma_e}{\partial\sigma_i}E_0^2, \tag{41}$$

where $p_i$ is the volume fraction of phase $i$ [$p_i = p$ or $(1 - p)$]. One can also use a *decoupling approximation* (Stroud and Wood, 1989) according to which,

$$\langle|E|^{1/n}\rangle_i \simeq \langle|E|^2\rangle_i^{1/2n}, \tag{42}$$

so that the right-hand side of Eq. (41) is only a function of $\langle|E|^2\rangle_i$. Therefore, Eq. (41), when written for both phases $a$ and $b$, forms a set of coupled self-consistent equations, the solution of which yields $\langle E^2\rangle_a/E_0^2$ and $\langle E^2\rangle_b/E_0^2$. Given these two quantities and Eq. (29), the effective generalized conductivity $g_e$ is then estimated.

As an example, consider a 2D system. With $f(\sigma) = p\delta(\sigma - \sigma_a) + (1 - p)\delta(\sigma - \sigma_b)$ and $Z = 4$, Eq. (29) yields

$$g_e = \frac{\sigma_e}{E_0^{1/n}}$$

$$= \frac{1}{2E_0^{1/n}}\left\{(1 - 2p)(X_b - X_a) + \left[(1 - 2p)^2(X_b - X_a)^2 + 4X_aX_b\right]^{1/2}\right\}, \quad X_i = g_i\langle|E|^2\rangle_i^{1/2n}, \tag{43}$$

and from Eq. (41) one obtains, for example,

$$\langle E^2\rangle_a = \frac{E_0^2}{2p}\left\{(2p - 1) + \frac{2X_b - (1 - 2p)^2(X_b - X_a)}{[(1 - 2p)^2(X_b - X_a)^2 + 4X_aX_b]^{1/2}}\right\}. \tag{44}$$

It can then be shown that Eq. (43) is identical with the Hashin–Shtrikman lower bound for $g_e$, derived by Ponte Castaneda *et al.* (1992) and described in Chapter 2. Numerical simulations of the problem indicated close agreement with the predictions of Eq. (43).

Two other methods that have been proposed for treating the problem of conductivity of a nonlinear material embedded in a matrix are the perturbation expansion and the variational approach. Normally, these methods are described as part of the

continuum approach to these problems. However, since they were developed for materials with percolation disorder, we describe them here, rather than in Chapter 2. What follows is a brief description of each method.

### 3.1.4  Perturbation Expansion

In this approach, which was developed by Gu and Yu (1992), Yu and Gu (1992), and Yu *et al.* (1993), the expansion parameter is the nonlinear conductance $g_h$ of the host or the matrix into which the inclusions are embedded. The electrostatic potentials $\varphi^i$ and $\varphi^h$ for the inclusion and the host are expanded as

$$\varphi^i = \varphi_0^i + g_h\varphi_1^i + g_h^2\varphi_2^i + \cdots \tag{45}$$

$$\varphi^h = \varphi_0^h + g_h\varphi_1^h + g_h^2\varphi_2^h + \cdots \tag{46}$$

If $\Upsilon = |\mathbf{E}|^{1/n}$, one writes down an expansion for $\Upsilon^h = \Upsilon_0^h + g_h\Upsilon_1^h + g_h^2\Upsilon_2^h + \cdots$, with a similar expansion for $\Upsilon^i$. For example, for $n = 1/2$ one obtains

$$\Upsilon^h = (\nabla\varphi_0^h)^2 + 2g_h(\nabla\varphi_0^h) \cdot (\nabla\varphi_1^h) + g_h^2(\nabla\varphi_1^h)^2 + \cdots \tag{47}$$

with a similar expression for $\Upsilon^i$. The current density functions $\mathbf{I}^h$ and $\mathbf{I}^i$ can also be expanded in powers of $g_h$, $\mathbf{I}^h = \mathbf{I}_0^h + g_h\mathbf{I}_1^h + g_h^2\mathbf{I}_2^h \cdots$, with a similar expression for $\mathbf{I}^i$, since they can be expressed in terms of $\varphi_j^h$ and $\varphi_j^i$. When all of the expansions are substituted into Eqs. (34), one obtains sets of simultaneous equations for the functions $\varphi_j^i$ and $\varphi_j^h$ for $j = 1, 2, \cdots$ Then, specifying the shape of the inclusion and the boundary conditions allows one to solve for these functions, and thus obtain the overall nonlinear effective conductivity of the material. However, such perturbation expansions are not very accurate, particularly for percolation disorder, unless many terms of the expansion are computed. In fact, they break down and predict unphysical results if the nonlinearity is strong, e.g., if the applied field $\mathbf{E}_0$ is very large, since in this case the linear response vanishes identically in some regions of the composite.

### 3.1.5  Variational Approach

Yu and Gu (1994,1995) considered a class of strongly nonlinear composites that follow Eq. (33) with $n = 1/2$, where the nonlinear conductance $g$ takes on different values in the inclusions and in the host. Their approach is different from what we described in Chapter 2, and is closer to what is of interest to us in the present chapter. Yu and Gu considered the dilute limit in which a single cylindrical inclusion of volume $\Omega_i$ is inserted in a host medium with a larger volume $\Omega$. With $n = 1/2$, Eqs. (33)–(35) yield

$$\nabla \cdot [g(\mathbf{x})|\nabla\varphi(\mathbf{x})|^2\nabla\varphi(\mathbf{x})] = 0. \tag{48}$$

One now invokes the variational principle (see Chapter 2) to minimize the energy functional,

$$\mathcal{H}[\varphi] = \int_\Omega \mathbf{I} \cdot \mathbf{E}(\mathbf{x})d\Omega, \tag{49}$$

with respect to an arbitrary variation $\delta\varphi(\mathbf{x})$ away from the solution of Eq. (48), provided that $\delta\varphi$ vanishes on the surface of the inclusions. When the minimum condition is satisfied by a trial function $\hat{\varphi}$, the effective nonlinear conductivity is obtained from

$$g_e E_0^4 \Omega = \hat{\mathcal{H}} = \int_\Omega g(\mathbf{x})|\hat{\mathbf{E}}(\mathbf{x})|^4 d\Omega, \tag{50}$$

where $\hat{\mathbf{E}} = \nabla\hat{\varphi}$. Thus, it remains to develop suitable trial potential functions $\hat{\varphi}$.

The trial functions must be selected so as to satisfy the symmetry of the system and the boundary conditions that are imposed on it. Thus, if the inclusions are cylindrical, then, the trial functions, similar to Eqs. (36) and (37), are expansions in $\cos m\theta$ (with $m = 1, 3, 5, \cdots$), whereas for spherical inclusions one must use Legendre functions. If the trial functions are selected to be Eqs. (36) and (37) (which involve only the parameter $b_\alpha$), then Eq. (39) is obtained again. Yu and Gu (1995) improved the accuracy of the method by using higher-order terms in the expansions. Hence, for a cylindrical inclusion of radius $R$, they used

$$\varphi_i(r, \theta) = (c_{11}r + c_{13}r^3/R^2 + c_{15}r^5/R^4)\cos\theta$$
$$+ (c_{31}r + c_{33}r^3/R^2 + c_{35}r^5/R^4)\cos 3\theta$$
$$+ (c_{51}r + c_{53}r^3/R^2 + c_{55}r^5/R^4)\cos 5\theta, \quad r < R, \tag{51}$$

for the inclusion phase, and

$$\varphi_h(r, \theta) = r\cos\theta + (b_{11}R^2/r + b_{13}R^4/r^3 + b_{15}R^6/r^5)\cos\theta$$
$$+ (b_{31}R^2/r + b_{33}R^4/r^3 + b_{35}R^6/r^5)\cos 3\theta$$
$$+ (b_{51}R^2/r + b_{53}R^4/r^3 + b_{55}R^6/r^5)\cos 5\theta, \quad r > R, \tag{52}$$

where the external voltage has been set to be, $E_0 = 1$. Thus, the problem involves determining 18 variational parameters, the $b_i$ and $c_i$. By using the boundary condition for the potential $\varphi$ on the surface of the cylinder (at $r = R$), three relations between the 18 coefficients are found. Then, Eq. (49) is used to compute $\mathcal{H}$, and the result is then minimized with respect to the remaining 15 parameters. Compared to the case in which Eqs. (36) and (37) are used, this procedure with 18 variational parameters improves the accuracy of the predictions by about 10%.

## 3.1.6  Exact Duality Relations

In Chapters 4 and 5 of Volume I we described duality relations for the effective conductivity of linear materials. We now consider the same relations for nonlinear materials that are characterized by Eq. (33). Note that, in the notation of Eq. (33), the limit $n = \infty$ corresponds to the linear conduction case [whereas Eq. (1) reduces to the linear problem in the limit $n = 1$]. Recall that duality relations exist only for 2D systems, and therefore only such materials (for example, thin films) are considered here. We consider composites in which the nonlinear conductivity $g$

varies from phase to phase, but the exponent $n$ is the same for all the components. The duality relations that we describe here is due to Levy and Kohn (1998).

Consider a two-phase composite in which the local conductivities are defined by

$$g_j(|\mathbf{V}|) = g_j|\mathbf{V}|^{1/n}, \quad j = 1, 2. \tag{53}$$

The dual composite is another two-phase material with the *same* morphology, but with phases that have the following local conductivity,

$$g_j(|\mathbf{I}|) = \frac{1}{g_j(|\mathbf{V}|)} = g_j^{-n/(n+1)}|\mathbf{I}|^{-1/(n+1)}, \quad j = 1, 2. \tag{54}$$

The effective conductivities of the two materials are expressed as

$$g^*[g_1(|\mathbf{V}|), g_2(|\mathbf{V}|); V_0] = g_e V_0^{1/n}, \tag{55}$$

and

$$g_d^*[g_1(|\mathbf{I}|), g_2(|\mathbf{I}|); I_0] = g_e^{(d)} I_0^{-1/(n+1)}, \tag{56}$$

where $g_e^{(d)}$ is the effective conductivity of the dual composite, and

$$I_0 = g^*[g_1(|\mathbf{V}|), g_2(|\mathbf{V}|); V_0]V_0, \tag{57}$$

is the magnitude of the current that flows through the primal composite, which is also the magnitude of the volume-averaged electric field in the dual composite. Since the effective conductivities of the dual materials satisfy, $g_e V_0^{1/n} = 1/[g_e^{(d)} I_0^{-1/(n+1)}]$, we obtain an exact duality relation for heterogeneous (2D) materials made of power-law conductors:

$$g_e^{n/(n+1)} = \frac{1}{g_e^{(d)}}. \tag{58}$$

We may consider the consequences of duality for percolation composites by studying two limiting cases:

(1) A mixture of good conductors (nonlinear conductivity $g_M$, exponent $n$) and perfect insulators (nonlinear conductivity $g_I = 0$). Then, an equation similar to (17) must hold near the percolation threshold $p_c$ of the good conductor.
(2) A mixture of normal conductors (nonlinear conductivity $g_I$, exponent $n$) and superconductors (nonlinear conductivity $g_M = \infty$). Then, similar to linear resistor networks of conductors-superconductors for which one has, near $p_c$, $g_e \sim (p_c - p)^{-s}$, we expect to have

$$g_e \sim (p_c - p)^{-s_n}, \tag{59}$$

where $s_n = s(n)$ is the analogue of the exponent $s$, defined above. Therefore, if we take phase 2 to be a perfect insulator, then, the dual composite is a mixture of normal conductors and superconductors. Using Eq. (59), we then find that

(Straley and Kenkel, 1984; Levy and Kohn, 1998)

$$\mu(n) = \frac{n}{n+1} \, s\left(-\frac{1}{n+1}\right), \tag{60}$$

which, in the limit $n \to \infty$, reduces to the well-known relations, $\mu = s$, for 2D linear percolation conductivity which was already mentioned in Chapters 2 and 5 of Volume I. Let us emphasize that the exponent $\mu_n = \mu(n)$ used in Eq. (60) is slightly different from that in Eq. (17).

When the ratio of the conductivities of the two components is finite, we expect, similar to linear resistor networks studied in Chapters 2, 5, and 6, to have a scaling representation of $g_e$:

$$g_e \sim g_M (p - p_c)^{\mu_n} \Phi_\pm(z), \quad z = \frac{g_I/g_M}{(p - p_c)^{\mu_n + s_n}}, \tag{61}$$

where the plus (minus) sign is for $p > p_c$ ($p < p_c$). Thus, returning to our two-phase composite with conductivities $g_1$ and $g_2$, we find that when $g_1 \gg g_2 > 0$, then, the primal composite has an effective nonlinear conductivity given by

$$g_e \sim g_1 (p - p_c)^{\mu_n} \Phi\left[\frac{g_2/g_1}{(p - p_c)^{\mu_n + s_n}}\right]. \tag{62}$$

The dual of this composite has local nonlinear conductivities $g_1^{-n/(n+1)} \ll g_2^{-n/(n+1)}$, and therefore

$$g_e^{(d)} \sim g_2^{-n/(n+1)} (p - p_c)^{\mu_n} \Phi_{-n/(n+1)}\left[\frac{(g_2/g_1)^{n/(n+1)}}{(p - p_c)^{\mu_n + s_n}}\right]. \tag{63}$$

Therefore, using the duality relation, Eq. (59), we find that the scaling functions for the primal composite and its dual satisfy an exact relation:

$$\Phi_{-1/(n+1)}[z^{n/(n+1)}] = \left[\frac{z}{\Phi_{1/n}(z)}\right]^{n/(n+1)}, \tag{64}$$

with the understanding that if the left-hand side of Eq. (64) uses $\Phi_{-1/(n+1)}$ with a plus sign [see Eq. (61)], then, the right-hand side uses $\Phi_{1/n}$ with a minus sign, and vice versa.

### 3.1.7  Scaling Properties

The critical exponent $\mu$, defined by Eq. (18), that characterizes the power-law behavior of the effective linear conductivity $g_e$ of percolation composites near the percolation threshold [see Eq. (2.74)], can be expressed as

$$\mu = (d - 2)\nu + \zeta, \tag{65}$$

where $\zeta$ is the linear resistance exponent (that is, the resistance $R$ of a sample of linear size $L < \xi_p$ scale as, $R \sim L^{\zeta/\nu}$), $\nu$ is the critical exponent of percolation

correlation length, $\xi_p \sim |p - p_c|^{-\nu}$, and $d$ is the dimensionality of the composite, with $\nu = 4/3$ and $0.88$ for $d = 2$ and $3$, respectively. From our analysis of conduction in Bethe lattices with power-law conductors presented above and Eqs. (12), (16) and (17), it should be clear to the reader that for any $d$-dimensional random resistor network with power-law conductors, the exponent $\mu_n$, defined by Eq. (17), which is the analogue of $\mu$, must depend on $n$. This is indeed the case. One can rewrite Eq. (65) in a more general form (Kenkel and Straley, 1982)

$$\mu_n = \mu(n) = (d - 1)\nu + \frac{1}{n}[\zeta(n) - \nu], \qquad (66)$$

indicating explicitly that the $n$-dependence of $\mu$ must be through the resistivity exponent $\zeta$ as $\nu$ is a purely topological property, independent of the transport process. Numerical simulations and scaling analyses discussed below show that this is indeed the case. In fact, extensive analysis of random resistor networks with power-law conductors indicates that, for certain limits and values of $n$, the exponent $\tilde{\zeta}(n) = \zeta(n)/\nu$ is related to various topological properties of the network. We now describe these relations which provide insight into the $n$-dependence of $\zeta(n)$ and hence $\mu_n$.

In general, as Eq. (66) indicates, $\mu_n$ is larger than $\mu$, and therefore near $p_c$ the conductivity curve for power-law transport is flatter than that of the linear transport. Several exact relations between $\zeta(n)$ and the topological exponents of percolation networks have been derived. We present the proof of one of these relations to give the reader some idea about how they are derived. Blumenfeld and Aharony (1985) proved that

$$\tilde{\zeta}(n = \infty) = D_r, \qquad (67)$$

where $D_r = 1/\nu$ is the fractal dimension of the red bonds in the sample-spanning cluster, i.e., those that, if cut, split the cluster into two parts. If $M_r$ is the number of the red bonds, then for length scale $L < \xi_p$ the fractal dimension $D_r$ is defined by, $M_r \sim L^{D_r}$. To prove this relation, consider a two-terminal blob of bonds (a subcluster of multiply-connected conducting bonds) near $p_c$, and suppose that the current through the blob is $I$, while the voltage drop between its two terminals is $V$. Thus, the resistance $R$ of the blob is given by, $R = V^{1/n}/I$. Now, if we select any transport path between the two ends of the blob, we can write, $V = \sum_j R_j i_j^n$, where $R_j$ is the resistance of bond $j$ along the path, and $i_j$ is its current. Therefore,

$$R = \left[ \sum_j R_j \left( \frac{i_j}{I} \right)^n \right]^{1/n}. \qquad (68)$$

However, $i_j < I$, and therefore $(i_j/I)^n$ should vanish as $n \to \infty$, implying that the blob resistance will be zero, and thus all of the resistance of the cluster (material) is offered by the red bonds, hence proving Eq. (67). By similar arguments Blumenfeld and Aharony (1985) also proved that

$$\tilde{\zeta}(n = 0^+) = D_{min}, \qquad (69)$$

where $D_{min}$ is the fractal dimension of the minimum or chemical path between two points of a percolation cluster, i.e., the shortest path between the two points. Thus, for $L < \xi_p$, the minimum length $L_{min}$ scales with $L$ as, $L_{min} \sim L^{D_{min}}$, with $D_{min} \simeq 1.13$ and 1.34 in 2D and 3D, respectively. Moreover, Blumenfeld *et al.* (1986) showed that

$$\tilde{\zeta}(n = 0^-) = D_{max}. \tag{70}$$

Here $D_{max}$ is the fractal dimension associated with the *longest* self-avoiding walk (that is, a random walk in which the walker never visits any point more than once) between the two terminals of the percolation network; if $L_{max}$ is the length of the walk, then $L_{max} \sim L^{D_{max}}$. Blumenfeld *et al.* (1986) also proved that

$$\tilde{\zeta}(n = -1) = D_{bb}, \tag{71}$$

with $D_{bb}$ being the fractal dimension of the backbone of percolation clusters. Note, however, that it has not been possible to relate $\zeta(n = 1)$ to any of the topological exponents. Blumenfeld *et al.* (1986) also proved that $\zeta(n)$ decreases monotonically with $n$, and therefore $d\zeta(n)/dn \leq 0$, with the equality holding at $n = \infty$. Using values of the various exponents and fractal dimensions given in Table 2.3 of Volume I, we see that in 2D, $\zeta(n = \infty) = 1$, and $\zeta(n = -1) \simeq 2.18$, whereas in 3D $\zeta(n = \infty) = 1$, and $\zeta(n = -1) \simeq 1.6$. Therefore, $\zeta(n)$ is a slowly-varying function of $n$.

In addition to direct numerical simulations, there are at least two other methods for estimating $\mu_n$ and its dependence on $n$. These methods are generalizations of those discussed in Chapter 5 for the linear conductivity, and in what follows we describe them briefly.

### 3.1.7.1    Series Expansion Analysis

Meir *et al.* (1986) used a series expansion method to calculate $\zeta(n)$ for several values of $n$. As discussed in Chapter 5 of Volume I for linear conduction, in this method one defines a percolation susceptibility $\chi_p$ by

$$\chi_p = \left\langle \sum_j s_{ij} \right\rangle, \tag{72}$$

where $s_{ij} = 1$ if the two sites $i$ and $j$ belong to the same percolation cluster and $s_{ij} = 0$ otherwise, and the averaging is over all configurations of the occupied sites (probability $p$) and unoccupied ones [probability $(1 - p)$]. We now define a resistive susceptibility $\chi_R(n; \mathcal{C})$ for a cluster $\mathcal{C}$ of sites via

$$\chi_R(n; \mathcal{C}) = \sum_{i \in \mathcal{C}} \sum_{j \in \mathcal{C}} R_{ij}(n), \tag{73}$$

where $R_{ij}(n)$ is the nonlinear resistance between sites $i$ and $j$. Then, the total resistive susceptibility $\chi_R(n)$, defined by

$$\chi_R(n) = \left\langle \sum_j R_{ij}(n) \right\rangle, \tag{74}$$

is obtained by summing $\chi_R(n; \mathcal{C})$ over all cluster, weighting each cluster by its probability of occurrence. This is usually done in terms of cumulants, whereby one writes

$$\chi_R(n) = \sum_{\mathcal{C}} N(\mathcal{C}; d) p^{n_b(\mathcal{C})} \chi_R^c(n; \mathcal{C}). \tag{75}$$

In this equation $n_b(\mathcal{C})$ is the number of bonds in the cluster, $N(\mathcal{C}; d)$ is the number of ways *per site* a diagram, topologically equivalent to $\mathcal{C}$, can be realized on a $d$-dimensional simple-cubic lattice, and the sum is over all topologically *inequivalent* diagrams $\mathcal{C}$. Moreover, $\chi_R^c(n; \mathcal{C})$ is the cumulant defined by

$$\chi_R^c(n; \mathcal{C}) = \chi_R(n; \mathcal{C}) - \sum_{\gamma \in \mathcal{C}} \chi_R^c(n; \gamma), \tag{76}$$

where the sum is over all subdiagrams $\gamma$ of $\mathcal{C}$. Then, the average resistance $\langle R(n) \rangle$ is defined by

$$\langle R(n) \rangle = \frac{\chi_R}{\chi_p} \sim |p - p_c|^{-\zeta(n)}. \tag{77}$$

Therefore, the procedure for series analysis of resistance of random resistor networks with power-law conductors is as follows. For each cluster $\mathcal{C}$, and fixed values of $n$ (the power-law exponent) and $n_b$ (the number of bonds in the cluster), the resistance $R_{ij}(n)$ is computed (by solving the Kirchhoff's equations). These computations are carried out for all such clusters, from which $\chi_R(n; \mathcal{C})$ and hence $\chi_R(n)$ are obtained. Writing

$$\chi_R(n) = \sum_k \sum_l A(k, l) d^l p^k, \tag{78}$$

one obtains a power series in $p$ for $\chi_R(n)$. Since, in practice, the number of possible cluster configurations increases very rapidly with $n_b$, the computed power series cannot be very long. For example, Meir *et al.* (1986) calculated the first 11 terms of the series. Another series is obtained for $\chi_p$, the computation of which is very simple since it involves only counting of the number of clusters' configurations. The resulting two power series are then analyzed by a Padé approximation method, from which the average resistance $\langle R(n) \rangle$ and hence the resistivity exponent $\zeta(n)$ are computed. Using the results of Meir *et al.* (1986) and Eq. (66), we present in Figure 3.4 the variations of $\mu_n = \mu(n)$ with $n$. This figure indicates that $\mu_n$ decreases very rapidly with increasing.

### 3.1.7.2   Field-Theoretic Approach

Harris (1987) developed a field-theoretic approach to power-law transport, a generalization of what we described in Chapter 5 of Volume I for the linear conduction problem, and derived an $\epsilon$-expansion (where $\epsilon = 6 - d$, with $d$ being the dimensionality of the system) for $\zeta(n)$ which, to linear order in $\epsilon$, is given by

$$\zeta(n) = 1 + \frac{\epsilon}{42} \left[ 1 - \frac{7(n-1)}{72} \right] + O(\epsilon^2). \tag{79}$$

FIGURE 3.4. Dependence of the power-law conductivity exponent $\mu(n)$ on the power-law exponent $n$ (after Sahimi, 1993a, plotted based on the results of Meir *et al.*, 1986).

Since (Harris *et al.*, 1975), $\nu = 1/2 + 5\epsilon/84 + O(\epsilon^2)$, we obtain, using Eq. (66),

$$\mu(n) = \frac{5}{2} - \frac{17\epsilon}{84} + \frac{1}{n}\left\{\frac{1}{2} - \frac{\epsilon}{84}\left[3 + \frac{7(n-1)}{36}\right]\right\}, \tag{80}$$

which reduces, in the limit $n = 1$, to Eq. (5.233) of Volume I for linear conductivity. Such $\epsilon$-expansions, while predicting the correct general trends in the $n$- and $d$-dependence of the exponent $\mu(n)$, are not very accurate for the practical cases of $d = 2$ and 3.

### 3.1.8 Resistance Noise, Moments of Current Distribution, and Scaling Properties

As discussed in Section 5.16 of Volume I for linear conduction, in a conducting composite resistance noise manifests itself as voltage fluctuations, when the sample is subjected to constant current bias, or as current fluctuations in content voltage bias. The low-frequency power spectrum of the resistance fluctuations often varies as $1/f$, where $f$ is the frequency. This is the so-called flicker or $1/f$ noise. [In the literature on this subject, frequency is often denoted by $f$, instead of $\omega$, so that the resistance noise is often referred to as $1/f$ noise. Therefore, we depart in this section from our standard notation in this book, and use $f$, instead of $\omega$, to denote the frequency so as not to confuse the reader.] The magnitude of resistance

noise depends on the morphology of the conducting sample. Resistance noise was studied in Chapter 5 of Volume I for the case of linear composites. In this section, we consider the same problem for power-law conductors described by Eq. (1).

Consider a sample composite in which each of the conducting elements of the nonlinear resistors has the same average value, but is fluctuating independently with a correlation $\langle \delta r_a \delta r_b \rangle = \rho^2$, where $r_a$ and $r_b$ are two resistances. Then, similar to what was discussed in Chapter 5 of Volume I, the relative noise $S_R$ is calculated from

$$S_R = \frac{\langle \delta R \delta R \rangle}{R^2} = \frac{\rho^2}{r^2} \frac{\sum_b i_b^{2(n+1)}}{(\sum_b i_b^{n+1})^2}, \tag{81}$$

where $R$ is the resistance of the sample, $i_b$ is the current in the bonds, and the sums are over all the current-carrying bonds. Note that the voltage noise $S_V = \langle \delta V \delta V \rangle / V^2$ itself is given by, $S_V \sim I^{2n}$, and that for a homogeneous, $d$-dimensional lattice of linear size $L$, $S_R = (\rho^2/r^2)/L^d$.

For the sample-spanning percolation cluster at $p_c$ (or, equivalently, at length scales $L < \xi_p$ above $p_c$) the resistance noise scales with $L$ as

$$S_R \sim L^{-b_n}, \tag{82}$$

where $b_n = b(n)$ is the analogue of the exponent $b$ for linear conduction, Eq. (5.250) of Volume I. Rammal and Tremblay (1987) showed that

$$\tilde{\zeta}_n \leq b_n \leq D_{bb}, \quad b_n \leq 2\tilde{\zeta}_n - D_r, \tag{83}$$

where, as before, $\tilde{\zeta}_n = \tilde{\zeta}(n) = \zeta(n)/\nu$, and $D_{bb}$ and $D_r$ are the fractal dimensions of the backbone and the red bonds, respectively. As discussed in Chapter 5 of Volume I, these bounds are also satisfied in the linear conduction case. Near the percolation threshold $p_c$,

$$S_R \sim (p - p_c)^{-\kappa_n}, \tag{84}$$

where, similar to the case of linear conduction, $\kappa_n = \kappa(n)$ is a completely new exponent independent of all the percolation exponents. Of course, $\kappa_n$ and $b_n$ are related, $\kappa_n = \nu(d - b_n)$, and therefore the above bounds for $b(n)$ can be immediately converted to bounds for $\kappa_n$.

While $S_R$ is related to the 4th moment of the current distribution, one can, similar to linear conduction discussed in Chapter 5 of Volume I, construct the general moments $M_q(\mathbf{x}, \mathbf{x}')$ of the current distribution between two points $\mathbf{x}$ and $\mathbf{x}'$;

$$M_q(\mathbf{x}, \mathbf{x}') = \sum_b i_b^{(n+1)q}, \tag{85}$$

where, as before, the sum is over all the current-carrying bonds of the network. Then, for self-similar morphologies, such as the sample-spanning percolation cluster at $p_c$ (or at length scales $L < \xi_p$ above $p_c$), one can define an infinite hierarchy of exponents $\tau_q(n)$ for $|\mathbf{x} - \mathbf{x}'| \sim L$:

$$M_q(\mathbf{x}, \mathbf{x}') \sim L^{-\tilde{\tau}_q(n)}, \tag{86}$$

Similar to the case of linear conduction, the exponents $\tau_q(n)$ are independent of each other. Moreover, Rammal and Tremblay (1987) proved that $\tau_0 - \tau_q(n)$ is a decreasing convex function of $q$ that satisfies the following inequalities,

$$\tau_{q-1}(n) \leq \tau_q(n) \leq \frac{q}{q-1}\tau_{q-1}(n) - \frac{1}{q-1}\tau_0, \tag{87}$$

where the last of the two inequalities is valid only for $q \geq 1$. For the sample-spanning percolation cluster at length scales $L < \xi_p$, one has, $\tilde{\tau}_q(n) = \tau_q(n)/\nu$. Rammal and Tremblay (1987) obtained approximate (but not particularly accurate) estimates of these exponents.

## 3.2  Nonlinear Transport Caused by a Large External Field

Another type of nonlinear transport process arises as a result of applying a large external potential gradient or driving force to a disordered material. Examples are abundant and include flux lines in superconductors (see, for example, Larkin and Ovchinnikov, 1979; Brass *et al.*, 1989; Feigel'man and Vinokur, 1990; Fisher *et al.*, 1991), various fluid flow phenomena in porous materials (for reviews see, for example, Sahimi, 1993b,1995b), and sliding charge-density waves (see, for example, Fisher, 1985; Gorkov and Grüner, 1989). Dielectric breakdown, to be studied in Chapters 5 and 6, also belongs to this class of phenomena.

In general, one must distinguish between two different types of systems in which transport is driven by a large external field. In one type the disorder is weak, and thus the interactions between the transport carriers produce an elastic structure that will be distorted but will not break. Charge-density waves, and invasion of a porous material by a wetting front belong to this class of systems. In the second type, disorder is strong and the elastic medium can break up, giving rise to transport processes that are plastic or fluid-like. An important example is strongly-pinned vortex lines in the mixed state of superconducting films. This type of systems, unlike the first type, has not received the attention that it deserves, despite its practical importance, and is the subject of this section.

When a large potential gradient or driving force is imposed on a material, it induces bias in it in the sense that, in a $d$-dimensional system there will be an "easy" or longitudinal direction which is the direction of the external potential gradient, and along which transport takes place "easier" than the remaining $(d-1)$ transverse directions. This bias also induces *anisotropy* in the material so that one must introduce *two* correlation lengths, instead of one as in isotropic systems, which are the longitudinal correlation length $\xi_L$ and the transverse correlation length $\xi_T$ (see Figure 3.5). It is not unreasonable to assume that there is a critical value of the external potential or force $F_c$ such that for $F \geq F_c$ macroscopic transport occurs. Suppose now that an external driving force $F > F_c$ is imposed on the system. The dimensionless potential, $\chi = (F - F_c)/F_c$, plays the same role as $(p - p_c)$ in percolation. Because $F_c$ represents a type of critical point or

FIGURE 3.5. A strong external potential induces dynamic anisotropy in a material, giving rise to two correlation lengths $\xi_L$ and $\xi_T$. Circle denotes the point at which the potential is applied to the system (after Sahimi, 1993a).

threshold, it is not unreasonable to assume that near $F = F_c$ one must have

$$\xi_L \sim |F - F_c|^{-\nu_L}, \tag{88}$$

$$\xi_T \sim |F - F_c|^{-\nu_T}. \tag{89}$$

The problem studied here has certain similarities with *directed* percolation (Kinzel, 1983; Duarte, 1986,1990,1992; Duarte *et al.*, 1992). In directed percolation, the bonds of a network are directed and diode-like. Transport along such bonds is allowed only in one direction. If the direction of the external potential is reversed, then there can be no macroscopic transport in the new direction. Similar to the present problem, in directed percolation one also needs *two* correlation lengths to characterize the shape of the percolation clusters. However, there is an important difference between what we study here and directed percolation: The anisotropy in our system is *dynamically* induced, whereas the bias and anisotropy in directed percolation are *static* and fixed.

   An example of such nonlinear systems is the model proposed by Narayan and Fisher (1994) (see also the somewhat related model proposed by Herrmann and Sahimi, 1993, and Herrmann *et al.*, 1993). They considered a randomly-rough surface onto which a fluid or a charge carrier is poured into isolated "lakes," such that initially a sample-spanning cluster of connected lakes does not exist. The surface is then slowly tilted at an angle $\theta$, such that the fluid spills out of the filled lakes and feeds unfilled lakes further downhill. For $\theta < \theta_c$, where $\theta_c$ is the critical value of the tilt angle, the filled lakes cluster together. The characteristic size of such clusters increases as $\theta$ does, and diverges at $\theta = \theta_c$. Above $\theta_c$ the system becomes depinned, so that the fluid or the charge carrier can flow from the top to the bottom of the system. Near and above $\theta_c$ the transport process is highly inhomogeneous and confined to narrow and well-separated channels, somewhat

similar to Figure 3.5. Note that, under the influence of gravity, a force builds up at the terminus of a cluster, rather than being uniform everywhere in it. Therefore, when $\theta$ increases, clusters grow from their terminus sites, with a higher probability of growing if they are already large. This implies that, the dominating flow paths cannot be determined by a local analysis that searches for weak links in the system. Rather, one must consider the entire system, i.e., the phenomenon is *non-local*.

The above description is a continuum one, but has a well-defined lattice counterpart. In the lattice model, the sites represent the lakes, while the bonds are the transport paths that connect the lakes. A force $F$ is imposed on the lattice, and it suffices for each site $i$ to have outlets connecting it only to its $d$ nearest neighbors $i_\alpha$ in the next plane downhill, where $d$ is the dimensionality of the system. It is assumed that the current flowing in a path depends only on the depth above the lip of the lake it emerges from. Thus, a barrier $b_{i\alpha}$ is assigned to each outlet $\alpha$ emerging from a site $i$ which controls the current flowing through the outlet. The barriers are selected randonmly and independently from a distribution. At each site $i$ of the lattice there is a depth of fluid $h_i$. The current $I_{i\alpha}$ flowing through an outlet $\alpha$ from a site $i$ is zero if $h_i < b_{i\alpha} - F$, and

$$I_{i\alpha} = (h_i - b_{i\alpha} + F)^\varpi \quad \text{if} \quad h_i > b_{i\alpha} - F. \tag{90}$$

The exponent $\varpi$ characterizes the transport over the barrier lip. Narayan and Fisher (1994) presented arguments that indicate that $\varpi = 3 + d/2$ for a $d$-dimensional system. Note that an increase in $F$ is equivalent to uniformly lowering all the barriers $b_{i\alpha}$.

Narayan and Fisher (1994) argued that $\xi_T \sim \sqrt{\xi_L}$. That is, we can imagine that the consecutive events in which the bonds are filled with the flowing current are in fact consecutive steps of a random walk in the $(d-1)$ transverse directions. If so, the longitudinal direction acts as the time axis, and therefore the distance that the random walker travels in the transverse direction should increase with the square root of time (the usual law of random walks), implying that $\xi_T \sim \sqrt{\xi_L}$, and thus $v_T = v_L/2$. The random-walk argument can also be used to estimate the upper critical dimension $d_u$ of the system at and above which the mean-field theory is exact. The clusters perform random walks in the $(d-1)$-dimensional transverse space, with the longitudinal direction acting as the time coordinate. From the theory of random walks (Hughes, 1995) we know that if $d - 1 > 2$, then two walks that start out close to each other have a finite probability of not crossing each other, whereas for $d - 1 < 2$ they are certain to cross. Therefore $d_u - 1 = 2$ and hence $d_u = 3$. This immediately implies another significant difference between this model and directed percolation for which $d_u = 5$ (Obukhov, 1980), and also with isotropic percolation for which $d_u = 6$.

Narayan and Fisher (1994) studied various topological and transport properties of this model. One surprising aspect of this phenomenon is that the critical exponents that characterize the power-law behavior of the properties of interest above and below, but near, the threshold $F_c$ are *not* equal. Consider first the system below the threshold. We write $\xi_L \sim |F - F_c|^{-\nu_L^b}$, where superscript $b$ signifies the fact that the critical exponent is associated with the regime below the threshold. The

fraction of the sites $P_b(F)$ which are in clusters of length $\sim \xi_L$ scales as

$$P_b(F) \sim \xi_L^{-\tilde{\beta}_{nn}} \sim |F - F_c|^{\beta_{nn}}, \tag{91}$$

where $\tilde{\beta}_{nn} = \beta_{nn}/v_L^b$. The mean distance $\langle \ell_p \rangle$ travelled by a charge carrier from its initial position at $F = 0$ (also called the polarization density) scales as

$$\langle \ell_p \rangle \sim |F - F_c|^{1-\gamma_{nn}}. \tag{92}$$

The clusters of the sites (lakes) are fractal objects at length scales $L \ll \xi_L$ with a fractal dimension $D_f$. The two exponents $\beta_{nn}$ and $\gamma_{nn}$ are related through the following scaling law,

$$\gamma_{nn} = v_L(1 - \tilde{\beta}_{nn}), \tag{93}$$

both above and below the threshold $F_c$. One can show that, in the mean-field approximation, i.e., at $d = 3$, one has

$$v_L^b = \frac{3}{2}, \quad D_f = \frac{4}{3}, \quad \tilde{\beta}_{nn} = \frac{2}{3}, \quad \gamma_{nn} = \frac{1}{2}. \tag{94}$$

Consider next the regime above the threshold. An important property is the fraction $P_a(F)$ of sites that feed charge carriers into the transport paths, i.e., the analogue of $P_b(F)$ above the threshold. Near $F_c$,

$$P_a(F) \sim (F - F_c)^{\Gamma}, \tag{95}$$

and it is clear that, $\Gamma = \tilde{\beta}_{nn} v_L^a$, where superscript $a$ indicates that the critical exponent is associated with the regime above the threshold. In general, one has the following scaling laws (Narayan and Fisher, 1994)

$$v_L^a = \frac{1 + \Gamma}{d - 1}, \tag{96}$$

$$D_f = \frac{1}{2}(d + 1) - \frac{\beta_{nn}}{v_L^a}, \tag{97}$$

Near $F_c$ the mean current density $\langle I \rangle$ flowing through the system obeys the following power law

$$\langle I \rangle \sim (F - F_c)^{\mu_{nn}}. \tag{98}$$

The transport exponent $\mu_{nn}$ is then given by

$$\mu_{nn} = \frac{1}{2}(1 + \varpi)(1 + \Gamma). \tag{99}$$

Scaling law (99) is an interesting feature of this model for two reasons. First, it implies that, unlike percolation, in this model the transport exponent is related to the topological exponent $\Gamma$. Secondly, it indicates a sort of non-universality, since $\varpi$ is a local or microscopic quantity. In the mean-field approximation

$$v_L^a = \frac{3}{4}, \quad \Gamma = \frac{1}{2}, \quad \mu_{nn} = \frac{3}{4}(1 + \varpi). \tag{100}$$

Note that $v_L^b \neq v_L^a$. In 1D the problem can be solved exactly and one obtains, $\tilde{\beta}_{nn} = 0$, $v_L^b = 2$, and $\gamma_{nn} = 2$ (note that in 1D only the regime below the threshold

is physically meaningful). Since the upper critical dimension is $d_u = 3$, $d = 2$ is the only physical dimension for which exact results are not known. Numerical simulations of Narayan and Fisher (1994) yielded the following estimates

$$\nu_L^b \simeq 1.76, \quad \nu_L^a \simeq 1.41, \quad \tilde{\beta}_{nn} \simeq 0.29, \quad \Gamma \simeq 0.41, \quad D_f \simeq 1.21. \qquad (101)$$

Note the significant difference between $\nu_L^b$ and $\nu_L^a$. Note also that, similar to conventional percolation, all the exponents can be estimated from any two exponents, e.g., $\nu_L^a$ (or $\nu_L^b$ below the threshold) and $\Gamma$. The low value of $D_f$ implies that, a large external field and the associated dynamical bias and anisotropy give rise to transporting paths that are essentially restricted to a narrow cone (see Figure 3.5). Moreover, the fractal dimensions $D_f$ is considerably smaller than that of 2D percolation clusters, $D_f = 91/48 \simeq 1.896$. This can be understood if we consider the problem on the Bethe lattice, i.e., the mean-field limit. In this lattice any large external potential makes the network completely directed, since there are no closed loops in the lattice. As a result, the backbone is made of directed branches that have a quasi-1D structure, and thus the fractal dimension of the backbone is, $D_{bb} = 1$ (for percolation $D_{bb} = 2$), implying that only a small subset of all the bonds participate in the transport process.

## 3.3   Weakly Nonlinear Composites

We now consider a more general composite in which a material with nonlinear $I - V$ characteristics is embedded randomly in a host with either linear or nonlinear $I - V$ response. To our knowledge, the suggestion for theoretical consideration of such composites was first made by Fleming and Grimes (1979) and Mantese *et al.* (1981) (see also Yagil *et al.*, 1994, for an interesting experimental study of this problem). A concrete step toward this goal was taken by Gefen *et al.* (1986) who proposed and studied the following problem. Consider a random resistor network near the percolation threshold $p_c$, which is driven by an external current $I$. If $I$ is sufficiently weak, then the response of the system is linear, and its linear conductivity $g^{(\ell)}$ follows a power law similar to Eq. (18). If the external current $I$ is gradually increased, then for some critical current $I_c$ the conductivity of the system deviates significantly from its linear value $g_e^{(\ell)}$. Gefen *et al.* (1986) suggested that if $L$, the linear size of the sample, is greater than $\xi_p$, the percolation correlation length, then

$$I_c \sim \left[ g_e^{(\ell)} \right]^x, \qquad (102)$$

and that $x = 3/2$ in 2D. To confirm this prediction, they measured the electrical conductivity of thin gold films near $p_c$ and found that $x \simeq 1.47$, in good agreement with their prediction. If, however, $L \ll \xi_p$, then $I_c$ would depend on $L$ and Gefen *et al.* (1986) proposed that

$$I_c(L) \sim \left[ g_e^{(\ell)}(L) \right]^{-y}. \qquad (103)$$

Both $x$ and $y$ are supposed to be universal.

To explain their theoretical predictions and experimental measurements, Gefen *et al.* (1986) considered a percolation resistor network in which each conducting bond satisfied the following relation between the current $i$ flowing through it and the voltage $v$:

$$v = r_\ell i - r_n i^n, \tag{104}$$

where $r_\ell$ and $r_n$ are, respectively, the linear and nonlinear resistances, and $n > 1$. Note that, in materials with inversion symmetry, the lowest value of $n$ is 3. For small enough $i$, the second term of the right-hand side of Eq. (104) is much smaller than the first term, and therefore the resistor behaves linearly. For sufficiently large $i$ the second term becomes important, and the resistor is nonlinear. The critical current $i_c$ at which the crossover occurs is found by equating the two terms of the right-hand side of Eq. (104), resulting in

$$i_c = \left(\frac{r_\ell}{r_n}\right)^{1/(n-1)}. \tag{105}$$

Composites that are described by Eq. (104), or by similar equations (see below), are what we refer to as *weakly nonlinear materials*, since the leading order term is still linear.

Let us now discuss important properties of nonlinear composites modeled as a system of nonlinear elements with an $I - V$ characteristic that is described by Eq. (104) or by a similar equation. We do not discuss numerical simulations of such phenomena which, although somewhat difficult, is conceptually straightforward and requires no particular explanation

### 3.3.1   Effective-Medium Approximation

As the first problem in this class of composites, we describe the development of an effective-medium approximation (EMA) for predicting the macroscopic behavior of the composite. As usual, we use the terminology of a resistor network, although all the discussions presented here are also applicable to continuum models (with spherical inclusions). Consider a resistor network in which a fraction $1 - p$ of the bonds are linear conductors with an $I - V$ characteristic given by, $i = g_A v$, where $g_A$ is the conductance. The rest of the bonds, with a fraction $p$, are weakly nonlinear conductors with a current-voltage characteristic given by

$$i = g_B v + g^{(n)} v^3, \tag{106}$$

which is another version of Eq. (104), written explicitly for the current $i$ (rather than the voltage $v$). We assume that $g^{(n)} v^2 / g_B \ll 1$. To derive an EMA for this problem (Stroud and Hui, 1988; Zeng *et al.*, 1988; Zeng, Hui, Bergman and Stroud, 1989; Hui, 1990a; Yang and Hui, 1991) we replace the resistor network by a uniform effective network of identical conductors with a current-voltage characteristic given by

$$I = g_e^{(\ell)} v + g_e^{(n)} v^3, \tag{107}$$

where $g_e^{(\ell)}$ and $g_e^{(n)}$ are the effective linear and nonlinear response of the network, respectively. In general, as our discussion throughout this book should have made it clear, the effective linear conductivity $g_e^{(\ell)}$ in a binary random network with components $g_A$ and $g_B$ can always be written as

$$g_e^{(\ell)} = F(g_A, g_B, p), \tag{108}$$

where $F$ is a function which, in general, depends on the geometry of the system. Then the effective nonlinear response $g_e^{(n)}$ of the system is given by

$$g_e^{(n)} = \frac{g^{(n)}}{p} \left[ \frac{\partial F}{\partial g_e^{(\ell)}} \right]^2. \tag{109}$$

That is, the effective nonlinear response is estimated based on an estimate of the effective conductivity of the same material but in the linear regime. Recall from Chapter 2 that the same sort of idea was developed by Ponte Castañeda (1992b) in the context of the continuum models. The derivation of Eq. (109) will be discussed in detail in Section 3.4, where we describe the derivation of a similar equation for the dielectric constant of the same type of composites. Therefore, if the function $F$ can somehow be calculated, $g_e^{(n)}$ will also be determined from Eq. (109). Since $F$ is an estimate of the effective conductivity of a linear binary composite, we may use the EMA, Eq. (24) (or, for example, the Maxwell–Garnett or any other approximation), for linear resistor networks which for our binary network with $f(g) = p\delta(g - g_B) + (1 - p)\delta(g - g_A)$ is given by

$$(1 - p)\frac{g_A - g_e^{(\ell)}}{g_A + g_e^{(\ell)}(Z/2 - 1)} + p\frac{g_B - g_e^{(\ell)}}{g_B + g_e^{(\ell)}(Z/2 - 1)} = 0. \tag{110}$$

Thus, the procedure for calculating $g_e^{(n)}$ by an EMA is as follows. One first solves Eq. (110) for $g_e^{(\ell)}$. This equation, which is quadratic in $g_e^{(\ell)}$, defines the function $F$. Having determined $g_e^{(\ell)}$, one utilizes Eq. (109) to calculate $g_e^{(n)}$. Figures 3.6 and 3.7 compare the results of computer simulations in the square network in two limiting cases with the EMA predictions. The numerical results in Figure 3.6, which are for $g_A = 10$, $g_B = 20$, and $g^{(n)} = 0.1$, are in excellent agreement with the EMA predictions. The reason for the agreement is that the difference $g_B - g_A$ is not large and thus, as discussed in Chapter 5 of Volume I, the function $F$ (i.e., the EMA estimate) provides accurate predictions for $g_e^{(\ell)}$. On the other hand, the numerical results shown in Figure 3.7, which are for $g_A = 5000$, $g_B = 10$, and $g^{(n)} = 0.1$, agree only qualitatively with the EMA predictions because, as discussed in Chapter 5 of Volume I, in this case, due to the large difference between $g_A$ and $g_B$, $F$ (i.e., the EMA estimate) cannot provide accurate predictions for $g_e^{(\ell)}$, which is consistent with the general properties of the EMA.

It is clear that the development of an EMA for this class of composites involves two stages. More generally, one may consider composites with more complex $I - V$ characteristics and develop a similar, but multistage, procedure for an EMA-based computation of their effective transport properties. For example, one may

FIGURE 3.6. The effective nonlinear conductivity $g_e^{(n)}$, normalized by the effective conductivity of the system in the linear regime, versus the fraction $p$ of nonlinear conductors in the square network. Solid curve shows the EMA predictions, while the symbols show the results of numerical simulations for $g_B = 2g_A = 20$ and $g^{(n)} = 0.1$ (after Yang and Hui, 1991).



FIGURE 3.7. Same as in Figure 3.6, but for $g_A = 5000$, $g_B = 10$, and $g^{(n)} = 0.1$ (after Yang and Hui, 1991).

consider a composite material (Yu and Gu, 1993) in which a fraction $p$ of the system has an $I - V$ characteristic given by, $i = g_B v + g_{n1} v^3 + g_{n2} v^5$, while the rest of the composite, with fraction $1 - p$, is made of linear conductors, $i = g_A v$. One may compute the effective linear and nonlinear response of such a composite defined by, $I = g_e^{(\ell)} v + g_e^{(n1)} v^3 + g_e^{(n2)} v^5$, by first solving the EMA equation for the effective linear conductivity of the composite $g_e^{(\ell)}$. Then, an equation similar to (108) is used for computing the first nonlinear conductivity $g_e^{(n1)}$. The two conductivities $g_e^{(\ell)}$ and $g_e^{(n1)}$ so obtained are then used in a higher-order equation in order to compute $g_e^{(n2)}$.

## 3.3.2 Resistance Noise, Moments of Current Distribution, and Scaling Properties

To explain the experimental data of Gefen *et al.* (1986) (see above) for their weakly nonlinear conducting materials, Aharony (1987) established a relation between Gefen *et al.*'s problem and the distribution of currents in a *linear* random resistor network. Consider first the regime $L \ll \xi_p$, which is equivalent to $p = p_c$. The total dissipated power $\mathcal{P}$ in the network, the bonds of which have an $I - V$ characteristic given by Eq. (104), is

$$\mathcal{P} = \frac{1}{2} \sum_b r_\ell |i_b|^2 - \frac{1}{n+1} \sum_b r_n |i_b|^{n+1}, \qquad (111)$$

where $i_b$ is the current in bond $b$, which depends implicitly on $n$, and the sums are over all the conducting bonds of the network. Blumenfeld *et al.* (1986) had already proved that

$$\left. \frac{\partial \mathcal{P}}{\partial n} \right|_{r_n=0} = \frac{1}{n+1} \sum_b |i_b^0|^{n+1}, \qquad (112)$$

where $i_b^0 = i_b(r_n = 0)$. Therefore, to linear order in $r_n$, we can replace $i_b$ by $i_b^0$ and write

$$\mathcal{P} = \frac{1}{2} r_\ell M_1 I^2 - \frac{r_n}{n+1} M_{(n+1)/2} I^{n+1}, \qquad (113)$$

where $I$ is the total current in the network, and

$$M_q = \sum_b \left( \frac{i_b^0}{I} \right)^{2q}, \qquad (114)$$

is the $2q$th moment of the current distribution in the *linear* random resistor network. As already discussed in Section 3.1.8 [see Eq. (86)] for the case of strongly nonlinear composites, and in Section 5.16 of Volume I for linear systems, for $L \ll \xi_p$ the moments of the current distribution scale with $L$ as

$$M_q \sim L^{-\tilde{\tau}_q}, \qquad (115)$$

where all the $\tilde{\tau}_q$s are distinct. This means that the current distribution in a linear random resistor network is multifractal, i.e., each of its moments scales with $L$ with a distinct exponent, which is similar to the moments of the force distribution in elastic and superelastic percolation networks described in Chapter 8 of Volume I (see Stanley and Meakin, 1988, for a review of general properties multifractal systems and distributions). Therefore, the effective linear resistance $R_e^{(\ell)}$ of the network, which is obtained via $R_e^{(\ell)} = \partial^2 \mathcal{P} / \partial I^2$, shows deviations from a constant value for $n > 1$ and

$$I > I_c(L) \sim i_c \left[ \frac{M_1}{M_{(n+1)/2}} \right]^{1/(n-1)} \sim i_c L^{y \tilde{\tau}_1} \sim [g_e^{(\ell)}(L)]^{-y}, \qquad (116)$$

and therefore (Aharony, 1987)

$$y(n) = \frac{1 - \dfrac{\tilde{\tau}_{(n+1)/2}}{\tilde{\tau}_1}}{n - 1}. \qquad (117)$$

Since $\tilde{\tau}_q$ is a monotonic and convex function (see, for example, Blumenfeld *et al.*, 1986), so also is $y(n)$. For example, for $d = 2$ and 3 one has $y(3) \simeq 0.08$ and 0.06, and $y(0) \simeq 0.18$ and 0.1, respectively. This means that $0 < y(n) < y(1)$, and therefore the linear regime $I < I_c(L)$ extends to larger currents for larger linear sizes $L$, implying that even a narrow nonlinear regime will be *enhanced* (see also below) in a percolation network. A similar analysis for $L \gg \xi_p$ yields (Aharony, 1987)

$$x(n) = \frac{d - 1 - y(n)\tilde{\tau}_1}{d - 2 + \tilde{\tau}_1}, \qquad (118)$$

and therefore for $d = 2$ one finds that $x(n) = 1.03 - y(n)$. Since $y(n) > 0$, Eq. (118) does not agree with the experimental result of Gefen *et al.* (1986) for any $n$, and therefore a simple percolation network in which each conducting bond follows Eq. (104) cannot explain Gefen *et al.*'s data.

To study scaling properties of weakly nonlinear composites near the percolation threshold $p_c$, we must consider resistance and conductance fluctuations in *linear* resistors networks. Recall from Section 5.16 of Volume I that, for a percolation network near $p_c$, the relative linear resistance noise, $S_R = \langle \delta R \delta R \rangle / [R_e^{(\ell)}]^2$, follows the following power law [see also Eq. (84) for strongly nonlinear composites]

$$S_R \sim (p - p_c)^{-\kappa}, \qquad (119)$$

which defines the critical exponent $\kappa$. One can, in a similar fashion, consider conductance fluctuations $S_G$ of a linear superconducting percolation network below $p_c$. In this case

$$S_G \sim (p_c - p)^{-\kappa'}. \qquad (120)$$

It can be shown (Wright *et al.*, 1986) that in 2D, $\kappa = \kappa'$. Given Eqs. (119) and (120), we can discuss some of the scaling properties of weakly nonlinear composites near $p_c$.

Stroud and Hui (1988) considered a composite with the following characteristic,

$$\mathbf{I}(\mathbf{x}) = g^{(\ell)}(\mathbf{x})\mathbf{E}(\mathbf{x}) + g^{(n)}(\mathbf{x})|\mathbf{E}(\mathbf{x})|^n\mathbf{E}(\mathbf{x}), \tag{121}$$

where $n \geq 1$, and $g^{(\ell)}$ and $g^{(n)}$ are the linear and nonlinear conductivities of the medium, respectively, which depend, in general, on the spatial position $\mathbf{x}$, and the applied electric field (or voltage) $\mathbf{E}$. Equation (121) is just another version of (104), written explicitly for the current. As mentioned earlier, if one assumes that all the components in the disordered composite have inversion symmetry, then $n = 2$, which was the case studied by Stroud and Hui (1988). The volume-averaged current $\langle \mathbf{I} \rangle$ is defined by

$$\langle \mathbf{I} \rangle = g_e^{(\ell)}\mathbf{E}_0 + g_e^{(n)}|\mathbf{E}_0|^2\mathbf{E}_0, \tag{122}$$

with $\langle \mathbf{E} \rangle = \mathbf{E}_0$. Consider now the dissipated power for this composite which, in a continuum formulation, is given (for $n = 2$) by

$$\mathcal{P} = \int \mathbf{I} \cdot \mathbf{E}\, d\Omega = \Omega \left[ g_e^{(\ell)}|\mathbf{E}_0|^2 + g_e^{(n)}|\mathbf{E}_0|^4 \right]. \tag{123}$$

This equation, in which $\Omega$ is the volume of the composite, is the continuum analog of Eq. (111). Using Eq. (121), we rewrite Eq. (123) as

$$\mathcal{P} = \int \left[ g^{(\ell)}(\mathbf{x})\mathbf{E} \cdot \mathbf{E} + g^{(n)}(\mathbf{x})|\mathbf{E}|^4 \right] d\Omega = \mathcal{P}_2 + \mathcal{P}_4. \tag{124}$$

Then, to first order in $g^{(n)}(\mathbf{x})$, the second term of Eq. (124) is rewritten as,

$$\mathcal{P}_4 = \Omega\langle g^{(n)}(\mathbf{x})|\mathbf{E}|^4 \rangle_\ell = \langle \mathcal{P}_4 \rangle_\ell, \tag{125}$$

where the subscript $\ell$ indicates that the electric field must be calculated from the solution of the *linear* problem, i.e., in the limit, $g^{(n)}(\mathbf{x}) = 0$. In reality, the difference $\mathbf{E} - \mathbf{E}_\ell$ is of first order in $g^{(n)}$, and therefore will contribute to $\mathcal{P}_4$ only a second-order term. By a similar argument, one can show that

$$\mathcal{P}_2 = \langle \mathcal{P}_2 \rangle_\ell. \tag{126}$$

Therefore, to first order in $g^{(n)}(\mathbf{x})$, the effective conductivities $g_e^{(\ell)}$ and $g_e^{(n)}$ are given by (Stroud and Hui, 1988)

$$g_e^{(\ell)} = \frac{1}{\Omega|\mathbf{E}_0|^2} \int g^{(\ell)}(\mathbf{x})|\mathbf{E}_\ell|^2 d\Omega = \frac{\langle g^{(\ell)}|\mathbf{E}_\ell|^2 \rangle}{|\mathbf{E}_0|^2}, \tag{127}$$

$$g_e^{(n)} = \frac{1}{\Omega|\mathbf{E}_0|^4} \int g^{(n)}(\mathbf{x})|\mathbf{E}_\ell|^4 d\Omega = \frac{\langle g^{(n)}|\mathbf{E}_\ell|^4 \rangle}{|\mathbf{E}_0|^4}. \tag{128}$$

Observe that Eq. (128) is the same as (50) for strongly nonlinear composites. Equations (127) and (128) are manifestations of an important result: The effective linear and nonlinear conductivities of a weakly nonlinear composite can be calculated from the behavior of the electric field in the *linear* problem.

Utilizing a similar line of analysis, Stroud and Hui (1988) proved another important property of weakly nonlinear composites, namely, that to first order in

$g^{(n)}(\mathbf{x})$, $g_e^{(n)}$ is essentially given by the mean square conductivity fluctuations in a *linear* composite,

$$g_e^{(n)} = \frac{\Omega[\Delta g^{(\ell)}]^2}{c},\tag{129}$$

where $\Delta g^{(\ell)}$ is the root mean square conductivity fluctuations in the linear composite, and $c$ is a constant with dimensions of energy. Note that, since the conductivity fluctuations cause corresponding fluctuations in the current, which in turn are related to the 4th moment of the current distribution (see above), Eq. (129) is consistent with, but much more general than, Aharony's result, Eqs. (113)–(118), discussed above.

Using Eq. (129), one can now deduce the power-law behavior of the nonlinear conductivity $g_e^{(n)}$ near the percolation threshold $p_c$. According to Eq. (129), $g_e^{(n)}$ is given by conductivity, or resistivity, fluctuations of the linear conductivity problem. Therefore (Stroud and Hui, 1988), using Eq. (119), we can write

$$\frac{g_e^{(n)}}{[g_e^{(\ell)}]^2} \sim (p - p_c)^{-\kappa},\tag{130}$$

which, when combined with the power-law behavior of the effective linear conductivity $g_e^{(\ell)}$ near $p_c$, Eq. (18), yields

$$g_e^{(n)} \sim (p - p_c)^{2\mu-\kappa},\tag{131}$$

where $\mu$ is the critical exponent of the effective linear conductivity near $p_c$. Note that in a composite in which a fraction $p$ of the material is superconducting and the rest is made of weakly nonlinear conducting material, one has

$$g_e^{(n)} \sim (p_c - p)^{-2s-\kappa'}.\tag{132}$$

With the help of Eqs. (131) and (132), one can construct a general scaling representation for the effective conductivity of a composite, a fraction $p_M$ of which is a good weakly nonlinear conductor characterized by, $I = g_M^{(\ell)}V + g_M^{(n)}V^3$, while the rest of the composite, with a fraction $(1 - p_M)$, is a poor weakly nonlinear conductor which follows, $I = g_I^{(\ell)}V + g_I^{(n)}V^3$, with $g_M^{(\ell)} \gg g_I^{(\ell)}$ and $g_M^{(n)} \gg g_I^{(n)}$. Then, with $z = [g_I^{(\ell)}/g_M^{(n)}]/(p - p_c)^{\mu+s}$, $\Delta p = |p - p_c|$, and considering Eqs. (131) and (132), one can write (Levy and Bergman, 1994b)

$$g_e^{(n)} \simeq g_I^{(n)}\Delta p^{-2s-\kappa'}\Phi_I(z) + g_M^{(n)}\Delta p^{2\mu-\kappa}\Phi_M(z).\tag{133}$$

The properties of the two scaling functions $\Phi_I$ and $\Phi_M$ vary in three distinct regimes.

(1) In regime I, which is for $p_M > p_c$ and $|z| \ll 1$, the scaling function $\Phi_M$ must be constant in order for one to be able to obtain Eqs. (130) and (131). It is then straightforward to see that $\Phi_I$ must also be a constant.
(2) In regime II, which is for $p_M < p_c$ and $|z| \ll 1$, the scaling function $\Phi_I$ must be constant so that one can recover Eq. (132). It is not difficult to see

that in this case, $\Phi_M \sim z^4$. The morphology of the composite consists of a nearly insulating matrix (dominated by the $I$ phase) that contains conducting inclusions (made of the $M$ phase).

(3) In regime III, which is for $p_M \simeq p_c$ and $|z| \gg 1$, the scaling functions $\Phi_I$ and $\Phi_M$ must be such that the dependence of $g_e^{(n)}$ on $p$ is cancelled.

In regime I, the contribution of the good conductor to $g_e^{(n)}$ decreases as $p_M \to p_c^+$ (since $2\mu - \kappa > 0$), whereas the poor conductor's contribution increases. Therefore, if the contribution of the poor conductor happens to be dominant, we will have a *non-monotonic* dependence of $g_e^{(n)}$ upon $p_M$, with a maximum very close to $p_c$, in regime III, and a minimum somewhere above it, in regime I. On the other hand, in regime II ($p_M < p_c$), the contributions from both components increase as $p_M \to p_c^-$. Therefore, one cannot in general determine which component makes the dominant contributions to $g_e^{(n)}$ without specifying $g_I^{(\ell)}/g_M^{(\ell)}$, $g_I^{(n)}/g_M^{(n)}$ and $\Delta p$.

Based on such considerations, then, one can write

$$g_e^{(n)} \simeq \begin{cases} g_M^{(n)}\Delta p^{2\mu-\kappa} + g_I^{(n)}\Delta p^{-2s-\kappa'}, & \text{regime I,} \\ g_M^{(n)}[g_I^{(\ell)}/g_M^{(\ell)}]^4\Delta p^{-2\mu-4s-\kappa} + g_I^{(n)}\Delta p^{-2s-\kappa'}, & \text{regime II,} \\ g_M^{(n)}[g_I^{(\ell)}/g_M^{(\ell)}]^{(2\mu-\kappa)/(\mu+s)} + g_I^{(n)}[g_I^{(\ell)}/g_M^{(\ell)}]^{-(2s+\kappa')/(\mu+s)}, & \text{regime III.} \end{cases}$$

(134)

These scaling function representations are very similar, in their general form, to those for low-field Hall conductivity described in Section 5.17 of Volume I. Numerical simulations of Levy and Bergman (1994b) confirmed the validity of these scaling laws.

### 3.3.3  Crossover from Linear to Weakly Nonlinear Conductivity

Equations (122), (130) and (131) enable us to derive the critical current for the crossover from linear to weakly nonlinear regime. As discussed above, where we derived Eq. (105), the critical voltage $V_c$ or electric field $E_c$ is obtained by equating the two terms of the right hand side of Eq. (122). This yields

$$V_c \sim \left[\frac{g_e^{(\ell)}}{g_e^{(n)}}\right]^{1/2}, \tag{135}$$

from which the critical current $I_c$ is obtained (Blumenfeld and Bergman, 1991a):

$$I_c \sim \left[g_e^{(\ell)}\right]^{(1+\kappa/\mu)/2}. \tag{136}$$

We may interpret Eq. (136) as meaning that, the exponent $x$ defined by Eq. (102), is given by, $x = \frac{1}{2}(1 + \kappa/\mu)$. In 2D, where $\kappa \simeq 1.12$, Eq. (136) predicts that, $x \simeq 0.93$, which still does not agree with Gefen *et al.*'s measurement, $x \simeq 1.47$, but is closer to it than the prediction of Eq. (118).

More generally, let us consider a weakly nonlinear composite with percolation disorder. Specifically, we consider two limiting cases.

(1) A composite in which a (volume) fraction $p$ of the material is made of weakly nonlinear conductors that follow Eq. (121), while the rest of the composite, with a fraction $(1 - p)$, is insulating. Then, near the percolation threshold $p_c$, the critical current $I_c$ and voltage $V_c$ (or, equivalently, the critical electric field $E_c$) follow the following power laws,

$$I_c \sim (p - p_c)^w, \tag{137}$$

$$V_c \sim (p - p_c)^v. \tag{138}$$

(2) We also consider a composite a fraction $p$ of which is made of superconducting materials, while the rest of the system, with a fraction $(1 - p)$, is made of a weakly nonlinear conducting material with an $I - V$ (or $I - E$) characteristic that is given by Eq. (121). Then, we define the critical exponents $w'$ and $v'$ by

$$I_c \sim (p_c - p)^{w'}, \tag{139}$$

$$V_c \sim (p_c - p)^{v'}. \tag{140}$$

For the first limiting case, we use Eqs. (132) and (136) and the appropriate scaling laws for $g_e^{(\ell)}$ and $g_e^{(n)}$ to obtain

$$v = \frac{1}{2}(\kappa - \mu), \quad w = \frac{1}{2}(\kappa + \mu). \tag{141}$$

Since $\kappa + \mu > 0$ while $\kappa - \mu < 0$, Eq. (141) implies that, as $p_c$ is approached, the nonlinear effect is enhanced, so that very close to $p_c$, even a very small $I_c$ would be enough for a crossover from linear to weakly nonlinear conductivity behavior. For the second limiting case (Yu and Hui; 1994; see also Hui, 1990b, 1994) one has

$$v' = \frac{1}{2}(\kappa' + s), \quad w' = \frac{1}{2}(\kappa' - s), \tag{142}$$

where $s$ is the critical exponent that characterizes the power-law behavior of the effective linear conductivity of conductor-superconductor percolation composites near $p_c$, $g_e \sim (p_c - p)^{-s}$, and $\kappa'$ is defined by Eq. (120). Using the numerical estimates of the exponent $s \simeq 1.3$ and 0.73, and $\kappa' \simeq 1$ and 0.4 for $d = 2$ and 3, respectively, we find again that $I_c$ vanishes as $p \to p_c^-$, so that the nonlinear effect is enhanced.

More generally, if one replaces the insulating material with a linear material with conductivity $g_0$ (the first limiting case described above), and let $h = g_0/g^{(\ell)}$, then one has a general scaling equation for $I_c$ (Yu and Hui, 1994):

$$I_c = (p - p_c)^{(\kappa + \mu)/2} \Phi_I[h(p - p_c)^{-(s + \mu)}], \tag{143}$$

which is completely similar to Eqs. (61) and (62). The universal scaling function $\Phi_I(z)$ has the properties that, $\Phi_I(z) \to$ constant as $z \to 0$, while it behaves for

large $z$ as a power law in $z$. For length scales $L \ll \xi_p$ (which is equivalent to $p = p_c$), where $\xi_p$ is the percolation correlation length, one can write

$$I_c = h^{(\kappa+\mu)/2(s+\mu)} \Phi_I' [hL^{(s+\mu)/\nu}], \qquad (144)$$

where $\Phi_I'(z)$ is another universal scaling function such that $\Phi_I' \to$ constant as $z \to \infty$, while $\Phi_I'$ has a power-law dependence on $z$ for $z \to 0$. A similar scaling function representation can also be derived for $V_c$. Hence

$$V_c = (p_c - p)^{(\kappa'+s)/2} \Phi_V [h(p_c - p)^{-(\mu+s)}]. \qquad (145)$$

We note here that one may use the EMA to not only obtain estimates of the exponents $v$, $w$, $v'$ and $w'$, but also explicit expressions for the scaling functions $\Phi_I$, $\Phi_I'$, and $\Phi_V$. All one must do is using Eq. (24) to estimate $g_e^{(\ell)}$ and Eq. (109) to compute $g_e^{(n)}$. Then, it can easily be shown that, $w = v' = 1/2$.

The foregoing scaling laws are valid when one has cubic nonlinearity, i.e., when $n = 2$ in Eq. (121). Zhang (1996a) and Gao *et al.* (1999) generalized these results to any $n$. For the first limiting case, i.e., a composite of insulating and weakly nonlinear conducting materials near $p_c$, Gao *et al.* (1999) obtained the following estimates,

$$v = \frac{\nu d - \zeta - \mu}{2} + \frac{1 + (\nu D_{bb} - 1)^{-n/2} (\zeta - 1)^{n/2+1} - \zeta}{n}, \qquad (146)$$

$$w = \frac{\nu d - \zeta + \mu}{2} + \frac{1 + (\nu D_{bb} - 1)^{-n/2} (\zeta - 1)^{n/2+1} - \zeta}{n}, \qquad (147)$$

where $D_{bb}$ is the fractal dimension of the backbone of the percolation cluster, and $\zeta$ is the resistivity exponent defined by Eq. (65). For the case of a composite of superconducting and weakly nonlinear conducting materials, Zhang (1996a) obtained the following estimate,

$$v' = \frac{s}{2} + \frac{1}{2} \frac{2-n}{n} \nu d + \frac{\kappa'[(n+2)/2]}{n}, \qquad (148)$$

where $\kappa'[(n+2)/2]$ is the exponent associated with the conductance fluctuations below the percolation threshold defined above. In the limit $n = 2$ Eq. (148) reduces to (142). Numerical simulations for testing the validity of these predictions were reported by Levy and Bergman (1993, 1994b) and Zhang (1996b).

### 3.3.4  *Exact Duality Relations*

Similar to linear and strongly nonlinear conducting composites, weakly nonlinear heterogeneous materials also satisfy some exact duality relations in 2D which we now describe. These relations were derived by Levy and Kohn (1998), and parallel those already described for strongly nonlinear composites in Section 3.1.6.

Consider a two-phase weakly nonlinear composite material for which the $I - V$ characteristic is given by Eq. (121) (written in terms of the voltage **V** rather than

the electric field $\mathbf{E}$). The local conductivities of the two-phase material are given by

$$g_j(\mathbf{V}) = g_j^{(\ell)} + g_j^{(n)}|\mathbf{V}|^n, \quad j = 1, 2. \tag{149}$$

The dual composite is another two-phase material with the same microgeometry, but with its phases having the following local conductivity,

$$g_j(|\mathbf{I}|) = \frac{1}{g_j(|\mathbf{V}|)} = \frac{1}{g_j^{(\ell)}} - \frac{g_j^{(n)}}{[g_j^{(\ell)}]^{n+2}}|\mathbf{I}|^n. \tag{150}$$

The effective conductivities of the two components can be expressed as

$$g^*[g_1^{(\ell)}(|\mathbf{V}|), g_2^{(\ell)}(|\mathbf{V}|); V_0] = g_e^{(\ell)} + g_e^{(n)}V_0^n, \tag{151}$$

for the primal composite, and

$$g_d^*[g_1^{(\ell)}(|\mathbf{I}|), g_2^{(\ell)}(|\mathbf{I}|); I_0] = g_e^{(\ell,d)} + g_e^{(n,d)}I_0^n, \tag{152}$$

for the dual composite, with

$$I_0 = g^*\left[g_1^{(\ell)}(|\mathbf{V}|), g_2^{(\ell)}(|\mathbf{V}|); V_0\right]V_0, \tag{153}$$

being the magnitude of the current that flows through the primal composite, which is also the magnitude of volume-averaged electric field in the dual composite. All the notations have the same meaning as for the strongly nonlinear composites discussed earlier. To first order in the local nonlinear conductivity $g^{(n)}$, the effective conductivities satisfy

$$g_e^{(\ell)} + g_e^{(n)}V_0^n = \frac{1}{g_e^{(\ell,d)}} - \frac{g_e^{(n,d)}}{[g_e^{(\ell,d)}]^2}I_0^n. \tag{154}$$

This relation leads us to

$$g_e^{(\ell)} = \frac{1}{g_e^{(\ell,d)}}, \tag{155}$$

which is the same as the well-known duality relation for linear composites, and

$$g_e^{(n)}V_0^n = -\frac{g_e^{(n,d)}}{[g_e^{(\ell,d)}]^2}I_0^n. \tag{156}$$

Equation (156) implies immediately that for cubic nonlinearity ($n = 2$),

$$\frac{g_e^{(n)}}{[g_e^{(\ell)}]^2} = -\frac{g_e^{(n,d)}}{[g_e^{(\ell,d)}]^2}. \tag{157}$$

Similar to the case of strongly nonlinear composites described in Section 3.1.6, we can extend this analysis to weakly nonlinear materials near the percolation threshold and investigate its consequences. As discussed in Section 3.3.2, if we

have a mixture of good conductors [conductances $g_M^{(\ell)}$ and $g_M^{(n)}$] and perfect insulators, then near $p_c$ we expect to have [see Eq. (130)]

$$\frac{g_e^{(n)}}{[g_e^{(\ell)}]^2} \sim \frac{g_M^{(n)}}{[g_M^{(\ell)}]^2}(p - p_c)^{-\kappa}, \tag{158}$$

where $\kappa$ is the exponent for the resistance noise introduced and described above. Similarly, for a mixture of normal conductors [conductances $g_I^{(\ell)}$ and $g_I^{(n)}$] and superconductors near $p_c$, one must have [see Eq. (132)]

$$\frac{g_e^{(n)}}{[g_e^{(\ell)}]^2} \sim \frac{g_I^{(n)}}{[g_I^{(\ell)}]^2}(p_c - p)^{-\kappa'}, \tag{159}$$

where the exponent $\kappa'$ was also defined above. Using the duality relations described above, one can then show that $\kappa = \kappa'$ which, as discussed above and in Chapter 5 of Volume I, also holds for linearly conducting composites.

The foregoing discussions can be extended to the case in which the ratio $g_I/g_M$, for both the linear and nonlinear conductivities, is finite. In this case Eq. (133) should hold for the primal composite and its dual, both above and below the percolation threshold $p_c$. Then, using the above duality relations, one can show that the scaling functions $\Phi_{\pm,I}$ and $\Phi_{\pm,M}$ [where the plus (minus) sign is for $p > p_c$ ($p < p_c$)] and their dual counterparts satisfy the following relations

$$\Phi_M = \Phi_I^{(d)}, \quad \Phi_I = \Phi_M^{(d)}, \tag{160}$$

with the understanding that if the left-hand side of Eqs. (160) uses the scaling function with the plus sign, then, the right-hand side uses the function with the minus sign, and vice versa.

### 3.3.5  Comparison with the Experimental Data

The relevance of the above models of weakly nonlinear composites and their properties to modeling real materials was established by experimental studies of Lin (1992), who measured $I - V$ characteristics of $PrBa_2Cu_3O_{7-\delta}$, a compound thought for a long time to be superconducting, although it now appears that it is a normal conductor, even at very low temperatures. Figure 3.8 presents the results for four different experiments with the same compound, indicating highly nonlinear behavior beyond a current of about $I_c \simeq 0.02$ A. If we assume that Eq. (121) describes the $I - V$ behavior of the material, then one may estimate the exponent $n$ by fitting the data to this equation. Lin found that $n = 1$ and 2 both represent the data relatively well. When the critical current $I_c$ was plotted versus the linear conductivity $g_e^{(\ell)}$, the data shown in Figure 3.9 were obtained. The straight line passing through the data has a slope $x \simeq 0.6$. On the other hand, Eq. (136) predicts that, $x = (1 + \kappa/\mu)/2$, which implies that $x \simeq 0.93$ in both 2D and 3D, if we use $\mu \simeq 1.3$ and 2.0, and $\kappa \simeq 1.12$ and 1.60 in 2D and 3D, respectively. This estimate of $x$ does not agree with Lin's measurements. However, if we use

FIGURE 3.8. A typical nonlinear $I - V$ curve for a $PrBa_2Cu_3O_{7-\delta}$ compound at 300 K. Symbols show the data for four different samples (after Lin, 1992).



FIGURE 3.9. Logarithmic plot of the critical current $I_c$ versus the effective linear conductivity $g_e^{(l)}$ for four samples of $PrBa_2Cu_3O_{7-\delta}$. Solid circles show the data for an Ag-] added sample. The straight line represents $I_c \sim [g_e^{(l)}]^{0.6}$ (after Lin, 1992).

$\mu \simeq 2.5$ and $\kappa \simeq 5.14$ for the 3D Swiss-cheese model, i.e., the model in which spherical inclusions are distributed randomly in a uniform matrix, then Eq. (133) predicts that $x \simeq 0.74$, which is only about 20% larger than Lin's measurements which, given the scatter in the data shown in Figure 3.9, is quite acceptable.

## 3.4  Dielectric Constant of Weakly Nonlinear Composites

Most of our analysis of the effective conductivity of nonlinear composites is equally applicable to the problem of computing the effective dielectric function of the same materials, with the effective conductivities replaced by the effective dielectric constant $\epsilon_e$. Thus, in this section we summarize the most important results and discuss their ramifications for the static case. Frequency-dependent dielectric constant will be described in the next section.

Consider a two-component composite material in which each component is described by a weakly cubic nonlinear relation between the electric displacement field $\mathbf{D}$ and the electric field $\mathbf{E}$ given by

$$\mathbf{D}_i = \epsilon_i^{(\ell)} \mathbf{E}_i + \epsilon_i^{(n)} |\mathbf{E}_i|^2 \mathbf{E}_i, \quad i = 1, 2. \tag{161}$$

In the analysis that follows we assume that, $\epsilon_i^{(n)} |\mathbf{E}|^2 \ll \epsilon_i^{(\ell)}$. We wish to compute the effective nonlinear dielectric function $\epsilon_e^{(n)}$ defined by

$$\langle \mathbf{D} \rangle = \epsilon_e^{(\ell)} \langle \mathbf{E} \rangle + \epsilon_e^{(n)} \langle |\mathbf{E}|^2 \rangle \langle \mathbf{E} \rangle, \tag{162}$$

where $\epsilon_e^{(\ell)}$ is the effective linear dielectric function of the composite when the electric field is small enough, and $\langle \cdot \rangle$ denotes an average over the volume of the composite.

A general approximate scheme for this problem was proposed by Zeng *et al.* (1988) which we now summarize and discuss. As in the case of the effective conductivity, the linear effective dielectric function can always be written as

$$\epsilon_e^{(\ell)} = F\left[\epsilon_1^{(\ell)}, \epsilon_2^{(\ell)}, p_1\right], \tag{163}$$

which is the analogue of Eq. (108). Here $p_1$ is the volume fraction of the $\epsilon_1$ component, and $F$ is an estimate of the effective dielectric constant which, in general, depends on the morphology of the composite. We initially assume that only component 1 is nonlinear, so that $\epsilon_2 = \epsilon_2^{(\ell)}$, and therefore we can invoke an approximate nonlinear form of Eq. (163):

$$\epsilon_e = F(\epsilon_1, \epsilon_2, p_1), \tag{164}$$

where, $\epsilon_i = \epsilon_i^{(\ell)} + \epsilon_i^{(n)} \langle |\mathbf{E}_i|^2 \rangle$, and $\langle |\mathbf{E}_i|^2 \rangle$ is the mean square of the electric field in the $i$th component in the *linear* limit. We must keep in mind that Eq. (164) is valid only if $\epsilon_1$ and $\epsilon_2$ are constant in their respective component, implying that $\mathbf{E}$ is uniform in the nonlinear component.

The function $F$ is now expanded in a Taylor series around $\epsilon_e^{(\ell)}$:

$$\epsilon_e \simeq F\left[\epsilon_1^{(\ell)}, \epsilon_2, p_1\right] + F'\left[\epsilon_1^{(\ell)}, \epsilon_2^{(\ell)}, p_1\right]\epsilon_1^{(n)}\langle|E_1|^2\rangle, \tag{165}$$

where $F' = \partial F/\partial\epsilon_1$. However, one can express $F'$ *exactly* in terms of the average squared electric field in component 1 in the *linear* limit:

$$p_1\frac{\langle|E_1|^2\rangle}{E_0^2} = \left[\frac{\partial\epsilon_e^{(\ell)}}{\partial\epsilon_1}\right]_\ell \equiv F'\left[\epsilon_1^{(\ell)}, \epsilon_2^{(\ell)}, p_1\right], \tag{166}$$

where $E_0$ is the external field. Therefore,

$$\epsilon_e = \epsilon_e^{(\ell)} + \frac{\epsilon_1^{(n)}}{p_1}F'|F'|^2 E_0^2, \tag{167}$$

which means that, by the definition of the effective nonlinear dielectric function $\epsilon_e^{(n)}$, we obtain

$$\epsilon_e^{(n)} = \frac{\epsilon_1^{(n)}}{p_1}\left(\frac{\partial\epsilon_e}{\partial\epsilon_1}\right)_\ell\left|\frac{\partial\epsilon_e}{\partial\epsilon_1}\right|_\ell. \tag{168}$$

Equation (168) is the analogue of Eq. (109) for the nonlinear conductivity.

We can generalize this result to composites in which both components are weakly nonlinear. Hence, we write

$$\epsilon_e = \epsilon_e^{(\ell)} + \frac{\epsilon_1^{(n)}}{p_1}F_1'|F_1'|E_0^2 + \frac{\epsilon_2^{(n)}}{p_2}F_2'|F_2'|E_0^2, \tag{169}$$

where $F_i' = \partial\epsilon_e/\partial\epsilon_i$ ($i = 1, 2$). Therefore,

$$\epsilon_e^{(n)} = \frac{\epsilon_1^{(n)}}{p_1}F_1'|F_1'| + \frac{\epsilon_2^{(n)}}{p_2}F_2'|F_2'|. \tag{170}$$

Equation (170) also suggests an analogous generalization for the effective nonlinear conductivity $g_e^{(n)}$, which would then represent a generalization of Eq. (109). One can now use this general method of approximation and study its properties in certain limits.

### 3.4.1 Exact Results

There are a few simple morphologies for which $\epsilon_e^{(\ell)}$, and hence $\epsilon_e^{(n)}$, can be computed exactly. In one such morphology the two components are arranged in the form of cylinders that are parallel to the external field. The cylinders do not have to have circular cross sections. For this model,

$$\epsilon_e^{(\ell)} = p_1\epsilon_1^{(\ell)} + p_2\epsilon_2^{(\ell)}. \tag{171}$$

Then, it is not difficult to see that,

$$\epsilon_e^{(n)} = p_1\epsilon_1^{(n)} + p_2\epsilon_2^{(n)}. \tag{172}$$

The second morphology for which the effective nonlinear dielectric function can be exactly computed is one in which the components are arranged in the form of flat slabs perpendicular to the external field. For this case,

$$\epsilon_e^{(\ell)} = \frac{1}{p_1/\epsilon_1^{(\ell)} + p_2/\epsilon_2^{(\ell)}}, \tag{173}$$

from which one obtains, using Eq. (170),

$$\epsilon_e^{(n)} = p_1 \frac{\epsilon_1^{(n)}}{[p_1 + \epsilon_1^{(\ell)} p_2/\epsilon_2^{(\ell)}]^4} + p_2 \frac{\epsilon_2^{(n)}}{[p_2 + \epsilon_2^{(\ell)} p_1/\epsilon_1^{(\ell)}]^4}. \tag{174}$$

A perturbation expansion, similar to what we described in Section 3.1.4 for the effective conductivity of strongly nonlinear composites, was also developed by Yu *et al.* (1993).

### 3.4.2  Effective-Medium Approximation

As the reader probably knows by now, according to the EMA, the effective dielectric constant is one of the solutions of the following quadratic equation,

$$p_1 \frac{\epsilon_1^{(\ell)} - \epsilon_e^{(\ell)}}{\epsilon_1^{(\ell)} + (Z/2 - 1)\epsilon_e^{(\ell)}} + p_2 \frac{\epsilon_2^{(\ell)} - \epsilon_e^{(\ell)}}{\epsilon_2^{(\ell)} + (Z/2 - 1)\epsilon_e^{(\ell)}} = 0. \tag{175}$$

If both $\epsilon_1^{(\ell)}$ and $\epsilon_2^{(\ell)}$ are real and positive, then the physically relevant solution of the EMA is also the positive one. Equation (175) is now solved for $\epsilon_e^{(\ell)}$, from the solution of which the functions, $F_i = \partial \epsilon_e^{(\ell)}/\partial \epsilon_i^{(\ell)}$, are computed which, when substituted in Eq. (170), yield the EMA prediction for the effective nonlinear dielectric constant $\epsilon_e^{(n)}$.

### 3.4.3  The Maxwell–Garnett Approximation

In Chapter 2, as well as Section 4.9.4 of Volume I, we described the Maxwell–Garnett (MG) approximation for the effective linear conductivity and dielectric constant of composite materials based on the continuum models. As discussed there, the MG approximation is most appropriate for a heterogeneous solid in which one of the components plays the role of a matrix, while the other acts as an inclusion. Therefore, assuming that component 2 is the matrix, the MG approximation takes the following form:

$$\epsilon_e^{(\ell)} = \frac{\epsilon_1^{(\ell)}(2p_1 + 1) + 2\epsilon_2^{(\ell)}(1 - p_1)}{\epsilon_1^{(\ell)}(1 - p_1) + \epsilon_2^{(\ell)}(2 + p_1)} \epsilon_2^{(\ell)} \tag{176}$$

Using Eq. (176), the functions $F_i' = \partial \epsilon_e^{(\ell)}/\partial \epsilon_i$ are computed which, when substituted in Eq. (170), yield the MG estimate for the nonlinear dielectric function

$\epsilon_e^{(n)}$. Hui (1990a) extended the MG approximation to a more general composite for which, $\mathbf{D} = \epsilon^{(\ell)}\mathbf{E} + \epsilon^{(n)}|\mathbf{E}|^n\mathbf{E}$

We should emphasize again that in the type of nonlinear problems that we are discussing here, the geometry of the system and the boundary conditions are very important and have a profound influence on the overall behavior of the system. As a matter of fact, every result described so far is valid only for two-terminal systems, and essentially nothing is known for multi-terminal ones.

## 3.5   Electromagnetic Field Fluctuations and Optical Nonlinearities

In this section we continue the discussion that we began in Chapter 4 of Volume I and describe and discuss advances in understanding optical properties of disordered materials, and the effect that constitutive nonlinearities may have on such properties. The main conceptual framework for our discussions are the discrete models, in the form of disordered resistor networks. Hence, we are particularly interested in the optical properties of composite materials with percolation-type disorder. In general, as our discussions in Section 3.3 made it clear, disordered solid materials with percolation-type disorder are very sensitive to the magnitude of the external electric field because, (1) their macroscopic transport and optical properties are controlled by their backbone, i.e., the current-carrying part of the network, and (2) because of the sparse morphology of the backbone, and in particular its low fractal structure at length scale $L \ll \xi_p$ ($D_{bb} \simeq 1.675$ and 1.8 in 2D and 3D, respectively), the effect of the external field accumulates around its weak points, i.e., its red bonds which are those that, if cut, would split the backbone into two pieces. Therefore, such materials should have, and indeed do have, much larger nonlinear macroscopic response than those of their constitutes.

Even when there is no apparent constitutive nonlinearities in the conduction properties of the phases of a disordered material, percolation disorder may lead to nonlinear macroscopic response. An interesting manifestation of this phenomenon was provided by the AC and DC conductivities of a percolation composite of carbon particles embedded in a wax matrix (Bardhan, 1997). In this composite, neither the carbon particles nor the wax matrix exhibits any nonlinearity in their conduction properties; nevertheless, the macroscopic conductivity of the composite increases significantly when the applied voltage increases by only a few volts. Such a strong nonlinear response can be attributed to quantum tunneling between the conducting carbon particles, a distinct feature of electrical transport in disordered solids near the percolation threshold $p_c$.

Likewise, local fluctuations in the electromagnetic field and the resulting enhancement of nonlinear optical properties in disordered solids, such as metal-dielectric composites with percolation disorder, especially near $p_c$, constitute an important set of phenomena, since such composites have high potential for various applications. Nonlinear effects manifest themselves in two distinct ways:

(1) If the applied electric field or current exceeds a critical threshold, then, at zero frequency, strong nonlinearity results in the breakdown of the conducting elements of a composite. The critical field decreases to zero as the volume fraction of the conducting component approaches $p_c$, hence indicating that such composites become progressively more responsive to the external field as $p_c$ is approached. This phenomenon is what we have referred to as threshold nonlinearity; it will be studied in Chapter 5.

(2) Alternatively, although increasing the external voltage or current may not result in electric or dielectric breakdown of a composite, it can lead to very large enhancements of the nonlinearities as the volume fraction of the conducting component approaches $p_c$. We already described this phenomenon in Sections 3.1 and 3.3 in terms of the crossover from a linearly conductive material to a weakly nonlinear one, and our goal in this section is to do the same for optical properties of the same type of composite solids.

Following our discussions in Section 3.3, we consider in this section weak nonlinearities so that the field-dependent conductivity $g(E)$ can be written as a power series in the applied electric field $E$, with the leading term, i.e., the linear conductivity $g^{(\ell)}$, being much larger than the higher-order terms, a situation which is typical of various nonlinearities in the optical and infrared spectral ranges of interest to us. As discussed in Section 3.3, despite this weakness, such nonlinearities lead to qualitatively new phenomena, such as enhancement of higher harmonics in percolation composites, and the occurrence of bistable behavior of the composite (Bergman *et al.*, 1994; Levy *et al.*, 1995) in which the conductivity switches between two stable values. In such disordered materials, especially those that contain metal particles that are characterized by a dielectric constant with negative real and small imaginary parts, the fluctuations in the local field are strongly enhanced in the optical and infrared spectral ranges, leading to enhancement of various nonlinear properties. If the disorder in the morphology of such solid materials is of percolation-type, then they are potentially of great practical importance (see, for example, Flytzanis, 1992) as composites with intensity-dependent dielectric functions and, in particular, as nonlinear filters and optical bistable elements. The optical response of such nonlinear composites can be easily tuned by, for example, controlling the volume fraction and morphology of their constitutes.

More generally, optical properties of fractal aggregates of metal particles have been studied. These studies indicate that a fractal morphology results in very large enhancement of various nonlinear responses of the aggregates within the spectral range of their plasmon resonances. The typical size, $a \sim 10$ nm, of the metal particles in such fractal aggregates is much smaller than the wavelength $\lambda > 300$ nm in the optical and infrared spectral ranges. Since the average density of particles in fractal aggregates is much smaller than in non-fractal materials, and approaches zero with increasing size of the aggregates, it is possible to consider each particle in the aggregate as an elementary dipole and introduce the corresponding interaction operator. If this is done, then, solving the problem of the optical response of metal fractal aggregates reduces to diagonalizing the interaction operator for the

light-induced dipoles. If the size of the fractal aggregate is not very large, the diagonalization can be done numerically (and efficiently) and thus the local electric field can be calculated (see, for example, Stockman *et al.*, 1995, 1996; Stockman, 1997; Shalaev *et al.*, 1993; Markel *et al.*, 1999). Computations of this type indicate that large field fluctuations are localized in some small parts of the fractal aggregate and change with the wavelength. These predictions and numerical computations of large enhancements of optical nonlinearities in metal fractals have also been verified experimentally for degenerate four-wave mixing and nonlinear refraction and absorption. In these experiments, aggregation of silver particles (which were initially isolated) into fractal clusters led to six orders of magnitude enhancement of the efficiency of the nonlinear four-wave process and about three orders of magnitude enhancement in the nonlinear refraction and absorption. The localized and strongly fluctuating local fields in these fractal aggregates were imaged by means of the near-field scanning optical microscopy (Shalaev *et al.*, 1993; Markel *et al.*, 1999). A similar pattern was obtained for the field distribution in self-affine thin films (Shalaev *et al.*, 1996a,b; Safonov *et al.*, 1998). As discussed in Chapter 1, such self-affine films possess a fractal surface with different scaling properties in the plane of the film and normal to it.

Despite such progress, the distribution of the local field and the corresponding nonlinearities were, until recently, poorly understood for metal-dielectric composites with percolation-type disorder, especially in the most interesting spectral range where the plasmon resonances occur in the metal grains. As shown in Section 3.3, if a small volume fraction $p \ll 1$ of a nonlinear material is embedded in a linear host, the effective nonlinear response of the composite can be calculated explicitly. As one may expect, the nonlinearities are enhanced at the frequency $\omega_r$ corresponding to the plasmon resonance of a single metal grain. Numerical calculations (Stroud and Zhang, 1994; Zhang and Stroud, 1994) for a finite $p$ also indicate considerable enhancement in the narrow frequency range around $\omega_r$ and, moreover, the system sizes that can currently be used in the computations are not large enough for drawing quantitative conclusions about the nonlinear properties for frequencies $\omega \simeq \omega_r$. However, we should recognize that a small system size $L$ may act as an artificial damping factor that cuts off all the fluctuations in the local field when the spatial separation is larger than $L$, hence resulting in a corresponding decrease of the nonlinearities which may otherwise not be seen in a large enough sample.

An alternative method to numerical simulations is the effective-medium approximation (EMA) that has the virtue of mathematical and conceptual simplicity. We already described in Sections 3.1 and 3.3 such EMAs for nonlinear composites near $p_c$. As discussed there, for the static case the predictions of the nonlinear EMA (Wan *et al.*, 1996; Hui *et al.*, 1997) are in good agreement with numerical simulations for 2D percolation composite. However, despite this success, application of any type of nonlinear EMA is suspect for the frequency range corresponding to the plasmon resonances in metal grains. This is due to the fact that both computer simulations and experimental data for the field distribution in percolation composites indicate that the distribution contains sharp peaks that are separated

by distances that are much larger than the metal grain size. Thus, the local electric field cannot be assumed to be the same in all the metal grains of the composite, implying that the main pillar of the EMA, i.e., the assumption of a uniform field, fails for the frequency range corresponding to the plasmon resonance in the films.

To address this problem, a new theory of the distribution of the electromagnetic field and nonlinear optical processes in metal-dielectric composites was developed (Sarychev and Shalaev, 1999; Sarychev *et al.*, 1999). The theory is based on the concepts of percolation processes, and takes advantage of the fact that the problem of optical excitations in percolation composites can be mapped onto the Anderson localization problem. It predicts localization of surface plasmons (SP) in composites with percolation disorder, and describes in detail the localization pattern. It also indicates that the SP eigenstates are localized on length scales that are much smaller than the wavelength of an incident light. The eigenstates with eigenvalues that are close to zero (resonant modes) are excited most efficiently by the external field. Since the eigenstates are localized and only a small portion of them is excited by the incident beam, overlapping of the eigenstates can typically be neglected, a fact that significantly simplifies the theoretical analysis and allows one to derive relatively simple expressions for enhancement of linear and nonlinear optical responses.

The purpose of this section is to describe and summarize this progress. An excellent comprehensive review of this subject was presented by Sarychev and Shalaev (2000). This section is patterned closely after their review and represents a summary of their discussions. Since the languages of nonlinear currents/conductivities and nonlinear polarizations/susceptibilities, or dielectric constants, are completely equivalent, they will be used interchangeably in this section.

## 3.5.1  *Scaling Properties of Moments of the Electric Field*

As already demonstrated in Chapters 2, 5, and 6 of Volume I and earlier in the present chapter, in metal-dielectric percolation composites the effective static (DC or zero frequency) conductivity $g_e$ decreases with decreasing volume fraction $p$ of the metal component, and vanishes at $p = p_c$. Since for $p < p_c$ the effective DC conductivity $g_e = 0$, the material is dielectric-like. Therefore, a metal-insulator transition takes place at the percolation threshold $p_c$. However, although the transition at $p_c$ is second-order, the pattern of the fluctuations in percolation composites appears to be quite different from that for a second-order phase transition, the fluctuations of which are usually characterized by long-range correlations, with their relative magnitudes being of the order of unity. In contrast, for (DC) percolation conductivity, the local electric fields are concentrated on the edges of large metal clusters, so that the field maxima (large fluctuations or peaks) are separated by distances that are of the order of the percolation correlation length $\xi_p$. Since $\xi_p$ diverges at $p_c$ (recall that near $p_c$, $\xi_p \sim |p - p_c|^{-\nu}$), the implication is that the distance between the field maxima or peaks also increases as $p_c$ is approached.

To obtain insight into the high-frequency properties of metals, consider first a simple model—the Drude model (already utilized in Chapters 4 and 6 of Volume

I)—that reproduces semi-quantitatively the basic optical properties of a metal. According to this model, the dielectric constant $\epsilon_m$ of metal grains is given by

$$\epsilon_m(\omega) = \epsilon_b - \frac{(\omega_p/\omega)^2}{1 + i\omega_\tau/\omega}, \tag{177}$$

where $\epsilon_b$ is the contribution to $\epsilon_m$ due to the inter-band transitions, $\omega_p$ is the plasma frequency, and $\omega_\tau = 1/\tau \ll \omega_p$ is the relaxation rate (in Chapters 4 and 6 we took $\epsilon_b = 1$). In the high-frequency range considered here, losses in the metal grains are relatively small, $\omega_\tau \ll \omega$. Therefore, if we write, $\epsilon_m = \epsilon_m' + i\epsilon_m''$, then $|\epsilon_m'|/\epsilon_m'' \simeq \omega/\omega_\tau \gg 1$. Moreover, one has, $\epsilon_m' < 0$ for the frequencies $\omega < \tilde{\omega}_p$, where $\tilde{\omega}_p$ is the renormalized plasma frequency which is given by

$$\tilde{\omega}_p = \frac{\omega_p}{\sqrt{\epsilon_b}}. \tag{178}$$

Therefore, the metal conductivity, $g_m = -i\omega\epsilon_m/4\pi \simeq (\epsilon_b\tilde{\omega}_p^2/4\pi\omega)[i(1 - \omega^2/\tilde{\omega}_p^2) + \omega_\tau/\omega]$, is characterized by the dominant imaginary part for $\tilde{\omega}_p > \omega \gg \omega_\tau$, i.e., it is of inductive character. In this sense, the metal grains can be thought of as inductances $L$, while the dielectric gaps between the metal grains can be represented by capacitances $C$. Then, the percolation composite represents a set of randomly distributed $L$ and $C$ elements. The collective surface plasmons, excited by the external field, can be thought of as resonances in different $L - C$ circuits, and the excited surface plasmon eigenstates represent giant fluctuations of the local field.

### 3.5.1.1   Distribution of Electric Fields in Strongly Disordered Composites

Before embarking on discussing the properties of the distribution of local electric field in a composite, let us recall from Chapters 5 and 6 of Volume I how the dielectric constant of a disordered material is computed via a discrete, percolation-type model. Suppose that a percolation composite is illuminated by light and consider the local optical field distributions in the material. A typical metal grain size $a$ in the composite is much smaller than $\lambda$, the wavelength of the light in the visible and infrared spectral ranges. If so, then one can introduce a potential $\phi(\mathbf{r})$ for the local electric field and write the local current density $\mathbf{I}$ as, $\mathbf{I}(\mathbf{r}) = g(\mathbf{r})[-\nabla\phi(\mathbf{r}) + \mathbf{E}_0]$, where $\mathbf{E}_0$ is the external field, and $g(\mathbf{r})$ is the local conductivity at $\mathbf{r}$. In the quasi-static limit, computation of the field distribution reduces to finding the solution of the Poisson's equation since, due to current conservation, $\nabla \cdot \mathbf{I} = 0$, one has

$$\nabla \cdot \{g(\mathbf{r})[-\nabla\phi(\mathbf{r}) + \mathbf{E}_0]\} = 0, \tag{179}$$

where the local conductivity $g(\mathbf{r}) = g_m$ or $g_d$ for the metal and dielectric components, respectively. We rewrite Eq. (179) in terms of the local dielectric constant, $\epsilon(\mathbf{r}) = 4\pi i g(\mathbf{r})/\omega$, so that

$$\nabla \cdot [\epsilon(\mathbf{r})\nabla\phi(\mathbf{r})] = \mathcal{E}, \tag{180}$$

where $\mathcal{E} = \nabla \cdot [\epsilon(\mathbf{r})\mathbf{E}_0]$. The external field $\mathbf{E}_0$ can be real, while $\phi(\mathbf{r})$ is, in general, a complex function since $\epsilon_m$ is complex in the optical and infrared spectral ranges.

Since Eqs. (179) and (180) are difficult to solve analytically, one discretizes them in order to solve them by numerical simulations. If, for example, a standard 5-point (in 2D) or 7-point (in 3D) finite-difference discretization is used, then, a discrete model on a simple-cubic lattice is obtained in which the metal and dielectric particles are represented by metal and dielectric bonds of the lattice. Thus, Eq. (180), in discretized form, takes on the form of Kirchhoff's equations defined on a lattice. Assuming that the external electric field $\mathbf{E}_0$ is directed along the $z$-axis, one obtains

$$\sum_j \epsilon_{ij}(\phi_j - \phi_i) = \sum_j \epsilon_{ij} E_{ij} \tag{181}$$

where $\phi_i$ is the electric potential at site $i$ of the lattice, and the sum is over the nearest neighbors $j$ of the site $i$. For the bonds $ij$ in the $\pm z$-direction, the electromotive force $E_{ij}$ is given by, $E_{ij} = \pm E_0 a_0$ (where $a_0$ is the spatial period of the lattice), while $E_{ij} = 0$ for the other bonds that are connected to site $i$. Thus, the composite material is modeled by a resistor-capacitor-inductor network in which the bond permittivities $\epsilon_{ij}$ are statistically independent and $a_0$ is equal to the metal grain size, $a_0 = a$. In the case of a two-component metal-dielectric random composite, the permittivities $\epsilon_{ij}$ take values $\epsilon_m$ and $\epsilon_d$ with probabilities $p$ and $1 - p$, respectively. To make further progress, we use a simple-cubic lattice which has a very large but finite number of sites $N$ and rewrite Eq. (181) in a matrix form:

$$\mathcal{H}\boldsymbol{\phi} = \mathcal{E}, \tag{182}$$

where $\boldsymbol{\phi} = \{\phi_1, \phi_2, \ldots, \phi_N\}$, and the elements of the vector $\mathcal{E}$ are, $\mathcal{E}_i = \sum_j \epsilon_{ij} E_{ij}$. Here $\mathcal{H}$ is a $N \times N$ matrix such that for $i \neq j$, $\mathcal{H}_{ij} = -\epsilon_{ij} = \epsilon_d > 0$ and $\epsilon_m = (-1 + i\kappa)|\epsilon'_m|$ with probabilities $p$ and $1 - p$, respectively, and $\mathcal{H}_{ii} = \sum_j \epsilon_{ij}$, where $j$ refers to nearest neighbors of site $i$, and $\kappa$ is the usual loss factor, $\kappa = \epsilon''_m / |\epsilon'_m| \ll 1$. The diagonal elements $\mathcal{H}_{ii}$ are distributed between $2d\epsilon_m$ and $2d\epsilon_d$, where $d$ is the dimensionality of the space.

Similar to the dielectric constant, we write $\mathcal{H} = \mathcal{H}' + i\kappa\mathcal{H}''$, where $i\kappa\mathcal{H}''$ represents losses in the system. The Hamiltonian $\mathcal{H}'$ formally coincides with the Hamiltonian of the problem of metal-insulator transition (Anderson transition) in quantum systems, i.e., it maps the quantum-mechanical Hamiltonian for the Anderson transition problem with both on- and off-diagonal correlated disorder onto the present problem. Hereafter, we refer to $\mathcal{H}'$ as the Kirchhoff's Hamiltonian (KH). Thus, the problem of determining the solution of Kirchhoff's equation, Eq. (181) or (182), is equivalent to the eigenfunction problem for the KH, $\mathcal{H}'\Psi_n = \Lambda_n\Psi_n$, whereas the losses can be treated as perturbations.

Since $\epsilon'_m < 0$, and the permittivity $\epsilon_d$ of the dielectric matrix is positive, the set of the KH eigenvalues $\Lambda_n$ contains eigenvalues with real parts that are equal (or close) to zero. Then, eigenstates $\Psi_n$ that correspond to eigenvalues $\Lambda_n$ such that, $|\Lambda_n| \ll |\epsilon_m|$, and $|\epsilon_d|$, are strongly excited by the external field and are seen as giant field fluctuations, representing the resonant surface plasmon modes. If one assumes that the eigenstates excited by the external field are localized, then they should look like the peaks of the local field with the average distance between

them being about $a(N/n)^{1/d}$, where $n$ is the number of the KH eigenstates excited by the external field.

Consider now the special case when $\epsilon_m' = -\epsilon_d$, which corresponds to the plasmon resonance of individual particles in a 2D system. Since a solution to Eq. (181) does not change if $\epsilon_m$ and $\epsilon_d$ are multiplied by the same factor, we normalize the system and set $\epsilon_d = -\epsilon_m = 1$. We also suppose, for simplicity, that the metal concentration is $p = 0.5$. In this case, the eigenstates $\Psi_n$ are all localized. On the other hand, computer simulations (Müller *et al.*, 1997) showed that there is a transition from chaotic (Berry, 1977) to localized eigenstates for the 2D Anderson problem, with a crossover region between the two. Consider first the case when the metal volume fraction $p = p_c = 1/2$ for the 2D bond percolation problem. Then, the diagonal disorder in the KH is characterized by, $\langle \mathcal{H}_{ii}' \rangle = 0$, and $\langle \mathcal{H}_{ii}'^2 \rangle = 4$, which correspond to the chaos-localization transition (Müller *et al.*, 1997). Moreover, $\mathcal{H}'$ also possesses strong off-diagonal disorder, $\langle \mathcal{H}_{ij}' \rangle = 0$, which favors localization (see, for example, Verges, 1998). There is therefore strong evidence that the eigenstates $\Psi_n$ are localized for all $\Lambda_n$ in the 2D system, although one cannot rule out the possibility of inhomogeneous localization, similar to that obtained for fractal clusters (see, for example, Stockman *et al.*, 1994), or power-law localization (Kaveh and Mott, 1981; Kramer and MacKinnon, 1993).

In the case of $\epsilon_d = -\epsilon_m = 1$ and $p = 1/2$, all parameters in $\mathcal{H}'$ are of the order of unity, and therefore its properties do not change under the transformation $\epsilon_d \Longleftrightarrow \epsilon_m$. Therefore, the real eigenvalues $\Lambda_n$ are distributed symmetrically around zero in an interval of the order of one. The eigenstates with eigenvalues $\Lambda_n \simeq 0$ are effectively excited by the external field and represent the giant local field fluctuations. When $p$ decreases (increases), the eigenstates with eigenvalues $\Lambda_n \simeq 0$ are shifted from the center of the distribution toward its lower (upper) edge, which typically favors localization. Because of this, one may assume that in 2D the eigenstates, or at least those with eigenvalues $\Lambda_n \simeq 0$, are localized for all metal volume fractions $p$.

The situation in 3D is much more complex. Despite the great effort and the progress that has been made, the Anderson transition in 3D is not yet fully understood. Computer simulations (Kawarabayashi *et al.*, 1998) of Anderson localization in 3D [with $\epsilon_d = -\epsilon_m = 1$, $p = 1/2$, the diagonal matrix elements $w_{ii}$ distributed uniformly around 0, $-w_0/2 \leq w_{ii} \leq w_0/2$, and the off-diagonal elements $w_{ij} = \exp(i\phi_{ij})$, with phases $\phi_{ij}$ also distributed uniformly in $0 \leq \phi_{ij} \leq 2\pi$] show that in the center of the band the states are localized for the disorder $w_0 > w_c = 18.8$. In the 3D $\mathcal{H}'$ Hamiltonian discussed here, the diagonal elements are distributed as $-6 \leq \mathcal{H}_{ii} \leq 6$, and therefore the diagonal disorder is smaller than the critical disorder $w_c$, but the off-diagonal disorder is stronger than in the calculations of Kawarabayashi *et al.* (1998). It has been shown (Verges, 1998; Elimes *et al.*, 1998) that even small off-diagonal disorder strongly enforces localization, and thus one may conjecture that, in the 3D case, the eigenstates corresponding to the eigenvalues $\Lambda_n \simeq 0$ are also localized for all $p$.

If we express the potential $\phi$ in Eq. (182) in terms of the eigenfunctions $\Psi_n$ of $\mathcal{H}'$ as, $\phi = \sum_n A_n \Psi_n$, and substitute it in Eq. (182), we obtain the following

equation for coefficients $A_n$:

$$(i\kappa b_n + \Lambda_n)A_n + i\kappa \sum_{m\neq n} \left(\Psi_n \left|\mathcal{H}''\right| \Psi_m\right) A_m = \mathcal{E}_n, \tag{183}$$

where $b_n = \left(\Psi_n \left|\mathcal{H}''\right| \Psi_n\right)$, and $\mathcal{E}_n = (\Psi_n|\mathcal{E})$ is a projection of the external field onto the eigenstate $\Psi_n$. Since all the parameters in $\mathcal{H}'$ are of the order of unity, the $b_n$ are also of the order of unity and can be approximated by some constant $b \simeq 1$. Sarychev and Shalaev (2000) suggested that the eigenstates $\Psi_n$ are localized within spatial domains $\xi_A(\Lambda)$, where $\xi_A(\Lambda)$ is the Anderson localization length. Then, the sum in Eq. (183) is convergent and can be treated as a small perturbation. The first two coefficients in the approximation are then given by

$$A_n^{(0)} = \frac{\mathcal{E}_n}{\Lambda_n + i\kappa b}, \tag{184}$$

whereas

$$A_n^{(1)} = -i\kappa \sum_{m\neq n} \left(\Psi_n \left|\mathcal{H}''\right| \Psi_m\right) A_m^{(0)}. \tag{185}$$

In Eq. (185), the most important eigenstates in the sum, in the limit $\kappa \to 0$, are those with eigenvalues $|\Lambda_m| \leq b\kappa$. Since the eigenvalues $\Lambda_n$ are distributed in an interval of the order of unity, the spatial density of the eigenmodes with $|\Lambda_m| \leq b\kappa$ vanishes as $a^{-d}\kappa \to 0$ as $\kappa \to 0$, implying that $A_n^{(1)}$ is exponentially small, $|A_n^{(1)}| \sim |\sum_{m\neq n} \left(\Psi_n|\mathcal{H}''|\Psi_m\right) \mathcal{E}_m/b_m| \propto \exp\left\{-[a/\xi_A(0)]\kappa^{-1/d}\right\}$, and can be neglected when $\kappa \ll [a/\xi_A(0)]^d$. Then, the local potential $\phi$ is given by, $\phi(\mathbf{r}) = \sum_n A_n^{(0)}\Psi_n = \sum_n \mathcal{E}_n \Psi_n(r)/(\Lambda_n + i\kappa b)$, and the fluctuating part of the local field $\mathbf{E}_f = -\nabla\phi(\mathbf{r})$ is given by

$$\mathbf{E}_f(\mathbf{r}) = -\sum_n \mathcal{E}_n \left[\nabla\Psi_n(\mathbf{r})/(\Lambda_n + i\kappa b)\right], \tag{186}$$

where $\nabla$ is understood as a lattice operator. The average field intensity is then given by

$$\left\langle |\mathbf{E}|^2 \right\rangle = \left\langle |\mathbf{E}_f + \mathbf{E}_0|^2 \right\rangle = |\mathbf{E}_0|^2 + \left\langle \sum_m \sum_n \frac{\mathcal{E}_n \mathcal{E}_m^* [\nabla\Psi_n(\mathbf{r}) \cdot \nabla\Psi_m^*(\mathbf{r})]}{(\Lambda_n + i\kappa b)(\Lambda_m - i\kappa b)} \right\rangle, \tag{187}$$

where we used the fact that $\langle \mathbf{E}_f \rangle = \langle \mathbf{E}_f^* \rangle = 0$.

Consider now the eigenstates $\Psi_n$ with eigenvalues $\Lambda_n$ within a small interval $|\Lambda_n - \Lambda| \leq \Delta\Lambda \ll \kappa$ centered at $\Lambda$, which we denote them by $\Psi_n(\Lambda r)$. Recall that the eigenstates are assumed to be localized so that eigenfunctions $\Psi_n(\Lambda r)$ are well-separated in space, with the average distance $l$ between them being, $l(\Delta\Lambda) \sim a[\mathcal{N}(\Lambda)\Delta\Lambda]^{-1/d}$, where

$$\mathcal{N}(\Lambda) = \frac{a^d}{\Omega} \sum_n \delta(\Lambda - \Lambda_n), \tag{188}$$

is the dimensionless density of states for the Kirchhoff Hamiltonian (KH) $\mathcal{H}'$, and $\Omega$ is the system's volume. We assume that the metal volume fraction $p \simeq 1/2$, so that all quantities in the KH $\mathcal{H}'$ are about unity, and therefore the density of states $\mathcal{N}(\Lambda)$ is also about unity at $\Lambda = 0$. Hence, the distance $l(\Delta\Lambda)$ can be arbitrary large as $\Delta\Lambda \to 0$, while it is still much smaller than the system size. It is further assumed that the total number of eigenstates $\Psi_n(\Lambda\mathbf{r})$ is large. When $l(\Delta\Lambda) \gg \xi_A(\Lambda)$, the localized eigenfunctions $\Psi_n(\Lambda\mathbf{r})$ are characterized by spatial positions of their centers $\mathbf{r}_n$, so that $\Psi_n(\Lambda\mathbf{r}) = \Psi(\Lambda\mathbf{r} - \mathbf{r}_n)$ and Eq. (187) becomes

$$\left\langle |\mathbf{E}|^2 \right\rangle = |\mathbf{E}_0|^2 + \sum_{\Lambda_1}\sum_{\Lambda_2} \frac{\left\langle \sum_m \sum_n \mathcal{E}_n \mathcal{E}_m^* \left[ \boldsymbol{\nabla}\Psi(\Lambda_1, \mathbf{r} - \mathbf{r}_n) \cdot \boldsymbol{\nabla}\Psi^*(\Lambda_2, \mathbf{r} - \mathbf{r}_m) \right] \right\rangle}{(\Lambda_1 + i\kappa b)(\Lambda_2 - i\kappa b)},$$
(189)

where the first sums are over positions of the intervals $|\Lambda_n - \Lambda_1|$ and $|\Lambda_m - \Lambda_2|$ in the $\Lambda$ space, whereas the sums in the numerator are over spatial positions $\mathbf{r}_n$ and $\mathbf{r}_m$ of the eigenfunctions. For each realization of a macroscopically-homogeneous disordered film, the positions $\mathbf{r}_n$ of the eigenfunctions $\Psi(\Lambda\mathbf{r} - \mathbf{r}_n)$ take on new values that do not correlate with $\Lambda$. Therefore, we can independently carry out the averaging in the numerator in the second term of Eq. (189) over positions $\mathbf{r}_m$ and $\mathbf{r}_n$ of eigenstates $\Psi_m$ and $\Psi_n$. Since, $\langle \boldsymbol{\nabla}\Psi_n(\mathbf{r}) \rangle = 0$, we obtain

$$\left\langle \mathcal{E}_n \mathcal{E}_m^* \left[ \boldsymbol{\nabla}\Psi(\Lambda_1, \mathbf{r} - \mathbf{r}_n) \cdot \boldsymbol{\nabla}\Psi^*(\Lambda_2, \mathbf{r} - \mathbf{r}_m) \right] \right\rangle$$
$$\simeq \left\langle |\mathcal{E}_n|^2 |\boldsymbol{\nabla}\Psi(\Lambda_1, \mathbf{r} - \mathbf{r}_n)|^2 \right\rangle \delta_{\Lambda_1 \Lambda_2} \delta_{nm},$$
(190)

which, when substituted in Eq. (187), results in

$$\left\langle |\mathbf{E}|^2 \right\rangle = |\mathbf{E}_0|^2 + \sum_{\Lambda} \frac{\sum_n |\mathcal{E}_n|^2 \left\langle |\boldsymbol{\nabla}\Psi_n(\Lambda, \mathbf{r})|^2 \right\rangle}{\Lambda^2 + (b\kappa)^2}.$$
(191)

The localized eigenstates are not in general degenerate, so that the eigenfunctions $\Psi_n$ can be selected to be real, i.e., $\Psi_n = \Psi_n^*$ (where $*$ denotes the complex conjugate). Then, $|\mathcal{E}_n|^2 = |(\Psi_n|\mathcal{E})|^2 = \left| \sum_{i=1}^N \Psi_{n,i}\mathcal{E}_i \right|^2 \sim a^{-2d} |\int \Psi_n \mathcal{E} d\mathbf{r}|^2$, which, after using (180) and (181), yields

$$|\mathcal{E}_n|^2 \sim a^{4-2d} \left| \int \Psi_n(\mathbf{E}_0 \cdot \boldsymbol{\nabla}\epsilon) d\mathbf{r} \right|^2 = a^{4-2d} \left| \int \epsilon(\mathbf{E}_0 \cdot \boldsymbol{\nabla}\Psi_n) d\mathbf{r} \right|^2.$$
(192)

Since the local dielectric constants $|\epsilon|$ are of the order of unity, one can write, $\boldsymbol{\nabla}\Psi_n \sim \Psi_n/\xi_A(\Lambda)$, and therefore,

$$|\mathcal{E}_n|^2 \sim \frac{|\mathbf{E}_0|^2 a^4}{a^{2d}\xi_A^2(\Lambda)} \left| \int \Psi_n(\mathbf{r}) d\mathbf{r} \right|^2 \sim \frac{|\mathbf{E}_0|^2 a^4}{\xi_A^2(\Lambda)} \left| \sum_{i=1}^N \Psi_{n,i} \right|^2.$$
(193)

Using the fact that, $\langle \Psi_n | \Psi_n \rangle = \sum_{i=1}^N |\Psi_{n,i}|^2 = 1$, and that $\Psi_n$ are localized within $\xi_A(\Lambda)$, one obtains $\Psi_{n,i} \sim [\xi_A(\Lambda)/a]^{-d/2}$ in the localization domain which, when

substituted in Eq. (193), yields

$$|\mathcal{E}_n|^2 \sim |\mathbf{E}_0|^2 \, a^2 [\xi_A \Lambda)/a]^{d-2}. \tag{194}$$

One can estimate in a similar way,

$$\left\langle |\nabla \Psi_n(\Lambda, \mathbf{r})|^2 \right\rangle \sim \xi_A^{-2}(\Lambda) \left\langle |\Psi_n(\Lambda, \mathbf{r})|^2 \right\rangle \sim \xi_A^{-2}(\Lambda) N^{-1} \sum_{i=1}^{N} |\Psi_{n,i}|^2 \sim \xi_A^{-2}(\Lambda), \tag{195}$$

where $N = \Omega/a^d$ is the total number of sites. Using these estimates and taking into account the fact that the total number of the eigenstates within interval $\Delta \Lambda$ is equal to $N\mathcal{N}(\Lambda)\Delta \Lambda$, one finally obtains

$$\left\langle |\mathbf{E}|^2 \right\rangle \sim |\mathbf{E}_0|^2 + |\mathbf{E}_0|^2 \int \frac{\mathcal{N}(\Lambda)[a/\xi_A(\Lambda)]^{4-d}}{\Lambda^2 + (b\kappa)^2} d\Lambda. \tag{196}$$

Since all matrix elements in the Hamiltonian $\mathcal{H}'$ are of the order of unity (in fact, the off-diagonal elements are $\pm 1$), the density of states $\mathcal{N}(\Lambda)$ and localization length $\xi_A(\Lambda)$ vary significantly within an interval of the order of one, while the denominator in Eq. (191) has an essential singularity at $\Lambda = \pm i b\kappa$. Then, the second moment of the local electric field, $M_2 \equiv M_{2,0} = \left\langle |\mathbf{E}|^2 \right\rangle / |\mathbf{E}_0|^2$, is estimated as

$$M_2^\star \sim 1 + \mathcal{N}(a/\xi_A)^{4-d} \int \frac{1}{\Lambda^2 + (b\kappa)^2} d\Lambda \sim \mathcal{N}(a/\xi_A)^{4-d} \kappa^{-1} \gg 1, \tag{197}$$

provided that $\kappa \ll \mathcal{N}(a/\xi_A)^{4-d}$ [we set $\xi_A(\Lambda = 0) \equiv \xi_A$, $\mathcal{N}(\Lambda = 0) \equiv \mathcal{N}$ and $b \simeq 1$]. Thus, in this case, the field distribution is described as a set of the KH eigenstates localized within $\xi_A$, with its peaks having the amplitudes

$$E_m^\star \sim E_0 \kappa^{-1} (a/\xi_A)^2, \tag{198}$$

which are separated by the field correlation length $\xi_e^\star$,

$$\xi_e^\star \sim a(\mathcal{N}\kappa b)^{-1/d} \sim a(\mathcal{N}\kappa)^{-1/d}. \tag{199}$$

All the assumptions that led us to Eqs. (197)–(199) hold when $\xi_e^\star \gg \xi_A$, which is fulfilled in the limit $\kappa \to 0$.

Hereafter by superscript $^\star$ we mark the fields, while the spatial scales are given for the special case $-\epsilon_m' = \epsilon_d = 1$ considered here (note that $^\star$ should not be confused with the complex conjugation denoted by $*$), while for $\xi_A$ and $\mathcal{N}$ we omit the $^\star$ sign in order to avoid complex notations; it is implied that their values are always taken at $-\epsilon_m' = \epsilon_d = 1$, even if the case of $|\epsilon_m/\epsilon_d| \gg 1$ is considered.

The assumption that the localization length $\xi_A$ is proportional to the eigenstate size might not, in general, be true for the Anderson systems, although it has been confirmed well by numerical calculations for 2D percolation composites. It was also assumed that the metal volume fraction $p \simeq 1/2$, which corresponds to the 2D percolation threshold, and that the density of states $\mathcal{N}(\Lambda)$ is finite and about unity for $\Lambda = 0$. The latter assumption is, however, violated for small values of $p$ when the distribution of the eigenvalues shifts to the positive side of $\Lambda$, so that the eigenstates with eigenvalues $\Lambda \simeq 0$ are shifted to the lower edge of the distribution,

and the density of states $\mathcal{N}$ in Eq. (197) becomes a function of $p$. In the limit $p \to 0$, the number of states effectively excited by the external field is proportional to the number of metal particles, and hence $\mathcal{N}(p) \sim p$. The same consideration holds in the opposite limit, $p \to 1$, and therefore $\mathcal{N}(p) \sim 1 - p$. When $\mathcal{N}$ decreases, localization becomes stronger and one can write, $\xi_A(\Lambda = 0, p \to 0) \sim \xi_A(\Lambda = 0, p \to 1) \sim a$. When $p \to 0$ or $p \to 1$, the number of the field maxima decreases while the peaks become progressively sharper. Equation (197) also indicates that strong field fluctuations ($M_2 > 1$) exist in a metal-dielectric composite with $\epsilon_d = -\epsilon'_m$ in a wide range of concentrations,

$$\kappa < p < 1 - \kappa, \quad \kappa \ll 1. \tag{200}$$

Although the above local fields were estimated for the special case of $\epsilon_d = -\epsilon'_m$, all the above results, which are based on the assumption that the eigenstates of the Kirchhoff Hamiltonian are localized, hold in a more general case, when the real part $\epsilon'_m$ of the metal dielectric constant is negative and its absolute value is of the order of $\epsilon_d$. The important case of $|\epsilon_m| \gg \epsilon_d$ will be considered in the next subsection.

### 3.5.1.2  Moments of the Electric Field

Consider now the moments of the local electric field of arbitrary order, defined as

$$M_{n,m} = \frac{1}{\omega E_0^m |\mathbf{E}_0|^n} \int |E(\mathbf{r})|^n \, E^m(\mathbf{r}) \, d\mathbf{r}, \tag{201}$$

where, as above, $E_0 \equiv E^{(0)}$ is the amplitude of the external field, and $E(\mathbf{r})$ is the local field at $\mathbf{r}$. We denote, for simplicity, $M_{n,0} = M_n$, and assume that a volume-averaged quantity is equivalent to its ensemble-averaged value, i.e., $M_{n,m} = \langle |E|^n \, E^m \rangle / E_0^m |E_0|^n$.

The high-order moment $M_{2k,m} \propto \langle E^{k+m} E^{*k} \rangle$ represents a nonlinear optical process in which in one elementary act $k + m$ photons are added and $k$ photons are annihilated (see, for example, Boyd, 1992). This is because the complex-conjugated field in the general expression for the nonlinear polarization implies photon annihilation, so that the corresponding frequency enters the nonlinear susceptibility with a minus sign. Enhancement of the Kerr optical nonlinearity $G_K$ is proportional to $M_{2,2}$, the enhancement of the third-harmonic generation is given by $|M_{0,3}|^2$, and surface-enhanced Raman scattering is represented by $M_{4,0}$ (see below).

An important case is when $M_{n,m} \gg 1$, i.e., when the fluctuating part of the local electric field $\mathbf{E}_f$ is much larger than $\mathbf{E}_0$. Suppose, for simplicity, that $\mathbf{E}_0$ is real and that $|\mathbf{E}_0| = 1$. We can write, for the moment $M_{2p,2q}$ ($p$ and $q$ are integers), the following equation

$$M_{2p,2q} = \left\langle \sum_{n_1,n_2,\cdots n_{2p};m_1,m_2,\cdots m_{2q}} \frac{\mathcal{E}_{n_1}\mathcal{E}_{n_2}\left(\nabla\Psi_{n_1}\cdot\nabla\Psi_{n_2}^*\right)\cdots\mathcal{E}_{n_{2p-1}}\mathcal{E}_{n_{2p}}\left(\nabla\Psi_{n_{2p-1}}\cdot\nabla\Psi_{n_{2p}}^*\right)}{\left(\Lambda_{n_1}+ibk\right)\left(\Lambda_{n_2}-ibk\right)\cdots\left(\Lambda_{n_{2p-1}}+ibk\right)\left(\Lambda_{n_{2p}}-ibk\right)} \right.$$
$$\left. \times \frac{\mathcal{E}_{m_1}\mathcal{E}_{m_2}\left(\nabla\Psi_{m_1}\cdot\nabla\Psi_{m_2}\right)\cdots\mathcal{E}_{m_{2q-1}}\mathcal{E}_{m_{2q}}\left(\nabla\Psi_{m_{2q-1}}\cdot\nabla\Psi_{m_{2q}}\right)}{\left(\Lambda_{m_1}+ibk\right)\left(\Lambda_{m_2}+ibk\right)\cdots\left(\Lambda_{m_{2q-1}}+ibk\right)\left(\Lambda_{m_{2q}}+ibk\right)} \right\rangle,$$
$$\tag{202}$$

where $\langle \cdot \rangle$ denotes an ensemble average (which, as discussed above, is equivalent to the volume-average), and the sums are over all eigenstates of the KH $\mathcal{H}'$. We now average Eq. (202) over spatial positions of eigenstates $\Psi_n(\mathbf{r}) \equiv \Psi(\mathbf{r} - \mathbf{r}_n)$ to obtain

$$M_{2p,2q} \sim \sum_{\Lambda} \frac{\sum_{|\Lambda_n - \Lambda| \le \Delta\Lambda} |\mathcal{E}_n|^{2p} \, \mathcal{E}_n^{2q} \left\langle \left(\nabla\Psi_n \cdot \nabla\Psi_n^*\right)^p \left(\nabla\Psi_n \cdot \nabla\Psi_n\right)^q \right\rangle}{\left[\Lambda^2 + (bk)^2\right]^p (\Lambda + ibk)^{2q}},$$
(203)

where the summation in the numerator is over eigenfunctions $\Psi_n = \Psi(\Lambda, \mathbf{r} - \mathbf{r}_n)$ with eigenvalues within the interval $|\Lambda_n - \Lambda| \le \Delta\Lambda \ll \kappa$, while the external sum is over positions $\Lambda$ of the intervals that cover the entire range of eigenvalues $\Lambda_n$. Following the same line of arguments that was used for deriving Eq. (197), one can show that (Sarychev and Shalaev, 2000)

$$M_{2p,2q} \sim \int \frac{\mathcal{N}(\Lambda) \left[a/\xi_A(\Lambda)\right]^{4(p+q)-d}}{\left[\Lambda^2 + (b\kappa)^2\right]^p (\Lambda + ib\kappa)^{2q}} d\Lambda.$$
(204)

Assuming that the density of states $\mathcal{N}(\Lambda)$ and the localization length $\xi_A(\Lambda)$ are both smooth functions of $\Lambda$ in the vicinity of zero, and taking into account the fact that all parameters of the Hamiltonian $\mathcal{H}'$ for the case $\epsilon_d = -\epsilon_m' = 1$ are of the order of one, the following estimate for the moments of the local field is obtained

$$M_{n,m}^{\star} \sim \mathcal{N}(p)[a/\xi_A(p)]^{2(n+m)-d} \, \kappa^{-n-m+1},$$
(205)

for $n + m > 1$ and $m > 0$ (for simplicity we set $b = 1$). We remind the reader once again that $\mathcal{N}(p)$ and $\xi_A(p)$ should be understood as $\mathcal{N}(p) = \mathcal{N}(p, \Lambda = 0)$ and $\xi_A(p) = \xi_A(p, \Lambda = 0)$, i.e., they are given at the eigenvalue $\Lambda = 0$.

The maximum of the Anderson localization length $\xi_A(\Lambda)$ is typically at the center of the distribution of the eigenvalues $\Lambda$ (Kawarabayashi *et al.*, 1998). When $p \ne 1/2$, $\Lambda = 0$ moves from the center of the $\Lambda$-distribution toward its tails where the localization is typically stronger (i.e., $\xi_A$ is smaller). Therefore, it is plausible that $\xi_A(p)$ reaches its maximum at $p = 1/2$ and decreases toward $p = 0$ and $p = 1$, so that the absolute values of the moments of the local field may have a minimum at $p = 1/2$. In 2D composites the percolation threshold $p_c$ is typically close to 0.5. Therefore, in such composites the moments $M_{n,m}$ do have a local minimum at $p_c$ as a function of the metal volume fraction $p$, and the amplitudes of various nonlinear processes, while much enhanced, have a characteristic *minimum* at $p_c$. It is important to note that the magnitude of the moments in Eq. (205) do not depend on the number of annihilated photons in one elementary act of the nonlinear scattering. However, when all photons are added (i.e., when all frequencies enter the nonlinear susceptibility with the plus sign) and $n = 0$, one cannot estimate the moments $M_{0,m}$ by Eq. (205), since the integral in Eq. (204) is no longer determined by the poles at $\Lambda = \pm ib\kappa$. However, all the functions of the integrand are about unity and $M_{0,m} \sim O(1)$ for $m > 1$. The moment $M_{0,m}$ is an important quantity since it yields the enhancement $G_{nHG}$ of the $n$th order harmonic generation through the relation, $G_{nHG} = |M_{0,m}|^2$ (see below).

### 3.5.1.3    Field Fluctuations at Frequencies Below the Resonance

So far we have assumed that $|\epsilon_m|/\epsilon_d \simeq 1$, which corresponds to the plasmon resonance in the metal grains. To estimate the fluctuations in the local field in percolation composites for $|\epsilon_m|/\epsilon_d \gg 1$, the renormalization approach developed by Shalaev and Sarychev (1998), Sarychev and Shalaev (1999) and Sarychev *et al.* (1999) can be utilized. Let us briefly recall the main concepts of the renormalization method (see also Chapter 5 of Volume I). Consider a percolation composite with the metal volume fraction $p = p_c$. The system is divided into cubic cells of size $b$, each of which is considered as a new renormalized element. The cells are classified into two types: Those that contain a continuous path of metallic particles are considered as conducting, while those without such a sample-spanning cluster are considered as non-conducting, or dielectric. The effective dielectric constant $\epsilon_m(b)$ of a conducting cell decreases with increasing its size $b$ as, $\epsilon_m(b) \simeq (b/a)^{-\mu/\nu}\epsilon_m$, whereas the effective dielectric constant $\epsilon_d(b)$ of a dielectric cell increases with $b$ as $\epsilon_d(b) \simeq (b/a)^{s/\nu}\epsilon_d$, where $\mu$, $s$ and $\nu$ are the usual percolation critical exponents for the conductivity, dielectric constant, and percolation correlation length, respectively (see above and Chapters 2, 5 and 6 of Volume I). The cube size $b$ is now taken to be

$$b = b_r = a(|\epsilon_m|/\epsilon_d)^{\nu/(\mu+s)}. \tag{206}$$

Let us recall that the exponent $s$ also characterizes the power law behavior of the effective conductivity of a conductor-superconductor composite near the percolation threshold. Then, in the renormalized system the dielectric constant of the new elements either takes a value, $\epsilon_m(b_r) = \epsilon_d^{\mu/(\mu+s)} |\epsilon_m|^{s/(\mu+s)} (\epsilon_m/|\epsilon_m|)$, for the renormalized conducting cell, or $\epsilon_d(b_r) = \epsilon_d^{\mu/(\mu+s)} |\epsilon_m|^{s/(\mu+s)}$, for the renormalized dielectric cell. The ratio of the dielectric constants of these new elements is then, $\epsilon_m(b_r)/\epsilon_d(b_r) = \epsilon_m/|\epsilon_m| \simeq -1 + i\kappa$, where the loss-factor $\kappa = \epsilon_m''/|\epsilon_m| \ll 1$ is the same as in the original system. As discussed in Chapter 5 of Volume I, at $p = p_c$, the volume fraction of conducting and dielectric elements does not change under a renormalization transformation. Since the field distribution in a two-component system depends on the ratio of the dielectric permittivities of the components, after the renormalization the problem becomes equivalent to what was discussed above for the case $\epsilon_d = -\epsilon_m' = 1$. Taking into account the fact that the electric field renormalizes as $E_0^\star = E_0(b_r/a)$, one obtains from Eq. (198) the following expression for the field's peaks in the renormalized system:

$$E_m \simeq E_0(a/\xi_A)^2(b_r/a)\kappa^{-1} \simeq E_0(a/\xi_A)^2 \left(\frac{|\epsilon_m|}{\epsilon_d}\right)^{\nu/(\mu+s)} \left(\frac{|\epsilon_m|}{\epsilon_m''}\right), \tag{207}$$

where $\xi_A = \xi_A(p_c)$ is the localization length in the renormalized system. Each maximum of the field in the renormalized system is in a dielectric gap in a dielectric cube of linear size $b_r$ or in between two conducting cells of the size $b_r$ that are not necessarily connected to each other. There is not a characteristic length in the original system which is smaller than $b_r$, except the grain size $a$. Therefore, it is

plausible that the width of a peak of the local field in the original system is about $a$. Then, values of the field maxima $E_m$ do not change when returning from the renormalized system to the original one. Hence, Eq. (207) yields values of the field maxima in the original system.

Equation (207) provides the estimate for the local field extrema when the real part $\epsilon'_m$ of the dielectric constant is negative. For metals $\epsilon_m$ increases in absolute value with the wavelength, when the frequency $\omega < \tilde{\omega}_p$. Therefore, the field maxima $E_m(\omega)$ increase strongly with the wavelength. For a Drude metal the steep growth of the peaks $E_m(\omega)$ occurs for the frequencies $\omega < \tilde{\omega}_p$, when the dielectric constant $\epsilon_m$ can be approximated as

$$\epsilon_m(\omega < \tilde{\omega}_p) \simeq 2(\omega - \tilde{\omega}_p)\frac{\epsilon_b}{\tilde{\omega}_p} + i\frac{\epsilon_b \omega_\tau}{\tilde{\omega}_p}, \tag{208}$$

which, when substituted in Eq. (207), yields

$$E_m(\omega < \tilde{\omega}_p) \simeq E_0(a/\xi_A)^2 \left(\frac{2\epsilon_b \left|\omega - \tilde{\omega}_p\right|}{\tilde{\omega}_p}\right)^{(\nu+\mu+s)/(\mu+s)} \frac{\tilde{\omega}_p}{\omega_\tau \epsilon_b \epsilon_d^{\nu/(\mu+s)}}. \tag{209}$$

Since in a typical metal, $\omega_\tau \ll \tilde{\omega}_p$, the amplitudes of the field's peak first increase steeply and then saturate (see below) at $E_m \simeq E_0(a/\xi_A)^2(\epsilon_b/\epsilon_d)^{\nu/(\mu+s)}(\tilde{\omega}_p/\omega_\tau) \sim E_0\tilde{\omega}_p/\omega_\tau$, when $\omega \simeq 0.5\tilde{\omega}_p$. Therefore, the intensity maxima $I_m$ exceed the intensity of the incident wave $I_0$ by a factor $I_m/I_0 \sim (\tilde{\omega}_p/\omega_\tau)^2 \gg 1$.

Consider now the case $\omega \ll \omega_p$, when for a Drude metal

$$\epsilon_m(\omega \ll \omega_p) \simeq -\left(\frac{\omega}{\omega_p}\right)^2 \left(1 - i\frac{\omega_\tau}{\omega}\right), \quad \omega \gg \omega_\tau \tag{210}$$

which, when substituted in Eq. (207), yields

$$E_m(\omega \ll \omega_p) \simeq E_0\left(\frac{a}{\xi_A}\right)^2 \left(\frac{\omega_p}{\sqrt{\epsilon_d}\omega}\right)^{2\nu/(\mu+s)} \left(\frac{\omega}{\omega_\tau}\right). \tag{211}$$

For 2D percolation, the critical exponents are, $\mu = s \simeq \nu = 4/3$, and thus Eq. (211) yields, $E_m \sim E_0(a/\xi_A)^2\omega_p/(\sqrt{\epsilon_d}\omega_\tau) = E_0(a/\xi_A)^2(\tilde{\omega}_p/\omega_\tau)\sqrt{\epsilon_b/\epsilon_d} \sim E_0(\tilde{\omega}_p/\omega)$, which coincides with the estimate obtained from Eq. (209) for $\omega = 0.5\tilde{\omega}_p$, implying that the local field's peaks increase steeply when $\epsilon'_m$, the real part of $\epsilon_m$, is negative and then remains essentially constant in the wide frequency range, $\tilde{\omega}_p < \omega < \omega_\tau$.

For 3D percolation composites, we roughly have, $\nu \simeq (\mu + s)/3$, and thus Eq. (211) yields, $E_m \sim E_0(\epsilon_b/\epsilon_d)^{1/3}\tilde{\omega}_p^{2/3}\omega^{1/3}/\omega_\tau$, implying that the local field peaks increase up to $E_m/E_0 \sim \tilde{\omega}_p/\omega_\tau$ when $\epsilon'_m < 0$, and then decrease as $E_m/E_0 \sim (\tilde{\omega}_p/\omega_\tau)(\omega/\tilde{\omega}_p)^{1/3}$, if the frequency decreases further.

To obtain $M_{n,m}$ we consider first the spatial distribution of the maxima of the field for $|\epsilon_m| \gg \epsilon_d$. The average distance between the maxima in the renormalized system is $\xi_e^\star$, given by Eq. (199). Then, the average distance $\xi_e$ between the maxima

in the original system (provided that $\mathcal{N} \sim 1$) is

$$\xi_e \simeq (b_r/a)\xi_e^\star \sim a \left( \frac{|\epsilon_m|}{\epsilon_d} \right)^{\nu/(\mu+s)} \left( \frac{|\epsilon_m|}{\epsilon_m''} \right)^{1/d}, \qquad (212)$$

which in 2D (with $\mu = s \simeq \nu = 4/3$) reduces to a simple form,

$$\xi_e \sim a \frac{|\epsilon_m|}{\sqrt{\epsilon_d \epsilon_m''}}. \qquad (213)$$

In the renormalized system a typical "area" of a peak of the field corresponds to $\xi_A^d$, implying that in the original system each maximum is stretched over $(\xi_A/a)^d$ clusters of the size $b_r$. In each of these clusters the field maximum splits into $n(b_r)$ peaks of amplitude $E_m$, distributed along a dielectric gap in the dielectric square of size $b_r$. Since the gap area scales as the capacitance of the dielectric square, one has

$$n(b_r) \propto (b_r/a)^{d-2+s/\nu}, \qquad (214)$$

and therefore

$$M_{n,m} \sim (\xi_A/a)^d \left( \frac{E_m}{E_0} \right)^{n+m} \frac{n(b_r)}{(\xi_e/a)^d}$$

$$\sim \mathcal{N}(\xi_A/a)^{d-2(n+m)} \left( \frac{|\epsilon_m|}{\epsilon_d} \right)^{[(n+m-2)\nu+s]/(\mu+s)} \left( \frac{|\epsilon_m|}{\epsilon_m''} \right)^{n+m-1}, \qquad (215)$$

for $n + m > 1$ and $n > 0$. Since $|\epsilon_m| \gg \epsilon_d$ and $|\epsilon_m|/\epsilon_m'' \gg 1$, $M_{n,m} \gg 1$ in the visible and infrared spectral ranges. We emphasize that the localization length $\xi_A$ in Eq. (215) corresponds to the renormalized system with $\epsilon_d = -\epsilon_m' = 1$. The localization length in the original system, i.e., a typical size of the eigenfunction, is about $(b_r/a)\xi_A \gg a$, i.e., the eigenstates become macroscopically large when $|\epsilon_m|/\epsilon_d \gg 1$, and consist of sharp peaks separated in space by distances much larger than $a$.

It is then not difficult to show, using Eq. (215), that

$$M_{0,m} \sim M_{0,m}^\star (b_r/a)^m \frac{n(b_r)}{(\xi_e/a)^d} \sim \left( \frac{\epsilon_m''}{|\epsilon_m|} \right) \left( \frac{|\epsilon_m|}{\epsilon_d} \right)^{(m-2+s/\nu)\nu/(\mu+s)}, \qquad (216)$$

which holds when $M_{0,m} > 1$. In 2D, if we use $\mu = s \simeq \nu = 4/3$, Eqs. (215) and (216) are simplified to

$$M_{n,m} \sim \mathcal{N} \left[ \frac{|\epsilon_m|^{3/2}}{(\xi_A/a)^2 \sqrt{\epsilon_d \epsilon_m''}} \right]^{n+m-1}, \qquad (217)$$

for $n + m > 1$ and $n > 0$, and

$$M_{0,m} \sim \frac{\epsilon_m'' |\epsilon_m|^{(m-3)/2}}{\epsilon_d^{(m-1)/2}}, \qquad (218)$$

for $m > 1, n = 0$ and $(\epsilon_m|/\epsilon_d)^{(m-1)/2} > |\epsilon_m|/\epsilon_m''$. The moments $M_{n,m}(n \neq 0)$ are strongly enhanced in 2D Drude metal-dielectric composites since they reach the maximum value

$$M_{n,m} \sim \mathcal{N} \left[ \frac{\omega_p}{\omega_\tau \sqrt{\epsilon_d}(\xi_A/a)^2} \right]^{n+m-1}, \qquad (219)$$

when $\omega \ll \omega_p$. Thus, in a 2D percolation composite the moments $M_{n,m}$ are independent of frequency if $\omega \ll \omega_p$. For metals this typically takes place in the red and infrared spectral ranges. For example, for a semi-continuous silver film on a glass substrate, the moments $M_{n,m}$ can be estimated as, $M_{n,m} \sim [3 \times 10^2 (a/\xi_A)^2]^{n+m-1}$, for $\omega \ll \omega_p$.

It follows from Eq. (215) that for 3D metal-dielectric percolation composites, for which the dielectric constant of the metal component can be estimated by the Drude formula, the moments $M_{n,m}(n \neq 0)$ achieve their maximum at frequency $\omega_{\max} \simeq 0.5\tilde{\omega}_p$. Since, as mentioned above, for 3D percolation, $v/(\mu + s) \simeq 1/3$, the maximum value of $M_{n,m}$ is roughly given by

$$M_{n,m}(\omega = \omega_{\max}) \sim \mathcal{N}(\xi_A/a) \left[ (a/\xi_A)^2 \, (\epsilon_b/\epsilon_d)^{1/3} \, \tilde{\omega}_p/\omega_\tau \right]^{n+m-1}, \qquad (220)$$

whereas for $\omega \ll \omega_p$,

$$M_{n,m} \left( \omega \ll \omega_p \right) \sim \mathcal{N}(\xi_A/a) \left[ \frac{(a/\xi_A)^2 \omega_p^{2/3} \omega^{1/3}}{\epsilon_d^{1/3} \omega_\tau} \right]^{n+m-1}. \qquad (221)$$

Figure 3.10 compares the results of numerical and theoretical calculations for



FIGURE 3.10. Moments $M_{n,m}$ of the electric field in semicontinuous silver films versus the wavelength $\lambda$ at the percolation threshold. On the left are the moments $M_n = M_{n,0}$, from the bottom to the top, for $n = 2, 3, 4, 5$ and $6$. The solid curves are the predictions of the scaling theory, Eq. (215), while the symbols are the numerical simulation data. Shown on the right are the moments $M_{4,0}$ (upper solid curve predicted by the scaling theory, versus $*$, the numerical data), $M_{0,4}$ (upper dashed curve), $M_{2,0}$ (lower solid curve predicted by the scaling theory, versus $+$, the numerical data), and $M_{0,2}$ (lower dashed curve predicted by the scaling theory, versus circles, the numerical data) (after Sarychev and Shalaev, 2000).

$M_{n,m}$ in a 2D semi-continuous silver film on glass, indicating excellent agreement between the scaling theory and numerical simulations, where $\xi_A \simeq 2a$ was used. The small value of $\xi_A$ indicates that, at least in 2D, there is strong localization of surface plasmons in percolation composites. Note that, as discussed above, nonlinear optical processes are, in general, phase dependent with their phase dependence being through the term $E^m$ and their enhancement being $M_{n,m} = \langle |E/E_0|^n (E/E_0)^m \rangle$. According to the above analysis, $M_{n,m} \sim M_{n+m,0} \equiv M_{n+m}$, for $n \geq 1$. Thus, for example, enhancement of the Kerr-type nonlinearity, $I_K = M_{2,2}$, is proportional to the enhancement of the Raman scattering, $I_{RS} \simeq M_4$.

So far, it has been assumed that, when analyzing the case of $\epsilon'_m \ll 0$, the metal volume fraction $p$ equals $p_c$. We now consider the range $\Delta p = p - p_c$, where the above estimates for $M_{n,m}$ are valid. First, note that the above expressions for the local field and the average moments $M_{n,m}$ of the field hold for almost all values of $p$ given by Eq. (200) when $\epsilon_m \simeq -\epsilon_d$. The metal volume fraction range $\Delta p$ shrinks, however, where the local electric field is strongly enhanced and $\epsilon'_m \ll 0$. The above analysis was based on the finite-size scaling analysis (see Chapter 2 of Volume I for description of the finite-size scaling method), which holds provided that $l_r < \xi_p$, where $\xi_p$ is the percolation correlation length. Since at $p_c$ the correlation length $\xi_p$ diverges, these estimates are valid in the wide frequency range $\omega_\tau < \omega < \tilde{\omega}_p$, which includes the visible, infrared, and far-infrared spectral ranges for typical metals. For any particular frequency from this interval, one can estimate the range $\Delta p$, where the giant field fluctuations occur, by requiring that, $l_r = \xi_p$, which results in, $|\Delta p| \leq (\epsilon_d/|\epsilon_m|)^{1/(\mu+s)}$. Therefore, the local electric field fluctuates strongly for such volume fractions and its moments $M_{n,m}$ are much enhanced.

### 3.5.1.4   Computer Simulation

A number of EMAs, as well as position-space renormalization group (PSRG) methods, of the type described in Section 5.11 of Volume I, have been proposed for calculation of optical properties of semi-continuous disordered films. However, none of these methods allows one to calculate the field fluctuations and the effects resulting from them. Because semi-continuous metal films are of great theoretical and practical interest, it is important to study statistical properties of the electromagnetic fields in their near zone. To simplify the theoretical considerations, one may assume that the electric field is homogeneous in the direction perpendicular to the film plane, implying that the skin depth $\delta$ for the metal grains is large, $\delta \simeq c/(\omega\sqrt{|\epsilon_m|}) \gg a$, where $a$ is the grain size, so that the quasi-static approximation holds (see also Chapter 4 of Volume I). Note that the role of the skin effect can be very important, resulting, in many cases, in strong alterations of the electromagnetic response found in the quasi-static approximation (see, for example, Sarychev *et al.*, 1995; Levy-Nathansohn and Bergman, 1997). At the same time, the quasi-static approximation simplifies significantly theoretical considerations of the field fluctuations and describes well the optical properties of semi-continuous films, providing qualitative, and in some cases, quantitative, agreement with experimental data.

In the discussion that follows, the skin effect is neglected so that a semi-continuous film can be considered as a 2D material. In the optical frequency range, where the frequency $\omega$ is much larger than the relaxation rate $\tau^{-1}$ of the metallic component, a semi-continuous metal film can be thought of as a 2D $L - R - C$ lattice (see, for example, Brouers $et\ al.$, 1993). As before, the capacitance $C$ represents the gaps between metal grains that are filled by the dielectric material (substrate) with a dielectric constant $\epsilon_d$. The inductive elements $L - R$ represent the metallic grains that for the Drude metal have the dielectric function $\epsilon_m(\omega)$ given by Eq. (177). In the high-frequency range considered here, the losses in the metal grains are small, $\omega \gg \omega_\tau$. Therefore, $\epsilon_m' \gg \epsilon_m''$ (in modulus) and $\epsilon_m' < 0$ for frequencies $\omega < \tilde{\omega}_p = \omega_p/\sqrt{\epsilon_b}$. Thus, the metal conductivity is almost purely imaginary and the metal grains can be modeled as $L$-$R$ elements, with the active component being much smaller than the reactive one. If the skin effect cannot be neglected, i.e., if the skin depth $\delta < a$, the simple quasi-static presentation of a semi-continuous film as a 2D array of the $L - R$ and $C$ elements is not valid. One can still use the $L - R - C$ model in the other limiting case, when the skin effect is very strong ($\delta \ll a$). In this case, the losses in the metal grains are small, regardless of value of $\omega/\omega_\tau$, whereas the effective inductance for a metal grain depends on the grain size and shape rather than on the material constants for the metal.

It is instructive to consider first the properties of the film at $p = p_c$, where the duality relation (see above and also Chapters 4 and 5 of Volume I) predicts that, the effective dielectric constant $\epsilon_e$ in the quasi-static case is given exactly by, $\epsilon_e = \sqrt{\epsilon_d \epsilon_m}$. If we neglect the metal losses and set $\omega_\tau = 0$, the metal dielectric constant $\epsilon_m < 0$ for $\omega < \tilde{\omega}_p$. We also neglect possible small losses in a dielectric substrate, assuming that $\epsilon_d$ is real and positive, in which case $\epsilon_e$ is purely imaginary for $\omega < \tilde{\omega}_p$. Therefore, a film consisting of loss-free metal and dielectric grains is absorptive for $\omega < \tilde{\omega}_p$. The effective absorption in a loss-free film means that the electromagnetic energy is stored in the system and thus the local fields could increase without limit. In reality, due to losses the local fields in a metal film are, of course, finite. However, if the losses are small, one may expect very strong fluctuations in the field. To calculate Rayleigh and Raman scattering, and various nonlinear effects in a semi-continuous metal film, one must know the field and current distributions in the film.

Although, as discussed in Chapters 4 and 5 of Volume I, there are several very efficient numerical methods for calculating the effective conductivity of composite materials, they typically do not allow calculations of the field distributions. Brouers $et\ al.$ (1997) developed a PSRG method, a generalization of what was described in Chapter 5 of Volume I, using a square lattice of the $L - R$ (metal) and $C$ (dielectric) bonds. A fraction $p$ of the bonds were metallic ($L - R$ bonds) and had a conductivity $g_m = -i\epsilon_m\omega/4\pi$, while the dielectric ($C$) bonds, with a fraction $1 - p$, had a conductivity $g_d = -i\epsilon_d\omega/4\pi$. The applied field $E_0$ was $E_0 = 1$, whereas the local fields inside the system were of course complex quantities. In this method, after each RG transformation, an external field $E_0$ is applied to the system and the Kirchhoff's equations are solved in order to determine the fields and the currents in all the bonds of the transformed lattice. The self-dual PSRG cell

FIGURE 3.11. $2 \times 2$ renormalization group cells in two and three dimensions.

of Figure 3.11 was used which, because of its hierarchical structure, allows these equations to be solved exactly. Then, the one-to-one correspondence between the elementary bonds of the transformed lattice and the bonds of the initial square lattice was used for determining the field distributions, as well as the effective conductivity, of the initial lattice. The number of operations for obtaining the full distributions of the local fields is proportional to $b^2$ [to be compared with $O(b^7)$ operations needed in the transform-matrix method and $O(b^3)$ operations needed in the Lobb-Frank algorithm that was described in Sections 5.14.2 and 5.14.3 of Volume I]. The Drude formula for metal dielectric functions was used, and thin films of silver (for which $\epsilon_b = 5$, the plasma frequency $\omega_p = 9.1$ eV, and the relaxation frequency $\omega_\tau = 0.021$ eV) and gold (for which $\epsilon_b = 6.5$, $\omega_p = 9.3$ eV, and $\omega_\tau = 0.03$ eV), deposited on a glass substrate with the dielectric constant $\epsilon_d = 2.2$, were modeled.

All the numerical results obtained with this method were in agreement with the predictions of the scaling theory discussed above, as well as with experimental data, described below.

### 3.5.1.5    Comparison with the Experimental Data

Optical properties of metal-insulator thin films have been intensively studied, both experimentally and by computer simulations. Semi-continuous thin metal films are usually produced by thermal evaporation or sputtering of metals onto an insulating substrate. At first, small metallic grains are formed on the substrate. As the film grows, the metal volume fraction increases and irregularly-shaped clusters are formed on the substrate, resulting in 2D fractal morphologies. The size of these structures diverges at $p_c$ where a percolating cluster of metal is formed, and a continuous conducting path appears between two opposite ends of the sample.

The metal-insulator transition is very close to this point, even in the presence of quantum tunneling. At higher surface coverage, the film is mostly metallic, with voids of irregular shapes. As coverage increases further, the film becomes uniform. Optical properties of such metal-dielectric films exhibit anomalous phenomena that are absent for bulk metal and dielectric components. For example, the anomalous absorption in the near-infrared spectral range leads to unusual behavior of the transmittance and reflectance in that, the transmittance is much higher than that of continuous metal films, whereas the reflectance is much lower.

The predictions of the PSRG computations have been compared with the experimental data for gold-on-glass films at various wavelengths (Sarychev and Shalaev, 2000). There is good qualitative agreement between the two. The data for such disordered metal-dielectric films near $p_c$ suggest localization of optical excitations in small nm-scale hot spots. The hot spots of a percolation film represent very large local fields (fluctuations); spatial positions of the spots strongly depend on the light frequency. Near-field spectra observed and calculated at various points of the surface consist of several spectral resonances, the spectral locations of which depend on the probed site of the sample. These features are observable only in the near zone. In the far zone, one observes images and spectra in which the hot spots and the spectral resonances are averaged out. The local field enhancement is large, which is especially important for nonlinear processes of the $n$th order, and are proportional to the enhanced local fields to the $n$th power. This opens up a fascinating possibility for *nonlinear* near-field spectroscopy of single nano-particles and molecules.

### 3.5.2 Anomalous Light Scattering from Semicontinuous Metal Films

A quantitative analysis of the spatial distribution of the local field fluctuations, and light scattering induced by such fluctuations, are now carried out. The resonance frequency $\omega_r$, corresponding to the condition $\epsilon'_m(\omega_r) = -\epsilon_d$ is considered first which, for a Drude metal, is fulfilled at the frequency

$$\omega_r = \omega_p \sqrt{\frac{1}{\epsilon_b + \epsilon_d} - \left(\frac{\omega_\tau}{\omega_p}\right)^2} \simeq \frac{\omega_p}{\sqrt{\epsilon_b + \epsilon_d}}, \qquad (222)$$

where it has been assumed that $\omega_\tau = 1/\tau \ll \omega_p$, which is the case for a typical metal. Then, the metal dielectric function is, $\epsilon_m(\omega_r) = \epsilon_d(-1 + i\kappa)$, where the loss factor $\kappa$ is given by, $\kappa \simeq (1 + \epsilon_b/\epsilon_d)\omega_\tau/\omega_r \ll 1$. In modeling the distribution of the local field fluctuations, we take advantage of the fact that, since this distribution does not change when bond conductances are multiplied by the same factor, it is convenient to consider a lattice in which a bond conductance is $g_m = -1 + i\kappa$ with probability $p$ (the $L$ bonds) and $g_d = 1$ with probability $1 - p$ (the $C$ bonds). Since the absolute values of $g_d$ and $g_m$ are very close, the standard method based on the percolation theory and scaling analysis cannot be used for estimating the spatial distribution of the field. One may, however, use the PSRG method described

above to carry out the analysis, which yields interesting results. For example, using a system of size $b = 1024$, $p = p_c = 0.5$, and $\omega = \omega_r$, Sarychev and Shalaev (2000) calculated the electric field in all the bonds for $10^{-4} \leq \kappa \leq 10^{-1}$ with the external field being $E_0 = 1$. The distribution of the field intensity, $I(\mathbf{r}) = |E(\mathbf{r})|^2$, was found to be close to the well-known log-normal distribution, with its values spread over many orders of magnitude, even for $\kappa = 10^{-1}$. For $\kappa = 10^{-4}$, $I(\mathbf{r})$ was distributed essentially uniformly in $(0, 10^4)$. The average intensity, $\langle I \rangle = |E_0|^2 M_2$, increased as, $\langle I \rangle \propto \kappa^{-1}$, in agreement with Eq. (205). Thus, the field fluctuations lead to enhanced light scattering from the film.

It should be pointed out that the fluctuations considered here, and the corresponding light scattering, do not arise because of the fractal morphology of the metal clusters, but are due to the distribution of local resonances in a disordered metal-dielectric film, which is homogeneous on a macroscopic scale. The local intensity of the electric field is strongly correlated in space, and the distribution is dominated by the field correlation length $\xi_e$ introduced by Eqs. (199) and (212), and defined as the length scale over which the field fluctuations are small. As the $L-$ (metallic) component becomes loss-free ($\kappa \to 0$), $\xi_e$ diverges according to

$$\xi_e \sim \kappa^{-\nu_e}, \tag{223}$$

where $\nu_e$ is a new critical exponent which has been estimated by several numerical methods. For example, in 2D the PSRG method described above yields $\nu_e \simeq 0.45 \pm 0.05$, while the scaling theory, Eq. (208), predicts that $\nu_e = 1/d$, where $d$ is the space dimension, a result that was also conjectured by Hesselbo (1994). For small loses at resonance, the correlation length $\xi_e$ is the only relevant length scale of the system at $p_c$ since $|\epsilon_m|/\epsilon_d \simeq 1$.

### 3.5.2.1   Rayleigh Scattering

We consider now Rayleigh scattering induced by the giant field fluctuations (Brouers *et al.*, 1998) discussed above. Suppose that a semi-continuous film is illuminated by a wave normal to the film plane. The space between the metal grains is filled by a dielectric material. Therefore, the film can be considered as a 2D array of metal and dielectric grains that are distributed over the film's plane. The incident electromagnetic wave excites the surface current $\mathbf{I}$ in the film. Consider the electromagnetic field induced by these currents at some distant point $\mathbf{R}$. The origin of the coordinates is fixed at some point in the film. Then, the vector potential $\mathbf{A}(\mathbf{R})$ of the scattered field defined by, $\mathbf{H}(\mathbf{R}) = \nabla \times \mathbf{A}(\mathbf{R})$ [where $\mathbf{H}(\mathbf{R})$ is the magnetic field], arising from the surface current $\mathbf{I}(\mathbf{r})$, is such that if

$$\mathbf{A}(\mathbf{R}, \mathbf{r}) \, d\mathbf{r} = \frac{\mathbf{I}(\mathbf{r})}{c} \frac{\exp(ik \, |\mathbf{R} - \mathbf{r}|)}{|\mathbf{R} - \mathbf{r}|} \, d\mathbf{r}, \tag{224}$$

where $k = \omega/c$ is a wavevector, then $\mathbf{A}(\mathbf{R}) = \int \mathbf{A}(\mathbf{R}, \mathbf{r}) \, d\mathbf{r}$, where the integration is over the entire film. In experiments, the dimensions of the film are small enough that $r \ll R$, and therefore, $ik \, |\mathbf{R} - \mathbf{r}| \simeq ikR - ik(\mathbf{n} \cdot \mathbf{r})$, where $\mathbf{n}$ is the unit vector

in the direction of **R**. Thus,

$$\mathbf{H}(\mathbf{R}) \simeq \frac{ik}{cR} \exp(ikR) \int [\mathbf{n} \times \mathbf{I}(\mathbf{r})] \exp[-ik(\mathbf{n} \cdot \mathbf{r})] \, d\mathbf{r}, \qquad (225)$$

and the electric field is given by

$$\mathbf{E}(\mathbf{R}) = \frac{i}{k}[\nabla \times \mathbf{H}(\mathbf{R})] \simeq \frac{-ik}{cR} \exp(ikR) \int \{\mathbf{n} \times [\mathbf{n} \times \mathbf{I}(\mathbf{r})]\} \exp[-ik(\mathbf{n} \cdot \mathbf{r})] \, d\mathbf{r}. \tag{226}$$

It follows from Eqs. (225) and (226) that $\mathbf{H}(\mathbf{R})$ is perpendicular to $\mathbf{E}(\mathbf{R})$, and that $|\mathbf{E}(\mathbf{R})| = |\mathbf{H}(\mathbf{R})|$, implying that the scattered field can be considered locally as a plane wave when the distance from the film is large. The total intensity $I_t$ of the light scattered in the direction $\mathbf{n} = \mathbf{R}/|\mathbf{R}|$ is given by

$$I_t(\mathbf{n}) = \left(\frac{c}{4\pi}R^2\right)\frac{1}{2}\mathrm{Re}\{\langle[\mathbf{E}(\mathbf{R}) \times \mathbf{H}^*(\mathbf{R})]\rangle\} = \frac{c}{8\pi}R^2\langle\mathbf{E}(\mathbf{R}) \cdot \mathbf{E}^*(\mathbf{R})\rangle$$

$$= \frac{c}{8\pi}\frac{k^2}{c^2} \int \langle[\mathbf{n} \times \mathbf{I}(\mathbf{r}_1)] \cdot [\mathbf{n} \times \mathbf{I}^*(\mathbf{r}_2)]\rangle \exp[ik\mathbf{n} \cdot (\mathbf{r}_1 - \mathbf{r}_2)] \, d\mathbf{r}_1 d\mathbf{r}_2, \qquad (227)$$

where the angular brackets indicate an ensemble averaging. The semi-continuous metal films that are considered here are much larger than any characteristic intrinsic spatial scale, such as the field correlation length $\xi_e$, and therefore the ensemble average can be included in the integrations over the film area in Eq. (227) without changing the result. It is assumed, for simplicity, that the incident light is natural (unpolarized), and that its direction is perpendicular to the film plane. Then, the averaging $\langle[\mathbf{n} \times \mathbf{j}(\mathbf{r}_1)] \cdot [\mathbf{n} \times \mathbf{j}^*(\mathbf{r}_2)]\rangle$ should be carried out over the polarizations of the incident wave, yielding, $\langle\mathbf{I}(\mathbf{r}_1) \cdot \mathbf{I}^*(\mathbf{r}_2)\rangle[1 - \sin^2(\theta/2)]$, where $\theta$ is the angle between $\mathbf{n}$ and the normal to the film plane.

If we replace in Eq. (227) the local currents $\mathbf{I}(\mathbf{r})$ by their average values $\langle\mathbf{I}(\mathbf{r})\rangle$, we obtain the specular scattering $I_s$. The scattering $I(\theta) = I_t - I_s$ in all other directions is then obtained as

$$I(\theta) = \frac{c}{8\pi}\frac{k^2}{c^2}\left(1 - \frac{1}{2}\sin^2\theta\right) \int \left[\langle\mathbf{I}(\mathbf{r}_1) \cdot \mathbf{I}^*(\mathbf{r}_2)\rangle - |\langle\mathbf{I}\rangle|^2\right] \exp[ik\mathbf{n} \cdot (\mathbf{r}_1 - \mathbf{r}_2)] \, d\mathbf{r}_1 \, d\mathbf{r}_2. \tag{228}$$

The natural correlation length for the local field fluctuations, and therefore for the current-current correlations, is $\xi_e$. If $\xi_e \ll \lambda$, where $\lambda = 2\pi/k$ is the wavelength of the incident light, Eq. (228) is simplified by replacing the exponential by unity, hence yielding

$$I(\theta) = \frac{c}{8\pi}\frac{k^2}{c^2}\left(1 - \frac{1}{2}\sin^2\theta\right) |\langle\mathbf{I}\rangle|^2 \int \left[\frac{\langle\mathbf{I}(\mathbf{r}_1) \cdot \mathbf{I}^*(\mathbf{r}_2)\rangle}{|\langle\mathbf{I}\rangle|^2} - 1\right] d\mathbf{r}_1 \, d\mathbf{r}_2. \qquad (229)$$

Note that for macroscopically-homogeneous and isotropic films the current-current correlations $\langle\mathbf{I}(\mathbf{r}_1) \cdot \mathbf{I}^*(\mathbf{r}_2)\rangle$ depend only on $r = |\mathbf{r}_2 - \mathbf{r}_1|$. We now introduce the correlation function

$$C(r) = \frac{\langle\mathbf{I}(\mathbf{r}_1) \cdot \mathbf{I}^*(\mathbf{r}_2)\rangle}{|\langle\mathbf{I}\rangle|^2} - 1 = \frac{\mathrm{Re}\langle\mathbf{I}(0) \cdot \mathbf{I}^*(\mathbf{r})\rangle}{|\langle\mathbf{I}\rangle|^2} - 1. \tag{230}$$

in terms of which the intensity of the scattered light is given by

$$I(\theta) = A \frac{c}{8\pi} \frac{k^2}{c^2} \left(1 - \frac{1}{2}\sin^2\theta\right) |\langle \mathbf{I}\rangle|^2 2\pi \int_0^\infty C(r)r\,dr, \qquad (231)$$

where $A$ is the film area. $I(\theta)$ should be compared with the integral intensity (power) of the incident light, $I_0 = A(c/8\pi)|\mathbf{E}_0|^2$, where $|\mathbf{E}_0|$ is the amplitude of the incident wave. For the normal incident light, $\langle \mathbf{E}\rangle = \mathcal{T}\mathbf{E}_0$, where $\mathcal{T}$ is the transmittance of the film. For semi-continuous metallic films at $p = p_c$ one has $|\mathcal{T}|^2 \simeq 0.25$ in a wide spectral range from the visible to the far infrared spectral range (Yagil *et al.*, 1992). One also has, $\langle \mathbf{I}\rangle = ag_e\langle \mathbf{E}\rangle = ag_e\mathcal{T}\mathbf{E}_0$, where $g_e = -i\epsilon_e\omega/(4\pi)$ is the effective conductivity, and thickness of the film has been approximated by the size $a$ of a metal grain.

Substituting $\langle \mathbf{I}\rangle = g_e a\mathcal{T}\mathbf{E}_0$ in Eq. (231), the ratio, $\tilde{I}(\theta) = I(\theta)/I_0$ is obtained,

$$\tilde{I}(\theta) = \frac{(ka)^4}{8\pi} \left(1 - \frac{1}{2}\sin^2\theta\right) |\mathcal{T}\epsilon_e|^2 \frac{1}{a^2}\int_0^\infty C(r)r\,dr, \qquad (232)$$

which is independent of the film's geometry. It follows from Eq. (232) that the portion of the incident light that is not reflected, transmitted or adsorbed, but is scattered from the film is given by

$$I_{tot} = 2\pi\int \tilde{I}(\theta)\,\sin\theta\,d\theta = \frac{(ka)^4}{3}|\mathcal{T}\epsilon_e|^2\frac{1}{a^2}\int_0^\infty C(r)r\,dr. \qquad (233)$$

The behavior of $C(r)$ depends on the frequency, and also on the behavior of $|\mathcal{T}\epsilon_e|^2$, which achieves large values, $|\mathcal{T}\epsilon_e|^2 \gg 1$, in the infrared spectral range.

We can compare Eq. (233) with the scattering for the case when the metal grains interact with the electromagnetic field independently. The cross section $\sigma_R$ of Rayleigh scattering from a single metal grain is estimated as, $\sigma_R = (8\pi/3)(ka)^4a^2$ for $|\epsilon_m| \gg 1$. The portion of the light which would be scattered if the grains were independent is given by $S_{tot}^R \simeq p(8/3)(ka)^4$. Assuming $p = 1/2$, the following estimate is obtained for the enhancement $I_g = I_{tot}/S_{tot}^R$ of the scattering due to the field fluctuations,

$$I_g \sim \frac{|\mathcal{T}\epsilon_e|^2}{4a^2}\int_0^\infty C(r)r\,dr. \qquad (234)$$

If the integral in (234) is determined by the largest distances where field correlations are most important, i.e., where, $r \sim \xi_e$, the scattering can even diverge if losses vanish and $\xi_e \to \infty$. This is certainly the case for 2D metal-dielectric films. The above formalism holds if $I_{tot} \ll 1$. Otherwise, it is necessary to take into account the feedback effects, i.e., the interaction of the scattered light with the film.

### 3.5.2.2   Scaling Properties of the Correlation Function

Using the PSRG approach described above, Brouers *et al.* (1998) calculated the correlation function $C(r)$ for a $1024 \times 1024$ $L - C$ system for gold semi-continuous metal films at $p = p_c = 1/2$ and for the resonance frequency $\omega_r$, so

that $\epsilon'_m(\omega_r) = -\epsilon_d$, and for several values of $\kappa = \epsilon''_m/|\epsilon'_m|$. For each $\kappa$, the results were averaged over 100 different realizations of the system. Their calculations indicated that for $a < r < \xi_e$, the correlation function decays as (the distance $r$ is measured in units of the metal grain size $a$)

$$C^\star(r) \sim M_2^\star (r/a)^{-(1+\eta)} \sim \kappa^{-1} (r/a)^{-(1+\eta)} , \qquad (235)$$

where $M_2^\star$ [see Eq. (197)] is the second moment of the local field in the system with $\epsilon'_m = -\epsilon_d$, and $\eta = 0.8 \pm 0.1$ is a new critical exponent that determines the spatial correlation of the local electric field. If we substitute Eq. (235) in (232) and (233), we find that the integrals diverge at the upper limit, implying that the scattering is determined by values of the correlation function $C(r)$ at large distances, i.e., at $r \sim \xi_e$. This means that the field fluctuations with spatial distances of the field correlation length $\xi_e \gg 1$ are responsible for the anomalous scattering from semi-continuous films.

We now consider the dependence of scattering on the frequency of an incident electromagnetic wave. We first consider frequencies just below $\tilde{\omega}_p$ where the metal dielectric function can be estimated for a Drude model as, $\epsilon_m = \epsilon'_m + i\epsilon''_m \simeq 2\epsilon_b(\omega - \tilde{\omega}_p)/\tilde{\omega}_p + i\epsilon_b\omega_\tau/\tilde{\omega}_p$, i.e., $\epsilon'_m < 0$. For such frequencies, $|\epsilon'_m|/\epsilon_d \leq 1$, while the loss factor $\kappa \simeq \omega_\tau/2(\tilde{\omega}_p - \omega)$ decreases rapidly with frequency $\omega$, and in particular decreases below the renormalized plasma frequency $\tilde{\omega}_p$. For $|\epsilon'_m|/\epsilon_d \simeq 1$, the correlation function $C(r)$ is estimated by Eq. (235) which, when substituted in Eq. (233) and integrated up to $\xi_e \sim a\kappa^{-1/d}$, yields

$$I_{tot} \sim \frac{(ka)^4}{3}|\mathcal{T}|^2\epsilon_d\epsilon_b \left(\frac{\tilde{\omega}_p}{\omega_\tau}\right)^{1+(1-\eta)/d} \left(1 - \frac{\omega}{\tilde{\omega}_p}\right)^{2+(1-\eta)/d} , \qquad \omega < \tilde{\omega}_p,$$
$$(236)$$

where the exact result, $\epsilon_e(p = p_c) = \sqrt{\epsilon_d\epsilon_m}$, which is a result of the duality relation for 2D percolation systems (see above and also Chapters 4 and 5 of Volume I), was used.

Consider now the limit $\omega \ll \tilde{\omega}_p$, assuming again that $\omega \gg \omega_\tau$, for which the dielectric constant for a Drude metal is approximated as, $\epsilon_m \simeq (\omega_p/\omega)^2(-1 + i\omega_\tau/\omega)$, yielding $|\epsilon'_m|/\epsilon_d \simeq (\omega_p/\omega)^2/\epsilon_d \gg 1$ and $\kappa = \epsilon''_m/|\epsilon'_m| \simeq \omega/\omega_\tau \ll 1$. To estimate the correlation function $C(r)$, the system is divided into squares of size $b$ and the procedure described above is followed by taking $b = b_r$, where $b_r$ is given by Eq. (206). Then, the correlation function $C^\star$ in the renormalized system has the same form as Eq. (235), while in the original system, $C(r) \simeq (b_r/a)^{1+\eta}C^\star(r)$ for $r \gg b_r$, and $C(r) \propto r^{-\mu/\nu}$ for $r \ll b_r$. By matching these asymptotic expressions at $r = b_r$, the following *ansatz* emerges,

$$C(r) \sim \begin{cases} M_2^\star(b_r/r)^{\mu/\nu} \sim \kappa^{-1}(b_r/r)^{\mu/\nu}, & a \ll r < b_r, \\ M_2^\star(b_r/r)^{1+\eta} \sim \kappa^{-1}(b_r/r)^{1+\eta}, & b_r < r < \xi_e, \end{cases} \qquad (237)$$

where $\xi_e$ is given by Eq. (212). Equation (237) allows one to estimate the second moment of the local electric current, $M_j \equiv \langle|\mathbf{I}(\mathbf{r})|^2\rangle$, at $p_c$. From Eq. (230), one can write, $M_j = |\mathbf{E}_0|^2|g_e|^2C(0) = (\omega/4\pi)^2|\mathbf{E}_0|^2|\epsilon_e|^2C(0)$. At $p_c$ one has $\epsilon_e \sim \epsilon_d(\epsilon_m/\epsilon_d)^{s/(s+\mu)}$ (see Chapter 5 of Volume I). The correlation function $C(r)$

for $r \sim a$ is given by, $C(0) \sim C(a) \sim M_2^\star (l_r/a)^{\mu/\nu} \sim M_2^\star (|\epsilon_m|/\epsilon_d)^{\mu/(s+\mu)}$, and hence

$$M_j \sim (\omega/4\pi)^2 |\mathbf{E}_0|^2 \epsilon_d^2 M_2^\star \left( \frac{|\epsilon_m|}{\epsilon_d} \right)^{(2s+\mu)/(s+\mu)} \sim (\omega/4\pi)^2 |\mathbf{E}_0|^2 \epsilon_d |\epsilon_m| M_2, \tag{238}$$

where $M_2^\star \equiv M_{2,0}^\star$ and $M_2 \equiv M_{2,0}$, as defined earlier. Equation (238) holds for arbitrary spatial dimension.

We now consider light scattering from semi-continuous metal films for $\omega \ll \omega_p$, where the metal dielectric constant for a Drude metal is approximated as, $\epsilon_m \simeq -(\omega_p/\omega)^2(1 - i\omega_\tau/\omega)$. By substituting Eq. (237) into (233) and taking into account the fact that at $p_c$, $|\epsilon_e|^2 \simeq \epsilon_d |\epsilon_m| \simeq \epsilon_d (\omega_p/\omega)^2$ (using $\mu = s \simeq \nu = 4/3$), the following result is obtained

$$I_{tot} \sim \frac{(ka)^4}{3} |\mathcal{T}|^2 |\epsilon_e|^2 \kappa^{-1} b_r^{1+\eta} \xi_e^{1-\eta}$$

$$\sim \frac{(ka)^4}{3} |\mathcal{T}|^2 \kappa^{-1-(1-\eta)/2} |\epsilon_m|^2 \sim 0.1 \left( \frac{\omega_p a}{c} \right)^4 \left( \frac{\omega}{\omega_\tau} \right)^{1+(1-\eta)/2}, \tag{239}$$

where the experimental result, $|\mathcal{T}|^2 \simeq 0.25$, which holds for $p = p_c$ and $\omega_\tau \ll \omega \ll \omega_p$, was used. Thus, the scattering first increases as $\omega^{1+(1-\eta)/2}$ with increasing $\omega$ according to Eq. (239) and then vanishes as $(\tilde{\omega}_p - \omega)^{2+(1-\eta)/2}$ as $\omega \to \tilde{\omega}_p$ [see Eq. (236)].

The enhancement of the scattering due to the field fluctuations can be estimated from Eqs. (234) and (237) as, $I_g \sim |\mathcal{T}|^2 \epsilon_d |\epsilon_m| b_r^2 \kappa^{-1-(1-\eta)/2}/4$, which yields, for a Drude metal and $\omega \ll \omega_p$, the following equation

$$I_g \sim \frac{|\mathcal{T}|^2}{4} \left( \frac{\tilde{\omega}_p}{\omega} \right)^4 \left( \frac{\omega}{\omega_\tau} \right)^{1+(1-\eta)/2}. \tag{240}$$

Using typical values, $|\mathcal{T}|^2 = 1/4$ and $\epsilon_d = 2.2$, the enhancement $I_g$ can become as large as $5 \times 10^4$ at wavelength $\lambda = 1.5$ $\mu$m and continues to increase towards the far infrared spectral range. Note that Rayleigh scattering decreases as $\omega^4$ with decreasing frequency, whereas the anomalous scattering varies as, $I \sim \omega^{1+(1-\eta)/2} \simeq \omega^{1.1}$, and therefore the enhancement increases as $I_g \sim \omega^{-2.9} \sim \lambda^{2.9}$ in the infrared part of the spectrum.

### 3.5.3  Surface-Enhanced Raman Scattering

We now consider surface-enhanced Raman scattering (SERS), one of the most intriguing optical effects discovered over the past 20 years (Moskovits, 1985; Markel *et al.*, 1996; Kneipp *et al.*, 1997; Nie and Emory, 1997), and describe a theory of Raman scattering (see also Chapter 6 of Volume I) enhanced by strong fluctuations of the local fields (Brouers *et al.*, 1997). In rough thin films this phenomenon is commonly associated with excitation of surface plasmon oscillations which are typically considered in two limiting cases: (1) Oscillations in non-interacting

roughness features of various shapes, and (2) surface plasmon waves (polaritons) that laterally propagate along the metal surface. In practice, there are strong light-induced interactions between different features of a rough surface, and therefore plasmon oscillations should be treated as collective surface excitations (localized surface plasmons) that depend strongly on the surface morphology.

### 3.5.3.1  General Formulation

The formulation of the problem and the solution that are discussed here are due to Brouers *et al.* (1997), as described by Sarychev and Shalaev (2000). We consider optical properties of a semi-continuous metal film consisting of metal grains, randomly distributed on a dielectric substrate. The space between the metal grains are usually filled by dielectric material of the substrate. As before, the local conductivity $g(\mathbf{r})$ of the film takes on either the metallic value, $g(\mathbf{r}) = g_m$, in the metal grains, or the dielectric value, $g(\mathbf{r}) = -i\omega\epsilon_d/4\pi$, outside the metal grains, where $\omega$ is the frequency of the external field. We assume that the wavelength $\lambda$ is much larger than the grain size $a$, the linear size of the space between the grains, percolation correlation length $\xi_p$, and the local field correlation length. Hence, the local field $\mathbf{E}(\mathbf{r})$ is given by Eq. (179).

It is instructive to assume first that the external field $\mathbf{E}_0(\mathbf{r})$ is step-like, $\mathbf{E}_0(\mathbf{r}) = \mathbf{E}_1\delta(\mathbf{r} - \mathbf{r}_1)$, where $\delta(\mathbf{r})$ is the Dirac delta-function. The current density at an arbitrary point $\mathbf{r}_2$ is then given by

$$\mathbf{I}(\mathbf{r}_1, \mathbf{r}_2) = \mathbf{\Sigma}(\mathbf{r}_2, \mathbf{r}_1)\mathbf{E}_1, \tag{241}$$

where $\mathbf{\Sigma}(\mathbf{r}_2, \mathbf{r}_1)$ is the non-local conductivity matrix representing the system's response at point $\mathbf{r}_2$ to a source at the point $\mathbf{r}_1$, such that if an external field $\mathbf{E}_0(\mathbf{r})$ is applied to the system, the local current at the point $\mathbf{r}_2$ will be given by

$$\mathbf{I}(\mathbf{r}_2) = \int \mathbf{\Sigma}(\mathbf{r}_2, \mathbf{r}_1)\mathbf{E}_0(\mathbf{r}_1)\, d\mathbf{r}_1, \tag{242}$$

where the integration is over the total area of the system.

In view of our discussion in Chapters 5 and 6 of Volume I, it should be clear that $\mathbf{\Sigma}$ can be expressed in terms of the Green function $G$ of Eq. (179):

$$\nabla \cdot \{g(\mathbf{r}_2)[\nabla G(\mathbf{r}_2, \mathbf{r}_1)]\} = \delta(\mathbf{r}_2 - \mathbf{r}_1), \tag{243}$$

where a differentiation with respect to the coordinate $\mathbf{r}_2$ is assumed. Comparing Eqs. (179) and (243), the following equation for the element of $\mathbf{\Sigma}$ is obtained

$$\Sigma_{\alpha\beta}(\mathbf{r}_2, \mathbf{r}_1) = g(\mathbf{r}_2)g(\mathbf{r}_1)\frac{\partial^2 G(\mathbf{r}_2, \mathbf{r}_1)}{\partial r_{2,\alpha}\,\partial r_{1,\beta}}, \tag{244}$$

where the Greek indices denote $x$ and $y$. It is clear that, because of the symmetry of the Green function,

$$\Sigma_{\alpha\beta}(\mathbf{r}_1, \mathbf{r}_2) = \Sigma_{\beta\alpha}(\mathbf{r}_2, \mathbf{r}_1). \tag{245}$$

Since we assumed that the wavelength of the incident electromagnetic wave is much larger than all spatial scales in a semi-continuous metal film, the external

field $\mathbf{E}_0$ is constant in the film plane. The local field $\mathbf{E}(\mathbf{r}_2)$, induced by the external field $\mathbf{E}_0$, is obtained by using Eq. (144) for the non-local conductivity $\Sigma$,

$$\mathbf{E}(\mathbf{r}_2) = \frac{1}{g(\mathbf{r}_2)} \int \Sigma(\mathbf{r}_2, \mathbf{r}_1)\mathbf{E}_0 \, d\mathbf{r}_1, \tag{246}$$

and excite Raman-active molecules that are (assumed to be) uniformly distributed in the composite. Such molecules, in turn, generate the Stokes fields, $\mathbf{E}_s(\mathbf{r}_2) = \alpha_s(\mathbf{r}_2)\mathbf{E}(\mathbf{r}_2)$, oscillating at the shifted frequency $\omega_s$, where $\alpha_s(\mathbf{r}_2)$ is the ratio of the Raman and linear polarizabilities of the Raman-active molecules at $\mathbf{r}_2$. The Stokes fields $\mathbf{E}_s(\mathbf{r}_2)$ induce in the composite the currents $\mathbf{I}_s(\mathbf{r}_3)$ that are given by

$$\mathbf{I}_s(\mathbf{r}_3) = \int \Sigma(\mathbf{r}_3, \mathbf{r}_2)\mathbf{E}_s(\mathbf{r}_2) \, d\mathbf{r}_2. \tag{247}$$

Since the frequency $\omega_s$ is typically close to the external field's frequency, i.e., $|\omega - \omega_s|/\omega \ll 1$, the non-local conductivities $\Sigma$ appearing in Eqs. (246) and (247) are essentially the same.

The intensity $I$ of the electromagnetic wave scattered from any inhomogeneous material is proportional to the current fluctuations inside the system:

$$I \propto \left\langle \left| \int [\mathbf{I}(\mathbf{r}) - \langle \mathbf{I} \rangle] \, d\mathbf{r} \right|^2 \right\rangle, \tag{248}$$

where the integration is over the entire system, and $\langle \cdot \rangle$ denotes an ensemble average. For Raman scattering, $\langle \cdot \rangle$ also includes averaging over the fluctuating phases of the incoherent Stokes fields generated by Raman-active molecules. Therefore, the average current densities oscillating at $\omega_s$ is zero, $\langle \mathbf{I}_s \rangle = 0$, and hence the intensity $I_R$ of Raman scattering from a semi-continuous metal film is given by

$$I_R \propto \left\langle \left| \int \mathbf{I}(\mathbf{r}) \, d\mathbf{r} \right|^2 \right\rangle$$

$$= \int \left\langle \Sigma_{\alpha\beta}(\mathbf{r}_3, \mathbf{r}_2)\alpha_s(\mathbf{r}_2) E_\beta(\mathbf{r}_2) \Sigma^*_{\alpha\gamma}(\mathbf{r}_5, \mathbf{r}_4)\alpha^*_s(\mathbf{r}_4) E^*_\gamma(\mathbf{r}_4) \right\rangle \, d\mathbf{r}_2 \, d\mathbf{r}_3 \, d\mathbf{r}_4 \, d\mathbf{r}_5 \tag{249}$$

where a summation over repeating Greek indices is implied, and the integration is over the entire film plane. Equation (249) is now averaged over the fluctuating phases of the Raman polarizabilities $\alpha_s$. Because the Raman field sources are incoherent, we have $\langle \alpha_s(\mathbf{r}_2)\alpha^*_s(\mathbf{r}_4) \rangle = |\alpha_s|^2\delta(\mathbf{r}_2 - \mathbf{r}_4)$, and therefore

$$I_R \propto \int \left\langle \Sigma_{\alpha\beta}(\mathbf{r}_3, \mathbf{r}_2)\Sigma^*_{\mu\gamma}(\mathbf{r}_5, \mathbf{r}_2)\delta_{\alpha\mu}|\alpha_s|^2 E_\beta(\mathbf{r}_2) E^*_\gamma(\mathbf{r}_2) \right\rangle \, d\mathbf{r}_2 \, d\mathbf{r}_3 \, d\mathbf{r}_5. \tag{250}$$

If we now take advantage of the facts that, (1) a semi-continuous film is macroscopically homogeneous, and thus its Raman scattering is independent of the orientation of the external field $\mathbf{E}_0$; (2) due to (1), Eq. (250) can be averaged over the orientations of $\mathbf{E}_0$ without changing the result, and (3) the non-local conductivity $\Sigma$ is

independent of the field orientations and is symmetric, we obtain

$$\langle I_R \rangle \propto \frac{|\alpha_s|^2}{|E_0|^2} \int |g(\mathbf{r}_2)|^2 \, \langle |E(\mathbf{r}_2)|^2 \rangle_0 \langle |E(\mathbf{r}_2)|^2 \rangle_0 \, d\mathbf{r}_2, \tag{251}$$

where $\langle \cdot \rangle_0$ denotes the orientation averaging. It is not difficult to show that, for macroscopically-isotropic materials, Eq. (251) can be rewritten as

$$\langle I_R \rangle \propto \frac{|\alpha_s|^2}{|E_0|^2} \int |g(\mathbf{r}_2)|^2 |E(\mathbf{r}_2)|^4 \, d\mathbf{r}_2. \tag{252}$$

In the absence of any metal grains on the film, the local fields would not fluctuate and one would obtain

$$I_R^0 \propto \int |g_d|^2 |\alpha_s|^2 |E_0|^2 \, d\mathbf{r}_2. \tag{253}$$

Therefore, the enhancement $I_{RS} = I_R/I_R^0$ of Raman scattering due to presence of metal grains on a dielectric substrate is given by

$$I_{RS} = \frac{\langle |g(\mathbf{r})|^2 |E(\mathbf{r})|^4 \rangle}{|g_d|^2 |E_0|^4} = \frac{\langle |\epsilon(\mathbf{r})|^2 |E(\mathbf{r})|^4 \rangle}{\epsilon_d^2 |E_0|^4}. \tag{254}$$

Note that the derivation of Eq. (254) is essentially independent of the dimensionality and morphology of the material. Therefore, the enhancement $I_{RS}$ should hold for *any* heterogeneous material, provided that the field fluctuations take place inside of it. In particular, Eq. (254) yields the enhancement for Raman scattering from a rough metallic surface, provided that the wavelength is much larger than the roughness spatial scales. It can also be used for calculating the enhancements in a 3D percolation composite. The present theory indicates also that the main source for the Raman scattering is the currents excited by Raman molecules in metal grains, hence explaining why a large $I_{RS}$ is obtained even for relatively flat metal surfaces (Moskovitz, 1985).

### 3.5.3.2  Raman and Hyper-Raman Scattering in Metal–Dielectric Composites

Since, as discussed above, the local electric field in materials with percolation disorder is distributed mainly in the dielectric space between the metal clusters, the SERS enhancement $I_{RS}$ may be estimated as, $I_{RS} \sim M_{4,0} = \langle |E(\mathbf{r})/E_0|^4 \rangle$. Hence, in view of Eq. (215), we obtain

$$I_{RS} \sim \mathcal{N}(p)[\xi_A(p)/a]^{d-8} \left( \frac{|\epsilon_m|}{\epsilon_d} \right)^{(2\nu+s)/(\mu+s)} \left( \frac{|\epsilon_m|}{\epsilon_m''} \right)^3, \tag{255}$$

indicating that, when the states are delocalized, $\xi_A \to \infty$, $I_{RS}$ vanishes very rapidly. Equation (255) can now be used for investigating the frequency and volume fraction dependence of Raman scattering. For 2D metal-dielectric composites with the critical exponents, $\mu = s \simeq \nu = 4/3$, the Drude metal dielectric function can be used for frequencies $\omega \ll \omega_p$, and therefore Eq. (255) predicts

that, $I_{RS} \sim \mathcal{N}(p)[a/\xi_A(p)]^6(\omega_p/\omega_\tau)^3/\epsilon_d^{3/2}$, independent of the frequency. For example, for silver-on-glass percolation films at $p_c$, the Anderson localization length $\xi_A$ is about $\xi_A \simeq 2a$, the density of state, $\mathcal{N}(p_c) \simeq 1$, and therefore, $I_{RS} \sim 10^6$. For 3D composites at $\omega \ll \omega_p$, $I_{RS}$ decreases with decreasing $\omega$ as $I_{RS} \sim \mathcal{N}(p)(\xi_A/a)^{-5}\omega_p^2\omega/\omega_\tau^3 \sim 10^6\omega/\omega_p$, where the 3D critical exponents have been approximated as, $\nu \simeq s \simeq (\mu + s)/3$, and the data, $\omega_p = 9.1$ eV and $\omega_\tau = 0.021$ eV, for silver dielectric constant have been utilized.

Consider now hyper-Raman scattering when $n$ photons of frequency $\omega$ are converted to one hyper-Stokes photon of the frequency $\omega_{hRS} = n\omega - \omega_{sf}$, where $\omega_{sf}$ is the Stokes frequency shift corresponding to the frequency of molecule oscillations (electronic or vibrational). Thus, following the same line of reasoning outlined above, the surface enhancement of hyper-Raman scattering (SEHRS) $I_{hRS}$ is given by

$$I_{hRS} = \frac{\langle|g_{hRS}(\mathbf{r})|^2|\mathbf{E}_{hRS}(\mathbf{r})|^2|E(\mathbf{r})|^{2n}\rangle}{|g_d|^2\left|E_{0,hRS}\right|^2|E_0|^{2n}} = \frac{\langle|\epsilon_{hRS}(\mathbf{r})|^2|\mathbf{E}_{hRS}(\mathbf{r})|^2|E(\mathbf{r})|^{2n}\rangle}{|\epsilon_d|^2\left|E_{0,hRS}\right|^2|E_0|^{2n}},$$
(256)

where $\mathbf{E}_{hRS}(\mathbf{r})$ is the local field excited in the system by the uniform probe field $\mathbf{E}_{0,hRS}$, oscillating with $\omega_{hRS}$, and $g_{hRS}(\mathbf{r})$ and $\epsilon_{hRS}(\mathbf{r})$ are the local conductivity and dielectric constant at frequency $\omega_{hRS}$. For $n = 1$ Eq. (256) describes the conventional SERS. To estimate $I_{hRS}$, we must keep in mind that the spatial scales $b_r$ for the field maxima at the fundamental frequency $\omega$ and the hyper-Stokes frequency $\omega_{hRS}$ are significantly different. Therefore, the average in Eq. (256) can be decoupled and approximated as, $\langle|\epsilon_{hRS}(\mathbf{r})|^2|\mathbf{E}_{hRS}(\mathbf{r})|^2|E(\mathbf{r})|^{2n}\rangle \sim \langle|\epsilon_{hRS}(\mathbf{r})\mathbf{E}_{hRS}(\mathbf{r})|^2\rangle\langle|E(\mathbf{r})|^{2n}\rangle = \langle|\epsilon_{hRS}(\mathbf{r})\mathbf{E}_{hRS}(\mathbf{r})|^2\rangle M_{2n}|E_0|^{2n}$, where $M_{2n}(\omega)$ is the $2n$th moment. It follows from Eq. (238) that, $\langle|\epsilon_{hRS}(\mathbf{r})\mathbf{E}_{hRS}(\mathbf{r})|^2\rangle \sim \epsilon_d|\epsilon_m(\omega_{hRS})|M_2|E_{0,hRS}|^2$, where $M_2(\omega_{hRS})$ is the second moment of the field $\mathbf{E}_{hRS}(\mathbf{r})$. Using the expressions for $M_2$ and $M_{2n}$ given above, and taking into account the fact that for $p \simeq p_c$ the density of states $\mathcal{N}$ is about unity, one obtains the following equation for enhancement of hyper-Raman scattering,

$$I_{hRS} \sim (\xi_A/a)^{2d-4(1+n)}\left(\frac{|\epsilon_m(\omega_{hRS})|}{\epsilon_d}\right)^{(\mu+2s)/(\mu+s)}\left(\frac{|\epsilon_m(\omega_{hRS})|}{\epsilon_m''(\omega_{hRS})}\right)$$
$$\times\left(\frac{|\epsilon_m(\omega)|}{\epsilon_d}\right)^{[2\nu(n-1)+s]/(\mu+s)}\left(\frac{|\epsilon_m(\omega)|}{\epsilon_m''(\omega)}\right)^{2n-1},$$
(257)

with $n \geq 2$. For a Drude metal and frequencies $\omega \ll \tilde{\omega}_p$, $\omega_{hRS} \ll \tilde{\omega}_p$ the metal dielectric constant can be approximated as, $|\epsilon_m(\omega_{hRS})| \sim |\epsilon_m(\omega)| \sim (\omega_p/\omega)^2$ and $\epsilon_m''(\omega)/|\epsilon_m(\omega)| \sim \omega_\tau/\omega$, and therefore Eq. (257) becomes

$$I_{hRS} \sim (\xi_A/a)^{2d-4(1+n)}\left(\frac{\omega_p}{\omega}\right)^{2[2\nu(n-1)+3s+\mu]/(\mu+s)}\left(\frac{\omega}{\omega_\tau}\right)^{2n},$$
(258)

which, in 2D (using $\mu = s \simeq \nu = 4/3$), simplifies to

$$I_{hRS} \sim (a/\xi_A)^{4n}\left(\frac{\omega_p}{\omega}\right)^{2(n+1)}\left(\frac{\omega}{\omega_\tau}\right)^{2n} \sim (a/\xi_A)^{4n}\left(\frac{\omega_p}{\omega}\right)^2\left(\frac{\omega_p}{\omega_\tau}\right)^{2n}.$$
(259)

FIGURE 3.12. Comparison of experimental data (points with error bars) for normalized SERS, $\bar{I} = I_{RS}(p)/I_{RS}(p = p_c)$, for a semicontinuous silver film, versus the theoretical computations (curve) (after Sarychev and Shalaev, 2000).

### 3.5.3.3   Comparison with the Experimental Data

As discussed earlier, the localization radius $\xi_A$ of the eigenstates $\Psi_n$ with eigenvalues $\Lambda \simeq 0$ decreases when one shifts from $p = p_c$ toward $p = 0$ or $p = 1$, because the eigenvalue $\Lambda = 0$ shifts from the center of the $\Lambda$-distribution to its tails, where localization of the eigenstates is stronger. Therefore, according to Eq. (255), Raman scattering must have a *minimum* at $p_c$, as a result of which $I_{RS}(p)$ must have *two* maxima, with one maximum below $p_c$ and a second one above $p_c$. Figure 3.12 presents (Gadenne *et al.*, 1998) experimental data for the dependence of SERS on the metal volume fraction $p$, and compares them with the theoretical predictions. It is clear that there is good qualitative agreement between the predictions and the data. In particular, in agreement with the theory, there is a minimum near $p_c$.

## 3.5.4   Enhancement of Optical Nonlinearities in Metal–Dielectric Composites

The next subject we consider is enhancement in heterogeneous materials with percolation-type disorder of various nonlinear optical processes, such as the Kerr optical effect and generation of high harmonics.

### 3.5.4.1   Kerr Optical Nonlinearities

These are third-order optical nonlinearities that result in an additional term in the electric displacement **D** given by

$$D_i^{(3)}(\omega) = \epsilon_{ijkl}^{(3)}(-\omega, \omega, \omega, -\omega) E_j E_k E_l^*, \tag{260}$$

where $\epsilon_{ijkl}^{(3)}(-\omega, \omega, \omega, -\omega)$ is the third-order nonlinear dielectric constant (see,

for example, Boyd, 1992), $\mathbf{E}$ is an electric field at frequency $\omega$, and summation over repeated indices is implied. The Kerr optical nonlinearity results in nonlinear corrections, which are proportional to the light intensity, for the refractive index and the absorption coefficient.

We consider disordered materials that are macroscopically homogeneous and isotropic. For such materials, the third-order term in the average electric displacement is given by

$$\left\langle \mathbf{D}^{(3)}(\mathbf{r}) \right\rangle = \alpha |\mathbf{E}_0|^2 \mathbf{E}_0 + \beta |\mathbf{E}_0|^2 \mathbf{E}_0^*, \tag{261}$$

where $|\mathbf{E}_0|$ is the amplitude of the external electric field at frequency $\omega$, and $\alpha$ and $\beta$ are some constants. Note that, for an isotropic film, the second term in Eq. (261) results in change of the polarization of the incident light. Moreover, for the case of linear and circular polarization of the incident light, Eq. (261) can be simplified since for linear polarization the complex vector $\mathbf{E}_0$ reduces to a real vector. Then, $|\mathbf{E}_0|^2 \mathbf{E}_0 = E_0^2 \mathbf{E}_0$, and Eq. (261) becomes

$$\left\langle \mathbf{D}^{(3)}(\mathbf{r}) \right\rangle = \epsilon_e^{(3)} |E_0|^2 \mathbf{E}_0, \tag{262}$$

where the effective nonlinear dielectric constant $\epsilon_e^{(3)}$ is now a scalar quantity. Let us consider, for the sake of simplicity, the linearly polarized incident wave. We write Eq. (262) in terms of the nonlinear average current $\langle \mathbf{I}^{(3)}(\mathbf{r}) \rangle$ and the effective Kerr conductivity $g_e^{(3)} = -i\omega\epsilon_e^{(3)}/4\pi$:

$$\left\langle \mathbf{I}^{(3)}(\mathbf{r}) \right\rangle = g_e^{(3)} |E_0|^2 \mathbf{E}_0. \tag{263}$$

We consider first the limit in which the nonlinearities in metal grains $g_m^{(3)}$ and dielectric $g_d^{(3)}$ are approximately equal, $g_m^{(3)} \simeq g_d^{(3)}$, which can be caused by, for example, molecules that are uniformly covering a semi-continuous film. Then

$$\mathbf{I}(\mathbf{r}) = g^{(\ell)}(\mathbf{r})\mathbf{E}'(\mathbf{r}) + g^{(3)} \left|\mathbf{E}'(\mathbf{r})\right|^2 \mathbf{E}'(\mathbf{r}), \tag{264}$$

where $\mathbf{E}'(\mathbf{r})$ is the local fluctuating field. Then, current conservation law takes the following form

$$\nabla \cdot \left\{ g^{(\ell)}(\mathbf{r}) \left[ -\nabla\phi(\mathbf{r}) + \mathbf{E}_0 + \frac{g^{(3)}}{g^{(\ell)}(\mathbf{r})} \mathbf{E}'(\mathbf{r}) \left|\mathbf{E}'(\mathbf{r})\right|^2 \right] \right\} = 0, \tag{265}$$

where $-\nabla\phi(\mathbf{r}) + \mathbf{E}_0 = \mathbf{E}'(\mathbf{r})$ is the local field. The second and third terms of Eq. (265) can be thought of as a renormalized external field

$$\mathbf{E}_e(\mathbf{r}) = \mathbf{E}_0 + \mathbf{E}_f(\mathbf{r}) = \mathbf{E}_0 + \frac{g^{(3)}}{g^{(\ell)}(\mathbf{r})} \mathbf{E}'(\mathbf{r}) \left|\mathbf{E}'(\mathbf{r})\right|^2, \tag{266}$$

where the field $\mathbf{E}_f(\mathbf{r})$ may change over the film but its average, $\langle \mathbf{E}_f(\mathbf{r}) \rangle$, is collinear to $\mathbf{E}_0$, in which case the average current density $\langle \mathbf{I}(\mathbf{r}) \rangle$ is also collinear to $\mathbf{E}_0$ and

can be written as

$$\langle \mathbf{I} \rangle = \frac{\mathbf{E}_0}{E_0^2}(\mathbf{E}_0 \cdot \langle \mathbf{I} \rangle) = \frac{\mathbf{E}_0}{E_0^2}\frac{1}{A}\int \mathbf{E}_0 \cdot \mathbf{I}(\mathbf{r})\, d\mathbf{r}, \qquad (267)$$

where $A$ is the total area of the film, the integration is over the film area, and $E_0^2 \equiv (\mathbf{E}_0 \cdot \mathbf{E}_0)$. Expressing $\mathbf{I}(\mathbf{r})$ in terms of the non-local conductivity matrix defined by Eq. (241) yields

$$\langle \mathbf{I} \rangle = \frac{\mathbf{E}_0}{E_0^2}\frac{1}{A}\int [\mathbf{E}_0 \boldsymbol{\Sigma}(\mathbf{r},\mathbf{r}_1)\mathbf{E}_e(\mathbf{r}_1)]\, d\mathbf{r}\, d\mathbf{r}_1. \qquad (268)$$

If we integrate Eq. (268) over the coordinates $\mathbf{r}$ and use the symmetry of the non-local conductivity matrix $\boldsymbol{\Sigma}$, we obtain

$$\langle \mathbf{I} \rangle = \frac{\mathbf{E}_0}{E_0^2}\frac{1}{A}\int [\mathbf{I}_0\mathbf{r} \cdot \mathbf{E}_e(\mathbf{r})]\, d\mathbf{r}, \qquad (269)$$

where $\mathbf{I}_0(\mathbf{r})$ is the current induced at $\mathbf{r}$ by the constant external field $\mathbf{E}_0$. Using Eq. (266) and carrying out the integration, Eq. (269) becomes

$$\langle \mathbf{I} \rangle = E_0\left[g_e^{(\ell)} + \frac{\left\langle g^{(3)}\left[\mathbf{E}(\mathbf{r}) \cdot \mathbf{E}'(\mathbf{r})\right]\left|\mathbf{E}'(\mathbf{r})\right|^2\right\rangle}{E_0^2}\right], \qquad (270)$$

where $g_e^{(\ell)}$ and $\mathbf{E}(\mathbf{r})$ are the effective conductivity and local fluctuating field in the linear approximation (i.e., for $g^{(3)} \equiv 0$). Comparison of Eqs. (270) and (263) yields an expression for the effective Kerr conductivity:

$$g_e^{(3)} = \frac{\left\langle g^{(3)}\left[\mathbf{E}(\mathbf{r}) \cdot \mathbf{E}'(\mathbf{r})\right]\left|\mathbf{E}'(\mathbf{r})\right|^2\right\rangle}{E_0^2\left|\mathbf{E}_0\right|^2}. \qquad (271)$$

Equation (271) is general and applicable to weak as well as strong nonlinearities. In the former case, $\mathbf{E}'(\mathbf{r}) \simeq \mathbf{E}(\mathbf{r})$, and Eq. (271) becomes

$$g_e^{(3)} = \frac{\left\langle g^{(3)}E^2(\mathbf{r})|\mathbf{E}(\mathbf{r})|^2\right\rangle}{E_0^2|\mathbf{E}_0|^2}, \qquad (272)$$

yielding $g_e^{(3)}$ in terms of the *linear* local field. Note that Eq. (272) is the analogue of (128). In the absence of metal grains, $g_e^{(3)} = g^{(3)}$. Therefore, the enhancement $I_K$ of the Kerr nonlinearity is given by

$$I_K = \frac{\left\langle E^2(\mathbf{r})\left|\mathbf{E}(\mathbf{r})\right|^2\right\rangle}{E_0^2\left|\mathbf{E}_0\right|^2} = M_{2,2}, \qquad (273)$$

where $M_{2,2}$ is the fourth moment of the local field.

Equations (272) and (273) were derived assuming that $g_m^{(3)} \simeq g_d^{(3)}$. If $g_m^{(3)} \neq g_d^{(3)}$, the above analysis can be repeated in order to derive the following equation,

$$g_e^{(3)} = pg_m^{(3)}\frac{\left\langle E^2(\mathbf{r})\left|\mathbf{E}(\mathbf{r})\right|^2\right\rangle_m}{E_0^2\left|\mathbf{E}_0\right|^2} + (1-p)g_d^{(3)}\frac{\left\langle E^2(\mathbf{r})|\mathbf{E}(\mathbf{r})|^2\right\rangle_d}{E_0^2\left|\mathbf{E}_0\right|^2}, \qquad (274)$$

where $\langle\cdot\rangle_m$ and $\langle\cdot\rangle_d$ represent averaging over the metal and dielectric grains, respectively. Note that, for the case of cubic nonlinearity in the conductivity of materials with percolation disorder, Eq. (274) was already derived and discussed in Section 3.2 [see Eq. (128)]. For the case of Kerr conductivity, Eq. (274) was first derived by Shalaev *et al.* (1998). According to Eq. (273), the Kerr enhancement $I_K$ is proportional to the fourth power of the local field, averaged over the sample, which is similar to the case of SERS with the enhancement factor $I_{RS}$ given by Eq. (254). Note, however, that while $I_K$ is complex, $I_{RS}$ is a real and positive quantity.

The enhancement of the Kerr nonlinearity can be estimated analytically using the methods described above. Consider first the case when $g^{(3)}(\mathbf{r})$ in the dielectric component is of the same order of magnitude or larger than in the metal component. Then,

$$
\begin{aligned}
I_K &\sim \left|g_e^{(3)}/\left\langle g^{(3)}(\mathbf{r})\right\rangle\right| = \left|\epsilon_e^{(3)}/\left\langle \epsilon^{(3)}(\mathbf{r})\right\rangle\right| \sim |M_{2,2}| \\
&\sim \mathcal{N}(\xi_A/a)^{d-8} \left(\frac{|\epsilon_m|}{\epsilon_d}\right)^{(2\nu+s)/(\mu+s)} \left(\frac{|\epsilon_m|}{\epsilon_m''}\right)^3,
\end{aligned}
\tag{275}
$$

where Eq. (215) has been used for the moment $M_{2,2}$. For $\omega \ll \omega_p$, the Kerr enhancement for 2D composites is estimated as, $I_K \sim \mathcal{N}(\xi_A/a)^{-6}(\omega_p/\omega_\tau)^3$, if the Drude formula is used for the metal dielectric constant $\epsilon_m$. For example, as discussed above, for silver-on-glass semi-continuous films, Anderson localization length $\xi_A \simeq 2a$ and density of states, $\mathcal{N} \simeq 1$, and therefore, $I_K \sim 10^5 - 10^6$. As discussed by Sarychev and Shalaev (2000), for $d = 2$ a plot of $I_K$ versus the metal volume fraction $p$ has a two-peak structure, which is similar to the case of Raman scattering shown in Figure 3.12. However, in contrast to $I_{RS}$, the dip at $p = p_c$ is much more pronounced and is proportional to the loss factor $\kappa$, implying that at $p = p_c$ the enhancement is actually given by, $I_K \sim \kappa M_{2,2}$. This result is presumably a consequence of the special symmetry of a 2D self-dual system at $p = p_c$. If one moves slightly away from $p = p_c$, the enhancement $I_K$ increases such that, $I_K \sim |M_{2,2}| \sim I_{RS} \sim M_{4,0}$. The fact that the minimum at $p = p_c$ is much smaller for SERS than for the Kerr process is presumably related to the latter being a phase sensitive effect. Moreover, as already discussed above, the local field maxima are concentrated in the dielectric gaps where $|\epsilon_m| \gg \epsilon_d$. Therefore, Eq. (275) is valid when the Kerr nonlinearity is located mainly in such gaps.

Consider now the case when the Kerr nonlinearity is due to metal grains (see, for example, Ma *et al.*, 1998; Liao *et al.*, 1998). Provided that $\epsilon_m' \simeq -\epsilon_d$, the local electric fields are equally distributed in the metal and dielectric components, implying that the Kerr enhancement is still given by Eq. (274) with $|\epsilon_m|/\epsilon_d = 1$. However, if $\epsilon_m \gg \epsilon_d$, the local field will be concentrated in the dielectric space between the conducting clusters with a value $E_m$ given by Eq. (207). The total current $I_s$ of the electric displacement flowing in the dielectric space between two resonate metal clusters of size $b_r$ is given by, $I_s = aE_m\epsilon_e b_r^{d-2}$. Because of the current continuity, the same current should flow in the adjacent metal clusters where it is concentrated in a percolating channel. The electric field $E_{mc}$ in the metal

channel, which spans over the cluster, is given by, $E_{mc} \sim I_s/(\epsilon_m a^{d-1})$, where $a^{d-1}$ represents the cross-section of the channel. Then the $n$th moment of the local electric field in a metal cluster of size $b_r$ is, $\langle E_{mc}^n \rangle = E_{mc}^n \mathcal{L} a^{d-1}/b_r^d$, where $\mathcal{L} = a(\epsilon_m/\epsilon_e)b_r^{-d+2}$ is the effective length of the conducting channel. Keeping in mind that only a fraction $\kappa = \epsilon_m''/|\epsilon_m| \ll 1$ of the metal clusters of size $b_r$ are excited by the external electric field, we obtain, $M_n^{met} = \langle |E|^n \rangle_{met}/E_0^n = \kappa \langle E_{mc}^n \rangle/E_0^n$, for the moments of the electric field in the metal component,

$$M_n^{met} \sim \left( \frac{|\epsilon_m|}{\epsilon_m''} \right)^{n-1} \left( \frac{|\epsilon_m|}{\epsilon_d} \right)^{[(d-1)(n-2)\nu - \mu(n-1)]/(\mu+s)}, \qquad (276)$$

where Eq. (206) was used for the size $b_r$ of the resonant clusters. Then, enhancement $I_K^{met}$ of the Kerr nonlinearity is given by

$$I_K^{met} \sim M_4^{met} \sim \left( \frac{|\epsilon_m|}{\epsilon_m''} \right)^3 \left( \frac{|\epsilon_m|}{\epsilon_d} \right)^{[2(d-1)\nu - 3\mu]/(\mu+s)}. \qquad (277)$$

In 2D (for which $\mu = s \simeq \nu = 4/3$) Eq. (277) yields, $I_K^{met} \sim M_4^{met} \sim (|\epsilon_m|/\epsilon_m'')^3$ $(\epsilon_d/|\epsilon_m|)^{1/2}$. As expected, $I_K^{met} \ll I_K$, and in fact for 2D systems near $p_c$,

$$\frac{I_K}{I_K^{met}} \sim \left( \frac{|\epsilon_m|}{\epsilon_d} \right)^2. \qquad (278)$$

Since in optic and infrared spectral ranges, $|\epsilon_m| \gg \epsilon_d$, the enhancement due to the Kerr nonlinearity is much larger than when the initial nonlinearity is located in the dielectric gaps where the local fields are much larger than in the metal. It follows from Eq. (278) that the Kerr enhancement $I_K^{met}$ may become less than one, implying that, on average, the local electric field in the metal component can be smaller than the external field. For example, for semi-continuous silver films on a glass substrate, $I_K^{met} < 1$ for wavelengths $\lambda > 10\mu$m.

## 3.5.4.2  Enhancement of Nonlinear Scattering from Strongly Disordered Films

The next subject we consider is percolation-enhanced nonlinear scattering (PENS) from a random metal-dielectric film at the metal volume fraction $p$ near $p_c$. Specifically, we consider the enhanced nonlinear scattering which is due to local field oscillation at frequency $n\omega$, while a percolation metal-dielectric film is exposed to an electromagnetic wave of frequency $\omega$. Since at $p_c$ a self-similar fractal metal cluster forms and the metal-dielectric transition occurs in a semi-continuous metal film, optical excitations of the self-similar cluster result in giant, scale-invariant, field fluctuations. As before, we assume that a semi-continuous film is exposed to the light that propagates normal to the film, with the wavelength $\lambda$ larger than any intrinsic length scale in the film. The space between the metal grains are filled by the dielectric substrate so that a semi-continuous metal film can be thought of as a 2D array of metal and dielectric grains that are randomly distributed over a plane for which we consider $n$th-order harmonic generation (nHG) for an incident wave of frequency $\omega$. The nHG is generated by the semi-continuous metal film that is

covered by a layer possessing a nonlinear conductivity $g^{(n)}$. The layer can be made of nonlinear organic molecules, semi-conductor quantum dots, or a quantum well on top of a percolation film. The local electric field $\mathbf{E}_\omega(\mathbf{r})$, induced in the film by the external field $\mathbf{E}_0$, generates in the layer the $n\omega$ current $g^{(n)}\mathbf{E}_\omega E_\omega^{n-1}$. Note that, strictly speaking, this expression is valid only for the scalar nonlinear conductivity and odd $n$. However, for obtaining order-of-magnitude estimates, we can use this formula for arbitrary $n$. The external field, oscillating at frequency $\omega$, is still denoted as $\mathbf{E}_0$, though the frequency is indicated explicitly for other fields. The nonlinear current $g^{(n)}\mathbf{E}_\omega E_\omega^{n-1}$, in turn, interacts with the film and generates the initial $n\omega$ electric field with an amplitude $\mathbf{E}^{(n)} = g^{(n)} E_\omega^{n-1}\mathbf{E}_\omega/g^{(\ell)}$, where $g^{(\ell)}$ is the *linear* conductivity of the nonlinear layer at frequency $n\omega$. The electric field $\mathbf{E}^{(n)}$ can be thought of as an inhomogeneous external field exciting the film at frequency $n\omega$.

The nHG current $\mathbf{I}^{(n)}$ induced in the film by the initial field $\mathbf{E}^{(n)}$ can be determined in terms of the non-local conductivity matrix $\Sigma(\mathbf{r}, \mathbf{r}')$ introduced by Eq. (241):

$$I_\beta^{(n)}(\mathbf{r}) = \int \mathbf{\Sigma}_{\beta\alpha}^{(n)}(\mathbf{r}, \mathbf{r}')E_\alpha^{(n)}(\mathbf{r}') \, d\mathbf{r}', \qquad (279)$$

where $\mathbf{\Sigma}_{\beta\alpha}^{(n)}$ is the conductivity matrix at frequency $n\omega$, the integration is over the entire film area, the Greek indices represent $\{x, y\}$, and summation over repeated indices is implied. It is $\mathbf{I}^{(n)}$ that eventually generates the nonlinear scattered field at frequency $n\omega$. By using the standard approach of the scattering theory adopted to semi-continuous metal films (Brouers *et al.*, 1998), and assuming that the incident light is unpolarized, the integral scattering in all directions but the specular one is given by

$$I = \frac{4k^2}{3c} \int \left( \left\langle I_\alpha^{(n)}(\mathbf{r}_1)I_\alpha^{(n)*}(\mathbf{r}_2)\right\rangle - \left|\left\langle \mathbf{I}^{(n)}\right\rangle\right|^2 \right) \, d\mathbf{r}_1 \, d\mathbf{r}_2, \qquad (280)$$

where the integrations is over the entire area $A$ of the film, $k = \omega/c$, and $\langle\cdot\rangle$ indicates an ensemble average. As in the case of Rayleigh scattering, we have assumed that the integrand vanishes for $r \ll \lambda$, where $\mathbf{r} = \mathbf{r}_2 - \mathbf{r}_1$ [therefore, the term $\exp(i\mathbf{k}\cdot\mathbf{r})$ was omitted]. Using Eq. (279), we can write

$$\int \left\langle I_\alpha^{(n)}(\mathbf{r}_1)\, I_\alpha^{(n)*}(\mathbf{r}_2)\right\rangle \, d\mathbf{r}_1 \, d\mathbf{r}_2$$

$$= \int \left\langle \Sigma_{\gamma\beta}^{(n)}(\mathbf{r}_1, \mathbf{r}_3)\Sigma_{\delta\alpha}^{(n)*}(\mathbf{r}_2, \mathbf{r}_4)\delta_{\gamma\delta}\left\langle E_\beta^{(n)}(\mathbf{r}_3)E_\alpha^{*(n)}(\mathbf{r}_4)\right\rangle_0\right\rangle \prod_{i=1}^{4} d\mathbf{r}_i, \qquad (281)$$

where $\langle\cdot\rangle_0$ denotes an average over the light polarization. We now introduce the spatially uniform probe field $\mathbf{E}_{n\omega}^{(0)}$ which oscillates at frequency $n\omega$ and is assumed to be unpolarized. For the unpolarized light, $\delta_{\gamma\delta} = 2\left\langle E_{n\omega,gamma}^{(0)}E_{n\omega,\delta}^{(0)*}\right\rangle_0/|E_{n\omega}^{(0)}|^2$, which, when substituted in Eq. (281), the integration is carried out over the coordinates $\mathbf{r}_1$ and $\mathbf{r}_2$, and the averaging over independent polarizations of fields $E_{n\omega}^{(0)}$ and $E_0$ are performed, the following equation for the current-current correlation

function is obtained,

$$\int \left\langle I_\alpha^{(n)}\,(\mathbf{r}_1)\,I_\alpha^{(n)*}\,(\mathbf{r}_2) \right\rangle\,d\mathbf{r}_1\,d\mathbf{r}_2 =$$

$$\frac{1}{|\mathbf{E}_{n\omega}^{(0)}|^2} \left\langle g_{n\omega}^{(\ell)}(\mathbf{r}_3)g_{n\omega}^{*}(\mathbf{r}_4)\left[\mathbf{E}_{n\omega}(\mathbf{r}_3)\cdot\mathbf{E}_{n\omega}^{*}(\mathbf{r}_4)\right]\left[\mathbf{E}^{(n)}(\mathbf{r}_3)\cdot\mathbf{E}^{(n)*}(\mathbf{r}_4)\right]\right\rangle\,d\mathbf{r}_3\,d\mathbf{r}_4,$$

(282)

where $E_{n\omega}(r)$ is the local $n\omega$ field excited in the film by the probe field $\mathbf{E}_{n\omega}^{(0)}$, and $g_{n\omega}^{(\ell)}(r)$ is the film linear conductivity at frequency $n\omega$.

In macroscopically-homogeneous and isotropic films considered here, the integral in Eq. (282) does not depend on direction of the probe field $\mathbf{E}_{n\omega}^{(0)}$. Therefore, $\mathbf{E}_{n\omega}^{(0)}$ can be selected to be collinear with the external field $\mathbf{E}_0$. Moreover, $\langle\mathbf{I}^{(n)}\rangle$ is parallel to the external field $\mathbf{E}_0$. If the probe field $\mathbf{E}_{n\omega}^{(0)}$ is aligned with $\mathbf{E}_0$, we have, $|\langle\mathbf{I}^{(n)}\rangle|^2 = \left|\left(\mathbf{E}_{n\omega}^{(0)}\cdot\mathbf{I}^{(n)}\right)\right|^2/|\mathbf{E}_{n\omega}^{(0)}|^2$. Then, using Eq. (279), we can write

$$\left|\left\langle\mathbf{I}^{(n)}\right\rangle\right|^2 = \frac{1}{A|\mathbf{E}_{n\omega}^{(0)}|^2}\left|\int E_{n\omega,\beta}^{(0)}\Sigma_{\beta\alpha}^{(n)}(\mathbf{r}_1,\mathbf{r}_2)E_\alpha^{(n)}(\mathbf{r}_2)\,d\mathbf{r}_1\,d\mathbf{r}_2\right|^2. \quad (283)$$

If the integration over coordinate $\mathbf{r}_1$ is carried out, one obtains

$$\left|\left\langle\mathbf{I}^{(n)}\right\rangle\right|^2 = \frac{\left|\left\langle g_{n\omega}^{(\ell)}\left(\mathbf{E}_{n\omega}\cdot\mathbf{E}^{(n)}\right)\right\rangle\right|^2}{|\mathbf{E}_{n\omega}^{(0)}|^2}, \quad (284)$$

and therefore

$$I = \frac{8\pi k^2}{3c|\mathbf{E}_{n\omega}^{(0)}|^2}\left|\frac{g^{(n)}}{g^{(\ell)}}\right|^2 A\left\langle\left|g\mathbf{E}_{n\omega}^{(\ell)}\right|^2|\mathbf{E}_\omega|^2\,|E_\omega|^{2(n-1)}\right\rangle\int_0^\infty C^{(n)}(r)r\,dr, \quad (285)$$

where $C^{(n)}(r)$ is the nonlinear correlation function defined as

$$C^{(n)}(r) =$$

$$\frac{\left\langle g_{n\omega}^{(\ell)}(\mathbf{r}_1)g_{n\omega}^{*}(\mathbf{r}_2)[\mathbf{E}_{n\omega}(\mathbf{r}_1)\cdot\mathbf{E}_{n\omega}^{*}(\mathbf{r}_2)][\mathbf{E}^{(n)}(\mathbf{r}_1)\cdot\mathbf{E}^{(n)*}(\mathbf{r}_2)]\right\rangle - \left|\left\langle g_{n\omega}^{(\ell)}\left(\mathbf{E}^{(n)}\cdot\mathbf{E}_{n\omega}\right)\right\rangle\right|^2}{\left\langle\left|g_{n\omega}^{(\ell)}\mathbf{E}_{n\omega}\right|^2\,\left|\mathbf{E}^{(n)}\right|^2\right\rangle},$$

(286)

which, for macroscopically-homogeneous and isotropic films, depends only on the distance $r = |\mathbf{r}_1 - \mathbf{r}_2|$.

Equation (285) should be compared with the $n\omega$ signal $I_{n\omega}$ from the nonlinear layer on a dielectric film with no metal grains on it, $S_{n\omega} = (c\epsilon_d^2/2\pi)A\left|g^{(n)}/g^{(\ell)}\right|^2\left|\mathbf{E}_\omega^{(0)}\right|^2\left|\mathbf{E}_\omega^{(0)}\right|^{2(n-1)}$. Therefore, the enhancement factor for PENS, $I_{PENS} = I/I_{n\omega}$, is given by

$$I_{PENS} = \frac{(ka)^4}{3}\frac{\left\langle|\epsilon_{n\omega}\mathbf{E}_{n\omega}|^2\,|\mathbf{E}_\omega|^2\,|\mathbf{E}_\omega|^{2(n-1)}\right\rangle}{\epsilon_d^2\left|\mathbf{E}_{n\omega}^{(0)}\right|^2\,|\mathbf{E}_0|^2\,|\mathbf{E}_0|^{2(n-1)}}\frac{n^2}{a^2}\int_0^\infty C^{(n)}(r)r\,dr. \quad (287)$$

Note that for a homogeneous ($p = 0$ or $p = 1$) surface, $C^{(n)}(r) = 0$, and therefore, $I_{PENS} = 0$, so that the scattering occurs only in the reflected direction. According to Eq. (287), the enhancement $I_{PENS}$ is proportional to $\langle |\mathbf{E}|^{2(n+1)} \rangle$ which, for highly fluctuating local fields, is very large. Since a metal-dielectric transition at $p_c$ is similar to a second-order phase transition, one may anticipate that local field fluctuations are rather large and have long-range correlations near $p_c$. However, what is surprising is that the field fluctuations in the optical spectral range discussed here are quite different from those for a second-order phase transition. The reason may be the following. The fluctuations in the local electric field that result in PENS are of the resonant character and their variations can be over several orders of magnitude. Therefore, the field correlation function $C^{(3)}(r)$ decreases very rapidly for $r > a$, and has a *negative* minimum, *regardless* of the magnitude of the local field correlation length $\xi_e$; this *anticorrelation* occurs because the field maxima have different signs. Moreover, the power-low decrease of $C^{(3)}(r)$, which is typical for critical phenomena, occurs only in the tail and deviates from it for $r > \xi_e$. The magnitude of $\xi_e$ can be estimated from Eq. (213) as, $\xi_e(\lambda) \simeq 5, 20$ and $30$ (in units of $a$, the grain size) for $\lambda = 0.34, 0.53$ and $0.9$ $\mu$m, respectively. For a typical size of a metal grain in a semi-continuous film, $a \simeq 2 - 20$ nm, the intrinsic spatial scale of the local field inhomogeneity $\xi_e \ll \lambda$, as assumed in advance. Based on such considerations, the dimensionless integral $a^{-2} \int_0^\infty C^{(n)}(r) r \, dr$ should be of the order of one for all $n$. Thus, one may anticipate that, in contrast to harmonic generation from conventional metal surfaces, PENS is characterized by a broad-angle distribution, with the total (in all directions) scattering being much larger than the coherent scattering in the reflected direction.

To obtain a more accurate estimate of PENS, we note that the typical size $b_r(\omega) \sim a\sqrt{|\epsilon_m(\omega)|}$ of the local field maxima increases with decreasing $\omega$, and thus for a Drude metal, $b_r(\omega) \propto \omega^{-1}$ if $\omega \ll \omega_p$. Since the spatial scales for $\mathbf{E}_{n\omega}$ and $\mathbf{E}_\omega$ are different, the average $\langle [|\epsilon \mathbf{E}_{n\omega}|^2 |\mathbf{E}_\omega|^2 |E_\omega|^{2(n-1)}]^2 \rangle$ in Eq. (287) can be decoupled and approximated roughly as, $\left\langle \left( |\epsilon_{n\omega} \mathbf{E}_{n\omega}|^2 |\mathbf{E}_\omega|^2 |E_\omega|^{2(n-1)} \right)^2 \right\rangle \sim \left\langle |\epsilon_{n\omega} E_{n\omega}|^2 \right\rangle \left\langle |E_\omega|^{2n} \right\rangle \sim |\epsilon_{n\omega} \epsilon_d| M_{2,n\omega} M_{2n} |E_{n\omega}^{(0)}|^2 |E_0|^{2n}$, where $M_{2,n\omega}$ and $M_{2n}$ are the spatial moments of the local fields $\mathbf{E}_{n\omega}$ and $\mathbf{E}_\omega$, respectively. Using this decoupling in Eq. (287) and taking into account the fact that, as discussed above, the integral there is of the order of unity, Eq. (287) simplifies to

$$\frac{I_{PENS}}{(ka)^4} \simeq B \left| \frac{\epsilon_m(n\omega)}{\epsilon_d} \right| M_{2,n\omega} M_{2n}, \tag{288}$$

where $B$ is an adjustable pre-factor. Finally, using Eq. (217) for the moments $M_{2,n\omega}$ and $M_{2n}$, and assuming that the localization length $\xi_A \sim a$, and that the density of states $\mathcal{N} \sim 1$, one obtains

$$\frac{I_{PENS}}{(ka)^4} \simeq B \frac{|\epsilon_m(n\omega)|^{5/2} |\epsilon_m(\omega)|^{3(n-1)/2}}{\epsilon_d^{n+1} \epsilon_m''(n\omega) \epsilon_m''(\omega)^{2n-1}}, \tag{289}$$

where it was assumed that the generated frequency $n\omega$ is less than $\omega_p$, so that

$\epsilon'_m (n\omega) < 0$; otherwise, $I_{PENS} \simeq B(ka)^4 M_\omega^{(2n)}$, since the local $n\omega$ fields are not enhanced for $\epsilon'_m (n\omega) > 0$. For the Drude metal and $n\omega \ll \omega_p$, Eq. (289) is simplified to

$$I_{PENS} \sim B (ka)^4 \frac{1}{\epsilon_d^{n+1}} \left(\frac{\omega_p}{\omega_\tau}\right)^{2n} \left(\frac{\omega_p}{\omega}\right)^2 . \qquad (290)$$

Equation (290) states that PENS increases with increasing the order of a nonlinear process and decreases toward the infrared part of the spectrum as $I_{PENS} \propto \lambda^{-2}$, in contrast to the well-known $\lambda^{-4}$ law for Rayleigh scattering. Moreover, it is interesting to note that, for high-harmonic scattering, PENS is proportional $\lambda^{-2}$, independently of the order $n$ of optical nonlinearity.

### 3.5.4.3   Comparison with the Experimental Data

The diffusive scattering of the second harmonic from metal-dielectric films has been observed in experiments with $C_{60}$-coated semi-continuous silver films (Akt-sipetrov *et al.*, 1993) and from thin but *continuous* silver films (Kuang and Simon, 1995). One may argue that the diffusive scattering of $2\omega$ field is due to the anomalous fluctuations of local electric fields on the rough features of the surface with the spatial scale $a$ being much smaller than wavelength $\lambda$ of the incident light. If so, then the scattering data reported by Kuang and Simon (1995) are similar to PENS from percolation films.

  To summarize, large field fluctuations in random metal-dielectric composites near $p_c$ result in a new physical phenomenon: Percolation-enhanced nonlinear scattering which is characterized by giant enhancement and a broad-angle distribution.

## 3.6   Electromagnetic Properties of Solid Composites

In the preceding discussions, the skin effects in the metal grains was neglected. We now consider electromagnetic properties of metal-dielectric materials, characterized by percolation disorder and irradiated by a high-frequency electromagnetic field under the conditions that the skin effect in the metal grains is strong. The goal of this section is to show that electromagnetic properties of random composites can be understood in terms of the effective dielectric constant and magnetic permeability, provided that the wavelength of an incident wave is much larger than the intrinsic spatial scale of the system. The wavelength inside a metal component can be very small. The most interesting effects are expected in the limit of the strong skin effect. Thus, one must go beyond the quasi-static approximation employed in the analyses presented above.

  Propagation of electromagnetic waves in percolation composites with wavelength $\lambda < \xi_p$, where $\xi_p$ is the correlation length of percolation, may be accompanied by strong scattering. On the other hand, wave propagation for $\lambda \gg \xi_p$

can be described by Maxwell's equations with effective dielectric constant $\epsilon_e$ and effective magnetic permeability $K_e$. In order to calculate these effective parameters, the approach suggested by Panina *et al.* (1990), as developed by Sarychev and Shalaev (2000), is described. We restrict our attention to the optically-thin systems of size $\mathcal{L} \ll \lambda/\sqrt{|\epsilon_e K_e|}$, which are still macroscopically homogeneous, so that $\mathcal{L} \gg \xi_p$. We already described in Section 4.13 of Volume I the theoretical treatment of this problem for linear materials, and what follows is the extension of that discussion to nonlinear composites.

## 3.6.1  Effective-Medium Approximation

Suppose that a percolation composite is placed inside a resonator, where electromagnetic standing waves are excited. The change in the field when a composite is placed inside of the resonator is determined by superposition of the fields scattered from individual metal and dielectric particles that have dielectric constants $\epsilon_m$ and $\epsilon_d$, respectively. The interaction between the particles can be taken into account by an effective-medium approximation (EMA). As discussed in the previous sections, and in Volume I, in this method, the interaction of a given metal or dielectric particle with the rest of the system is determined by replacing the latter by a homogeneous medium with the effective parameters $\epsilon_e$ and $K_e$. Assuming that the composite grains are spherical, the electric fields $\mathbf{E}_{in,m}$ and $\mathbf{E}_{out,m}$, excited by the external electric field $\mathbf{E}_0$, are calculated inside and outside of a metal grain of size $a$, yielding the following equations (see also Chapter 4 of Volume I) for the electric field inside the metal grain:

$$\mathbf{E}_{in,m}(\mathbf{r}) = \mathbf{E}_{in,m0} + 4\pi \mathbf{L}(\mathbf{r}), \tag{291}$$

where

$$\mathbf{E}_{in,m0} = \frac{3\epsilon_e}{2\epsilon_e + \tilde{\epsilon}_m} \mathbf{E}_0, \tag{292}$$

and $\tilde{\epsilon}_m$ is the renormalized dielectric constant of the metal defined as

$$\tilde{\epsilon}_m = \epsilon_m \frac{2F(y_m a)}{1 - F(y_m a)}, \qquad F(x) = \frac{1}{x^2} - \frac{\cot(x)}{x}, \tag{293}$$

with $k = \omega/c$, $y_m = k\sqrt{\epsilon_m K_m}$, $a$ being the radius of a metal grain. The skin (penetration) depth $\delta$ is given by, $\delta = 1/\mathrm{Im}(y_m)$. When the metal conductivity $g_m$ is a real quantity (i.e., in the microwave and radio frequency range), the skin depth, $\delta = c/\sqrt{2\pi K_m g_m \omega}$. In Cartesian coordinate system with the $z$-axis directed along the field $\mathbf{E}_0$, the local electric field $\mathbf{L}$ in Eq. (291) is determined by

$$\nabla \times \mathbf{L}(\mathbf{r}) = \frac{1}{4\pi} \nabla \times \mathbf{E}_{in,m}(\mathbf{r}) = \frac{ik}{4\pi} \mathbf{B}_E, \tag{294}$$

where

$$\mathbf{B}_E = -3i\mathbf{E}_0 \frac{ak\epsilon_m \epsilon_e \sin(y_m r) F(y_m r)}{(2\epsilon_e + \tilde{\epsilon}_m) \sin(y_m a) [F(y_m a) - 1]} \left\{ \frac{y}{r}, -\frac{x}{r}, 0 \right\} \tag{295}$$

is a rotational magnetic induction generated in a metal particle by the electric

current. Therefore, the inside electric field consists of uniform curl-free part $\mathbf{E}_{in,m0}$ (i.e., $\nabla \times \mathbf{E}_{in,m0}$) and the rotational part $\mathbf{L}(\mathbf{r})$ that depends on the coordinate. The field outside the metal particle is given by

$$\mathbf{E}_{out,m} = \mathbf{E}_0 + a^3 \frac{\epsilon_e - \tilde{\epsilon}_m}{2\epsilon_e + \tilde{\epsilon}_m} \nabla \left( \frac{\mathbf{E}_0 \cdot \mathbf{r}}{r^3} \right). \tag{296}$$

The local wavelength inside a dielectric grain, $\lambda_d = \lambda/\sqrt{\epsilon_d}$, is assumed to be much larger than the grain size $a$. Then, the electric fields inside and outside a dielectric particle are given by the following well-known equations, already familiar from Chapters 4 and 5 of Volume I:

$$\mathbf{E}_{in,d} = \mathbf{E}_0 \frac{3\epsilon_e}{2\epsilon_e + \epsilon_d}, \tag{297}$$

$$\mathbf{E}_{out,d} = \mathbf{E}_0 + a^3 \frac{\epsilon_e - \epsilon_d}{2\epsilon_e + \epsilon_d} \nabla \left( \frac{\mathbf{E}_0 \cdot \mathbf{r}}{r^3} \right). \tag{298}$$

Similar equations can be obtained for the magnetic field excited by a uniform magnetic field $\mathbf{H}_0$ inside and outside a metal (dielectric) particle:

$$\mathbf{H}_{in,m} = \mathbf{H}_{in,m0} + 4\pi \mathbf{M}, \tag{299}$$

where

$$\mathbf{H}_{in,m0} = \frac{3K_e}{2K_e + \tilde{K}_m} \mathbf{H}_0, \tag{300}$$

and the renormalized metal magnetic permeability $\tilde{K}_m$ is given by

$$\tilde{K}_m = K_m \frac{2F(y_m a)}{1 - F(y_m a)}, \tag{301}$$

where the function $F$ is defined by Eq. (293). Note that the renormalized metal magnetic permeability $\tilde{K}_m$ is *not* equal to one, even if the metal is non-magnetic and the seed magnetic permeability $K_m = 1$. The local magnetic field $\mathbf{M}$ in Eq. (299) is the solution of

$$\nabla \times \mathbf{M} = \frac{1}{4\pi} \nabla \times \mathbf{H}_{in,m} = -\frac{ik}{4\pi} \mathbf{D}_H, \tag{302}$$

with

$$\mathbf{D}_H = 3i\mathbf{H}_0 \frac{ak K_m K_e \sin(y_m r) F(k_m r)}{(2K_e + \tilde{K}_m) \sin(y_m a)[F(y_m a) - 1]} \left\{ \frac{y}{r}, -\frac{x}{r}, 0 \right\}, \tag{303}$$

where $\mathbf{D}_H$ is the electric displacement induced in the metal particle by high-frequency magnetic field $\mathbf{H}_0$. The displacement $\mathbf{D}_H$ can be written as $\mathbf{D}_H = i(4\pi/\omega)\mathbf{I}$, where the eddy electric current $\mathbf{I}$ is called the *Foucault current*. The field $\mathbf{H}_{in,m0}$ is the potential (curl-free) part, while $\mathbf{M}$ is the rotational part of the local magnetic field. The magnetic field outside the metal particle is irrotational and equals to

$$\mathbf{H}_{out,m} = \mathbf{H}_0 + a^3 \frac{K_e - \tilde{K}_m}{2K_e + \tilde{K}_m} \nabla \left( \frac{\mathbf{H}_0 \cdot \mathbf{r}}{r^3} \right). \tag{304}$$

We assume, for simplicity, that the dielectric component of the composite is non-magnetic, i.e., the dielectric magnetic permeability $K_d = 1$. Then,

$$\mathbf{H}_{in,d} = \mathbf{H}_0 \frac{3K_e}{2K_e + 1}, \tag{305}$$

and

$$\mathbf{H}_{out,d} = \mathbf{H}_0 + a^3 \frac{K_e - 1}{2K_e + 1} \nabla \left( \frac{\mathbf{H}_0 \cdot \mathbf{r}}{r^3} \right). \tag{306}$$

As in all the EMAs described so far in this book, the effective parameters $\epsilon_e$ and $K_e$ are determined by the self-consistent condition that the fluctuations in the fields should vanish when averaged over all the (spherical) inclusions, i.e., $\langle \mathbf{E}_{out} \rangle = p\mathbf{E}_{out,m} + (1-p)\mathbf{E}_{out,d} = \mathbf{E}_0$, and $\langle \mathbf{H}_{out} \rangle = p\mathbf{H}_{out,m} + (1-p)\mathbf{H}_{out,d} = \mathbf{H}_0$, where $\langle \cdot \rangle$ indicates a volume averaging. Therefore, when these averagings are carried out, they results in the following equations

$$p\frac{\epsilon_e - \tilde{\epsilon}_m}{2\epsilon_e + \tilde{\epsilon}_m} + (1-p)\frac{\epsilon_e - \epsilon_d}{2\epsilon_e + \epsilon_d} = 0, \tag{307}$$

$$p\frac{K_e - \tilde{K}_m}{2K_e + \tilde{K}_m} + (1-p)\frac{K_e - 1}{2K_e + 1} = 0. \tag{308}$$

These equations are completely similar to the traditional EMAs discussed in the previous sections and Volume I. It can be seen that the skin effect results in renormalization of the dielectric constant and magnetic permeability of the conducting component. Specifically, the metal dielectric constant $\epsilon_m$ and magnetic permeability $K_m$ are replaced by $\tilde{\epsilon}_m$ and $\tilde{K}_m$ given by Eqs. (293) and (301), respectively. This fact has an important effect on the frequency dependence of the effective parameters. For example, it is commonly accepted that the effective conductivity $g_e = -i\omega\epsilon_e/(4\pi)$ of a composite is dispersion-free, when the conductivity of metal component $g_m$ is independent of frequency and $g_m \gg \omega$ (which is typical for the microwave and far-infrared ranges). Thus, as shown in Chapter 5 of Volume I [see Eq. (5.62) there], the traditional EMA predicts that, $g_e = g_m(3p - 1)/2$ for $p > p_c$. Equation (307) yields the same result for the effective conductivity $g_e$, but with the metal conductivity being renormalized according to Eq. (293), which results in, $g_e = g_m F(y_m a)(3p - 1)/[1 - F(y_m a)]$. Thus, the effective conductivity has a dispersive behavior, provided that the skin effect in metal grains is important. In the limit of very strong skin effect, $\delta \ll a$, the effective conductivity decreases with the frequency as, $g_e \sim g_m(\delta/a) \sim g_m/\sqrt{\omega}$.

Another interesting prediction is that percolation composites exhibit magnetic properties, even if such properties are absent in each component, i.e., even if $K_m = K_d = 1$. In this case, the real part $K_e'$ of the effective magnetic permeability $K_e$ is less than one and decreases with frequency, while its imaginary part $K_e''$ has its maximum at frequencies such that, $\delta \sim a$.

One can now show that the effective parameters $\epsilon_e$ and $K_e$ determine propagation of an electromagnetic wave in the metal-dielectric composites. The average

electric field is equal to

$$\langle \mathbf{E} \rangle = p\mathbf{E}_{in,m} + (1-p)\mathbf{E}_{in,d} = p\mathbf{E}_{in,m0} + 4\pi \langle \mathbf{L} \rangle + (1-p)\mathbf{E}_{in,d}. \quad (309)$$

When Eqs. (291) and (297) are substituted in Eq. (309) and Eq. (307) is taken into account, Eq. (309) simplifies to

$$\langle \mathbf{E} \rangle = \mathbf{E}_0 + 4\pi \langle \mathbf{L} \rangle, \quad (310)$$

where $\langle \cdot \rangle$ indicates an average over the volume of the system. Therefore, the irrotational part of the local field, being averaged over the volume, gives the field $\mathbf{E}_0$, while the second term of Eq. (310) results from the skin effect in metal grains. In a similar fashion, we obtain

$$\langle \mathbf{H} \rangle = p\mathbf{H}_{in,m} + (1-p)\mathbf{H}_{in,d} = \mathbf{H}_0 + 4\pi \langle \mathbf{M} \rangle, \quad (311)$$

where the rotational field $\mathbf{M}$ in the metal grains is given by Eq. (302), and $\mathbf{M} = \mathbf{0}$ in the dielectric grains.

Consider now the average electric displacement $\langle \mathbf{D} \rangle$ induced in the system by the electric field $\mathbf{E}_0$, which can be written as

$$\langle \mathbf{D} \rangle = \epsilon_m p\mathbf{E}_{in,m0} + 4\pi\epsilon_m \langle \mathbf{L} \rangle + (1-p)\epsilon_d\mathbf{E}_{in,d}. \quad (312)$$

It follows from Eq. (292) for $\mathbf{E}_{in,m0}$ and Eq. (294) for $\mathbf{L}$ that the sum, $\epsilon_m p\mathbf{E}_{in,m0} + 4\pi\epsilon_m\langle \mathbf{L} \rangle$, in Eq. (312) can be written as

$$\epsilon_m p\mathbf{E}_{in,m0} + 4\pi\epsilon_m\langle \mathbf{L} \rangle = \epsilon_m p \left( \frac{3\epsilon_e}{2\epsilon_e + \tilde{\epsilon}_m}\mathbf{E}_0 + \frac{4\pi}{\Omega}\int \mathbf{L}\, d\mathbf{r} \right)$$

$$= \epsilon_m p \left[ \frac{3\epsilon_e}{2\epsilon_e + \tilde{\epsilon}_m}\mathbf{E}_0 + i\frac{k}{2\Omega}\int (\mathbf{r} \times \mathbf{B}_E)\ d\mathbf{r} \right] = p\frac{3\epsilon_e\tilde{\epsilon}_m}{2\epsilon_e + \tilde{\epsilon}_m}\mathbf{E}_0, \quad (313)$$

where the integration is over the volume $\Omega = 4\pi a^3/3$ of a metal particle, and $\mathbf{B}_E$ is given by Eq. (295). Substitution of Eqs. (313) and (297) into (312) yields

$$\langle \mathbf{D} \rangle = \epsilon_e\mathbf{E}_0. \quad (314)$$

Therefore, the average electric displacement is proportional to the irrotational part of the local field, and the proportionality coefficient is exactly equal to the effective dielectric constant. In a similar way, we obtain

$$\langle \mathbf{B} \rangle = K_e\mathbf{H}_0. \quad (315)$$

Equations (314) and (315) can be considered as definitions of the fields $\mathbf{E}_0$ and $\mathbf{H}_0$. Indeed, if the local fields were known in the composite, the fields $\mathbf{E}_0$ and $\mathbf{H}_0$ could be determined from these equations. Then, Eqs. (314) and (315) can be used to determine the effective dielectric constant $\epsilon_e$ and the effective magnetic permeability $K_e$ of a composite. These equations replace the usual constitutive equations, $\langle \mathbf{D} \rangle = \epsilon_e\langle \mathbf{E} \rangle$ and $\langle \mathbf{B} \rangle = K_e\langle \mathbf{H} \rangle$, which hold only in the quasi-static case.

We now derive the governing equations for the macroscopic electromagnetism in metal-dielectric composites. Equation (314) provides the average electric displacement excited by the electric field $\mathbf{E}_0$, but the local magnetic field also excites the

Foucault currents. Adding the electric displacement $\mathbf{D}_H$ given by Eq. (303) to the average displacement given by Eq. (314) yields the complete electric displacement,

$$\langle \mathbf{D} \rangle_f = \epsilon_e \mathbf{E}_0 + \frac{4\pi i}{k} \langle \mathbf{\nabla} \times \mathbf{M} \rangle. \tag{316}$$

Note that the second term of Eq. (316) disappears when the skin effect vanishes, i.e., when $|y_m|a \to 0$. We are still assuming that the linear size of the sample is much smaller than the wavelength $\lambda$. Similarly, the average magnetic induction $\langle \mathbf{B} \rangle_f$ is given by

$$\langle \mathbf{B} \rangle_f = K_e \mathbf{H}_0 - \frac{4\pi i}{k} \langle \mathbf{\nabla} \times \mathbf{L} \rangle. \tag{317}$$

At this point, the Maxwell's equations are averaged over macroscopic volume $\Omega \sim \mathcal{L}^3$, centered at point $\mathbf{r}$, such that $\xi_p \ll \mathcal{L} \ll \lambda$, yielding

$$\langle \mathbf{\nabla} \times \mathbf{E} \rangle = ik \langle \mathbf{B} \rangle_f = ik K_e \mathbf{H}_0 + 4\pi \langle \mathbf{\nabla} \times \mathbf{L} \rangle, \tag{318}$$

$$\langle \mathbf{\nabla} \times \mathbf{H} \rangle = -ik \langle \mathbf{D} \rangle_f = -ik\epsilon_e \mathbf{E}_0 + 4\pi \langle \mathbf{\nabla} \times \mathbf{M} \rangle. \tag{319}$$

The order of the curl operation and the volume averages in Eqs. (318) and (319) can be interchanged, as is usually done for derivation of the macroscopic Maxwell's equations. For example, $\langle \mathbf{\nabla} \times \mathbf{E} \rangle = \mathbf{\nabla} \times [\langle \mathbf{E} \rangle (\mathbf{r})]$, where $(\mathbf{r})$ indicates that the differentiation is over the position $\mathbf{r}$ of the volume $\Omega$. Then, the Maxwell's equations, Eqs. (318) and (319), become

$$\mathbf{\nabla} \times \mathbf{E}_0(\mathbf{r}) = ik K_e \mathbf{H}_0(\mathbf{r}), \tag{320}$$

$$\mathbf{\nabla} \times \mathbf{H}_0(\mathbf{r}) = -ik\epsilon_e \mathbf{E}_0(\mathbf{r}), \tag{321}$$

which have the typical forms for macroscopic electromagnetism, describing propagation of electromagnetic waves in composite media.

It is important to recognize that all quantities in Eqs. (310), (313), (314), (320), and (321) are well-defined and do not depend on the assumptions made in the course of their derivation. Thus, for example, $\langle \mathbf{M} \rangle$ in Eq. (311) can be determined as a magnetic moment of the Foucault currents per unit volume, so that

$$\langle \mathbf{M} \rangle = \frac{ik}{8\pi\Omega} \int (\mathbf{r} \times \mathbf{D}_H) \, d\mathbf{r} = \frac{1}{2c\Omega} \int (\mathbf{r} \times \mathbf{j}_H) \, d\mathbf{r}, \tag{322}$$

where the integration now is over the volume $\Omega$. This definition of $\langle \mathbf{M} \rangle$ is in agreement with Eq. (302), except that it is not required that the currents $\mathbf{I}_H$ be the same in all the metal particles. In a similar way, one may write

$$\langle \mathbf{L} \rangle = \frac{ik}{8\pi\Omega} \int (\mathbf{r} \times \mathbf{B}_E) \, d\mathbf{r}, \tag{323}$$

where the integration is still over the volume $\Omega$, and $\mathbf{B}_E = -(4\pi i/k)\mathbf{\nabla} \times \mathbf{E}$, with $\mathbf{E}$ being the local electric field. Note that $\langle \mathbf{L} \rangle$ has no direct analogue in the classical electrodynamics, since there is no such thing as loop magnetic currents in atoms and molecules.

## 3.7 Beyond the Quasi-static Approximation: Generalized Ohm's Law

The analysis presented above cannot be used for describing the optical properties of semi-continuous films in the important case in which skin effects in the metal grains are strong. Sarychev *et al.* (1994,1995) attempted to extend the above theoretical analysis beyond the quasi-static approximation, which is based on the full set of Maxwell's equations. In their approach the quasi-static approximation is not used because the fields are not assumed to be curl-free inside the film. In this section we summarize their theoretical analysis and discuss its implications. We restrict ourselves to the case in which all the external fields are parallel to the plane of the film. This means that an incident wave, as well as the reflected and transmitted waves, are travelling in the direction perpendicular to the film plane. The analysis is focused on the electric and magnetic fields at certain distances *away* from the film and attempts to relate them to the currents *inside* the film. We assume that the film's heterogeneities are over length scales that are much smaller than the wavelength $\lambda$, but not necessarily smaller than the skin depth $\delta$, so that the fields away from the film are curl-free and can be expressed as gradients of potential fields. The electric and magnetic induction currents, averaged over the film thickness, obey the usual 2D continuity equations. Therefore, equations such as, $\nabla \times \mathbf{E} = \mathbf{0}$, and $\nabla \cdot \mathbf{I} = 0$, are the *same* as in the quasi-static case. The only differences are that the fields and the average currents are now related by new constitutive equations, and that there are magnetic as well as electric currents.

In contrast to the traditional analyses, it is not assumed that the electric and magnetic fields inside a semi-continuous metal film are curl-free and $z$-independent, where the $z$-coordinate is perpendicular to the film plane. Let us consider first a homogeneous conducting film with a uniform conductivity $g_m$ and thickness $d$, and assume constant electric field $\mathbf{E}_1$ and magnetic field $\mathbf{H}_1$ at some reference plane $z = -d/2 - l_0$ behind the film, as shown in Figure 3.13. Under these conditions, the fields depend only on the $z$-coordinate, and Maxwell's equations for a monochromatic field can be written as

$$\frac{d}{dz}\mathbf{E}(z) = -\frac{i\omega}{c}K(z)[\mathbf{n} \times \mathbf{H}(z)], \tag{324}$$

$$\frac{d}{dz}\mathbf{H}(z) = -\frac{4\pi}{c}g(z)[\mathbf{n} \times \mathbf{E}(z)], \tag{325}$$

with boundary conditions

$$\mathbf{E}(z = -d/2 - l_0) = \mathbf{E}_1, \quad \mathbf{H}(z = -d/2 - l_0) = \mathbf{H}_1, \tag{326}$$

where $\mathbf{E}_1$ and $\mathbf{H}_1$ are parallel to the film plane. Here, the conductivity $g(z)$ is equal to the metal conductivity $g_m$ inside the film ($-d/2 < z < d/2$) and to $g_d = -i\omega/4\pi$ outside the film ($z < -d/2$ and $z > d/2$). Similarly, the magnetic permeability $K(z) = K_m$ and 1 inside and outside the film, respectively; the unit vector $\mathbf{n} = \{0, 0, 1\}$ is perpendicular to the film plane. When solving Eqs. (324) and (325), we must take into account the fact that the electric and magnetic fields are continuous at the film boundaries. Then, electric ($\mathbf{I}_E$) and magnetic ($\mathbf{I}_H$) cur-

FIGURE 3.13. Schematics of the model used in the computations. Electromagnetic wave of wavelength $\lambda$ is incident on a thin metal-insulator film of thickness $d$. The wave is partially reflected and absorbed, and the remainder passes through the film (after Sarychev and Shalaev, 2000).

rents flowing in between the two planes at $z = -d/2 - l_0$ and $z = d/2 + l_0$ are calculated as

$$\mathbf{I}_E = -\frac{i\omega}{4\pi} \left[ \int_{-d/2-l_0}^{-d/2} \mathbf{E}(z) \, dz + \int_{-d/2}^{d/2} \epsilon_m \mathbf{E}(z) \, dz + \int_{d/2}^{d/2+l_0} \mathbf{E}(z) \, dz \right],$$

(327)

$$\mathbf{I}_H = \frac{i\omega}{4\pi} \left[ \int_{-d/2-l_0}^{-d/2} \mathbf{H}(z) \, dz + \int_{-d/2}^{d/2} K_m \mathbf{H}(z) \, dz + \int_{d/2}^{d/2+l_0} \mathbf{H}(z) \, dz \right],$$

(328)

where $\epsilon_m = 4i\pi g_m/\omega$ is the metal dielectric constant. We assume, for simplicity, that the magnetic permeability $K_m = 1$. Since the Maxwell's equations are linear, the local fields $\mathbf{E}(z)$ and $\mathbf{H}(z)$ are linear functions of the boundary values $\mathbf{E}_1$ and

$\mathbf{H}_1$ defined at the plane $z = -d/2 - l_0$:

$$\mathbf{E}(z) = a(z)\mathbf{E}_1 + c(z)(\mathbf{n} \times \mathbf{H}_1), \tag{329}$$

$$\mathbf{H}(z) = b(z)\mathbf{H}_1 + d(z)(\mathbf{n} \times \mathbf{E}_1). \tag{330}$$

By substituting Eq. (329) for $\mathbf{E}(z)$ and (330) for $\mathbf{H}(z)$ in Eqs. (229) and (230), we can express the currents $\mathbf{I}_E$ and $\mathbf{I}_H$ in terms of the boundary fields $\mathbf{E}_1$ and $\mathbf{H}_1$:

$$\mathbf{I}_E = s\mathbf{E}_1 + g_1(\mathbf{n} \times \mathbf{H}_1), \tag{331}$$

$$\mathbf{I}_H = m\mathbf{H}_1 + g_2(\mathbf{n} \times \mathbf{E}_1). \tag{332}$$

Note that, Eq. (331) implies that, in contrast to the usual constitutive equations, the current $\mathbf{I}_E$ (which flows between the planes $z = -d/2 - l_0$ and $z = d/2 + l_0$) depends not only on the external electric field $\mathbf{E}_1$, but also on the external magnetic field $\mathbf{H}_1$, and similarly for the current $\mathbf{I}_H$. These equations are referred to as the generalized Ohm's law (GOL). The Ohmic parameters $s$, $m$, $g_1$ and $g_2$ have the dimension of surface conductivity and depend on the frequency $\omega$, the metal dielectric constant $\epsilon_m$, the film thickness $d$, and the distance $l_0$ between the film and the reference plane $z = -d/2 - l_0$. We assume that the films are invariant under reflection through the plane $z = 0$. In this case (Sarychev $et\ al.$, 1995), $g_1 = g_2 = g$. The Ohmic parameter $g$ is then expressed in terms of parameters $s$ and $m$ as

$$g = -\frac{c}{4\pi} + \sqrt{\left(\frac{c}{4\pi}\right)^2 - ms}. \tag{333}$$

Thus, the GOL equations take the following forms

$$\mathbf{I}_E = s\mathbf{E}_1 + g(\mathbf{n} \times \mathbf{H}_1), \tag{334}$$

$$\mathbf{I}_H = m\mathbf{H}_1 + g(\mathbf{n} \times \mathbf{E}_1). \tag{335}$$

The Ohmic parameters $s$ and $m$ can be expressed in terms of the film refractive index $\eta = \sqrt{\epsilon_m}$ and its thickness $d$:

$$s = \frac{c}{8n\pi} \left\{ \exp(-idk\eta) [\eta \cos(adk) - i \sin(adk)]^2 - \exp(idk\eta) [\eta \cos(adk) + i \sin(adk)]^2 \right\}, \tag{336}$$

$$m = \frac{c}{8\eta\pi} \left\{ \exp(-idk\eta) [i \cos(adk) + \eta \sin(adk)]^2 - \exp(idk\eta) [-i \cos(adk) + \eta \sin(adk)]^2 \right\}, \tag{337}$$

where $k = \omega/c$. We still assume, for simplicity, that $\epsilon = 1$ outside the film, and introduce a dimensionless parameter $a \equiv l_0/d$. In these notations, the skin (penetration) depth $\delta$ is, $\delta = 1/k[\mathrm{Im}(n)]$. In the microwave spectral range, the metal conductivity is real while the dielectric constant $\epsilon_m$ is purely imaginary, so that, $\delta = c/\sqrt{2\pi g_m \omega}$. On the other hand, the dielectric constant is negative for a typical metal in the optical and infrared spectra ranges, and therefore, in this case, $\delta \simeq 1/k\sqrt{|\epsilon_m|}$.

In the case of laterally heterogeneous films, the currents $\mathbf{I}_E$ and $\mathbf{I}_H$, as well as the fields $\mathbf{E}_1$ and $\mathbf{H}_1$, are functions of the 2D vector $\mathbf{r} = \{x, y\}$. It follows from

Maxwell's equations that the fields and currents are connected by linear relations given by

$$\mathbf{I}_E(\mathbf{r}) = s\mathbf{E}_1 + g(\mathbf{n} \times \mathbf{H}_1), \tag{338}$$

$$\mathbf{I}_H(\mathbf{r}) = m\mathbf{H}_1 + g(\mathbf{n} \times \mathbf{E}_1), \tag{339}$$

in which $s$, $m$ and $g$ represent integral operators. The metal islands in semi-continuous films usually have an oblate shape, so that the grain diameter $D$ is much larger than the film thickness $d$. When the thickness $d$ of a conducting grain (or the skin depth $\delta$) and the distance $l_0$ are much smaller than the grain diameter $D$, the relations between the fields $\mathbf{E}_1$ and $\mathbf{H}_1$ on one hand and the currents on the other hand become completely local in Eqs. (338) and (339). The local Ohmic parameters $s = s(\mathbf{r})$, $m = m(\mathbf{r})$, and $g = g(\mathbf{r})$, given by Eqs. (333), (336) and (337), are determined by the local refraction index, $\eta(\mathbf{r}) = \sqrt{\epsilon(\mathbf{r})}$, where $\epsilon(\mathbf{r})$ is a local dielectric constant. Equations (338) and (339) are the local GOL for semi-continuous films. For binary metal-dielectric semi-continuous films the local dielectric constant is equal to either $\epsilon_m$ or $\epsilon_d$. The electric ($\mathbf{I}_E$) and magnetic ($\mathbf{I}_H$) currents given by Eqs. (338) and (339) lie in between the planes $z = -d/2 - l_0$ and $z = d/2 + l_0$, and satisfy the usual 2D continuity equation, $\nabla \cdot \mathbf{I}_E(\mathbf{r}) = 0$, and $\nabla \cdot \mathbf{I}_H(\mathbf{r}) = 0$, which follow from the 3D continuity equations when the $z$-components of $\mathbf{E}_1$ and $\mathbf{H}_1$ are neglected at the planes $z = \pm(d/2 + l_0)$, made possible by the fact that these components are small since the average fields $\langle \mathbf{E}_1 \rangle$ and $\langle \mathbf{H}_1 \rangle$ are parallel to the film plane. Since we are considering semi-continuous films with a scale of heterogeneities much smaller than the wavelength $\lambda$, the fields $\mathbf{E}_1(\mathbf{r})$ and $\mathbf{H}_1(\mathbf{r})$ are still the gradients of potential fields when considered as functions of $x$ and $y$ in the fixed reference plane $z = -d/2 - l_0$, i.e.,

$$\mathbf{E}_1(\mathbf{r}) = -\nabla\varphi_1(\mathbf{r}), \quad \mathbf{H}_1(\mathbf{r}) = -\nabla\psi_1(\mathbf{r}). \tag{340}$$

By substituting these expressions and Eqs. (338) and (339) in the continuity equation, one obtains

$$\nabla \cdot [s\nabla\varphi_1 + g(\mathbf{n} \times \nabla\psi_1)] = 0, \qquad \nabla \cdot [m\nabla\psi_1 + g(\mathbf{n} \times \nabla\varphi_1)] = 0. \tag{341}$$

Equations (341) must be solved with the conditions that

$$\langle \nabla\varphi_1 \rangle = \langle \mathbf{E}_1 \rangle, \quad \langle \nabla\psi_1 \rangle = \langle \mathbf{H}_1 \rangle, \tag{342}$$

where the constant fields $\langle \mathbf{E}_1 \rangle$ and $\langle \mathbf{H}_1 \rangle$ are external (given) fields, and $\langle \cdot \rangle$ indicates an average over coordinates $x$ and $y$.

Summarizing, the basic idea behind the GOL is as follows. The properties of a 3D heterogeneous layer, which are described by the complete set of Maxwell's equations, are reduced to a set of quasi-static equations in a 2D reference plane, with the price being the introduction of coupled electric/magnetic fields and currents and dependence on one adjustable parameter, namely, the distance $l_0$ from the reference plane. Comparison of numerical calculation and the GOL approximation for metal films with periodic corrugation (Levy-Nathansohn and Bergman, 1997) indicate that the GOL results are *not* sensitive to $l_0$. The original choice $l_0 = 0.25D$

(Sarychev *et al.*, 1995) [i.e., the parameter $a = D/4d$ in Eqs. (336) and (337)] allows one to reproduce most of the computer simulations' results, except when a surface polariton is excited in the corrugated film.

To simplify (341), the system of the basic equations, the electric and magnetic fields on both sides of the film must be analyzed, namely, one must consider these fields at a distance $l_0$ behind the film, $\mathbf{E}_1(\mathbf{r}) = \mathbf{E}(\mathbf{r}, -d/2 - l_0)$, $\mathbf{H}_1(\mathbf{r} = \mathbf{H}(\mathbf{r}, -d/2 - l_0)$, and at the same distance in front of the film, $\mathbf{E}_2(\mathbf{r}) = \mathbf{E}(\mathbf{r}, d/2 + l_0)$, and $\mathbf{H}_2(\mathbf{r}) = \mathbf{H}(\mathbf{r}, d/2 + l_0)$. Then, Maxwell's second equation, $\nabla \times \mathbf{H} = (4\pi/c)\mathbf{I}$, can be written as, $\oint \mathbf{H} \, d\mathbf{l} = (4\pi/c)(\mathbf{n}_1 \cdot \mathbf{I}_E)\Delta$, where $\mathbf{n}_1$ is perpendicular to the integration contour, and the integration is over a rectangular contour which has sides $d + 2l_0$ and $\Delta$, such that the sides with length $d + 2l_0$ are perpendicular to the film and those with length $\Delta$ are in the planes $z = \pm(d/2 + l_0)$. In the limit $\Delta \to 0$ this equation takes the following form

$$\mathbf{H}_2 - \mathbf{H}_1 = -\frac{4\pi}{c}(\mathbf{n} \times \mathbf{I}_E) = -\frac{4\pi}{c}[s(\mathbf{n} \times \mathbf{E}_1) - g\mathbf{H}_1]. \qquad (343)$$

The same procedure, when applied to Maxwell's first equation, $\nabla \times \mathbf{H} = ik\mathbf{H}$, yields

$$\mathbf{E}_2 - \mathbf{E}_1 = -\frac{4\pi}{c}(\mathbf{n} \times \mathbf{I}_H) = -\frac{4\pi}{c}[m(\mathbf{n} \times \mathbf{H}_1) - g\mathbf{E}_1]. \qquad (344)$$

Now, the electric field $\mathbf{E}_1$ can be expressed, using Eq. (343), in terms of the magnetic fields $\mathbf{H}_1$ and $\mathbf{H}_2$, while the magnetic field $\mathbf{H}_1$ can be expressed, using Eq. (344), in terms of the electric fields $\mathbf{E}_1$ and $\mathbf{E}_2$. If we substitute the resulting expressions in the GOL, Eqs. (338) and (339), and use Eq. (333), we obtain

$$\mathbf{I}_E = u\mathbf{E}, \qquad \mathbf{I}_H = w\mathbf{H}, \qquad (345)$$

where $\mathbf{E} = \frac{1}{2}(\mathbf{E}_1 + \mathbf{E}_2)$, $\mathbf{H} = \frac{1}{2}(\mathbf{H}_1 + \mathbf{H}_2)$, and parameters $u$ and $w$ are given by

$$u = -\frac{c}{2\pi}\frac{g}{m}, \qquad w = -\frac{c}{2\pi}\frac{g}{s}, \qquad (346)$$

implying that the GOL is diagonalized by introducing new fields $\mathbf{E}$ and $\mathbf{H}$, such that Eqs. (345) have the same form as constitutive equations of macroscopic electro-dynamics, but with the difference that the local conductivity has been replaced by the parameter $u$ and the magnetic permeability $K$ has been replaced by $-4i\pi w/\omega$. It is then straightforward to show that, the new Ohmic parameters $u$ and $w$ can be expressed in terms of the local refractive index $\eta = \sqrt{\epsilon(\mathbf{r})}$ as

$$u = -i\frac{c}{2\pi}\frac{\tan(Dk/4) + \eta\tan(dk\eta/2)}{1 - \eta\tan(Dk/4)\tan(dk\eta/2)}, \qquad (347)$$

$$w = i\frac{c}{2\pi}\frac{\eta\tan(Dk/4) + \tan(dk\eta/2)}{\eta - \tan(Dk/4)\tan(dk\eta/2)}. \qquad (348)$$

In these equations, the refractive index $\eta$ takes on values $\eta_m = \sqrt{\epsilon_m}$ and $\eta_d = \sqrt{\epsilon_d}$ for metal and dielectric regions of the film, respectively. In the quasi-static limit, when the optical thickness of the metal grains is small, $dk|\eta_m| \ll 1$, but the metal

dielectric constant is large in magnitude, $|\epsilon_m| \gg 1$, the following estimates are obtained

$$u_m \simeq -i\frac{\omega\epsilon_m}{4\pi}d, \qquad w_m \simeq i\frac{\omega}{4\pi}\left(d + \frac{1}{2}D\right), \qquad (d/\delta \ll 1) \qquad (349)$$

for the metal grains. In the opposite limit, when the skin effect is strong, i.e., when $\delta = 1/k[\mathrm{Im}(\eta_m)] \ll d$, and the electromagnetic field does not penetrate the metal grains, we have

$$u_m = i\frac{2c^2}{\pi D\omega}, \qquad w_m = i\frac{\omega D}{8\pi}. \qquad (350)$$

In the dielectric region, when the film is thin enough that, $dk\eta_d \ll 1$ and $\epsilon_d \sim 1$, we obtain

$$u_d = -i\frac{\omega\epsilon'_d}{8\pi}D, \qquad w_d = i\frac{\omega}{4\pi}\left(d + \frac{1}{2}D\right), \qquad (351)$$

where $\epsilon'_d = 1 + 2\epsilon_d d/D$.

Potentials for the fields $\mathbf{E}_2(\mathbf{r})$ and $\mathbf{H}_2(\mathbf{r})$ may be introduced for the same reason as for $\mathbf{E}_1(\mathbf{r})$ and $\mathbf{H}_1(\mathbf{r})$. Therefore, the fields $\mathbf{E}(\mathbf{r})$ and $\mathbf{H}(\mathbf{r})$ in Eqs. (345) can in turn be represented as gradients of some potentials, $\mathbf{E} = -\nabla\phi'$, and $\mathbf{H} = -\nabla\psi'$. By substituting these expressions into Eqs. (345) and then in the continuity equation, we obtain two equations for $\phi'(\mathbf{r})$ and $\psi'(\mathbf{r})$ that can be solved independently under the conditions that, $\langle\nabla\phi'_1\rangle = \langle\mathbf{E}\rangle \equiv \mathbf{E}_0$, and $\langle\nabla\psi'_1\rangle = \langle\mathbf{H}_1\rangle \equiv \mathbf{H}_0$, where the constant fields $\mathbf{E}_0$ and $\mathbf{H}_0$ are external (given) fields that are determined by the incident wave. When $\mathbf{E}$, $\mathbf{H}$, $\mathbf{I}_E$, and $\mathbf{I}_H$ are determined from the solution of these equations, the local electric and magnetic fields in the plane $z = -l_0 - d/2$ are given by, $\mathbf{E}_1 = \mathbf{E} + (2\pi/c)(\mathbf{n}\times\mathbf{I}_H)$, and $\mathbf{H}_1 = \mathbf{H} + (2\pi/c)(\mathbf{n}\times\mathbf{I}_E)$. Note that the field $\mathbf{E}_1(\mathbf{r})$ is usually measured in a typical near field experiment. Then, the effective parameters $u_e$ and $w_e$ are defined in a way similar to Eqs. (345), viz., $\langle\mathbf{I}_E\rangle = u_e\mathbf{E}_0 \equiv \frac{1}{2}u_e(\langle\mathbf{E}_1\rangle + \langle\mathbf{E}_2\rangle)$, which, when substituted in Eqs. (343) and (344) (which are averaged over the $\{x, y\}$ coordinates), yield

$$[\mathbf{n} \times (\langle\mathbf{H}_2\rangle - \langle\mathbf{H}_1\rangle)] = \frac{2\pi}{c}u_e(\langle\mathbf{E}_1\rangle + \langle\mathbf{E}_2\rangle), \qquad (352)$$

$$[\mathbf{n} \times (\langle\mathbf{E}_2\rangle - \langle\mathbf{E}_1\rangle)] = \frac{2\pi}{c}w_e(\langle\mathbf{H}_1\rangle + \langle\mathbf{H}_2\rangle). \qquad (353)$$

Suppose now that the wave enters the film from the right-half space, such that its amplitude is proportional to $\exp(-ikz)$. The incident wave is partially reflected and partially transmitted through the film. The electric field amplitude in the right-half space, away from the film, can be written as $\mathbf{e}[\exp(-ikz) + r\exp(ikz)]$, where $r$ is the reflection coefficient and $\mathbf{e}$ is the polarization vector. Then, the electric component of the electromagnetic wave well behind the film will be $\mathbf{e}[t\exp(-ikz)]$, where $t$ is the transmission coefficient. We assume for simplicity that the film has no optical activity, which means that the wave polarization $\mathbf{e}$ remains the same before and after the film. At the planes $z = d/2 + l_0$ and $z = -d/2 - l_0$ the average electric field is $\langle\mathbf{E}_2\rangle$ and $\langle\mathbf{E}_1\rangle$, respectively. On the other hand, the wave

away from the film is matched with the average fields in the planes $z = d/2 + l_0$ and $z = -d/2 - l_0$, i.e., $\langle \mathbf{E}_2 \rangle = \mathbf{e} \{\exp[-ik(d/2 + l_0)] + r \exp[ik(d/2 + l_0)]\}$ and $\langle \mathbf{E}_1 \rangle = \mathbf{e}\{t \exp[ik(d/2 + l_0)]\}$. The same matching, but with the magnetic fields, yields, $\langle \mathbf{H}_2 \rangle = (\mathbf{n} \times \mathbf{e}) \{-\exp[-ik(d/2 + l_0)] + r \exp[ik(d/2 + l_0)]\}$ and $\langle \mathbf{H}_1 \rangle = -(\mathbf{n} \times \mathbf{e})t \exp[ik(d/2 + l_0)]$ in the planes $z = d/2 + l_0$ and $z = -d/2 - l_0$, respectively. Substitution of these expressions for $\langle \mathbf{E}_1 \rangle$, $\langle \mathbf{E}_2 \rangle$, $\langle \mathbf{H}_1 \rangle$, and $\langle \mathbf{H}_2 \rangle$ in Eqs. (352) and (353) yields two scalar, linear equations for reflection ($r$) and transmission ($t$) coefficients, the solution of which yields the reflectance,

$$\mathcal{R} \equiv |r|^2 = \left| \frac{\frac{2\pi}{c} (u_e + w_e)}{\left(1 + \frac{2\pi}{c} u_e\right) \left(1 - \frac{2\pi}{c} w_e\right)} \right|^2, \qquad (354)$$

the transmittance

$$\mathcal{T} \equiv |t|^2 = \left| \frac{1 + \left(\frac{2\pi}{c}\right)^2 u_e w_e}{\left(1 + \frac{2\pi}{c} u_e\right) \left(1 - \frac{2\pi}{c} w_e\right)} \right|^2, \qquad (355)$$

and the absorbance

$$\alpha = 1 - \mathcal{T} - \mathcal{R} \qquad (356)$$

of the film. Therefore, the effective Ohmic parameters $u_e$ and $w_e$ determine completely the optical properties of heterogeneous films. This analysis indicates that, the problem of the field distribution and optical properties of the metal-dielectric films reduces to uncoupled quasi-static conductivity problems for which extensive theoretical analyses have already been carried out. Numerous analytical as well as numerical methods, developed for heterogeneous media with percolation disorder (see Chapters 4–6 of Volume I), can be employed for determining the effective parameters $u_e$ and $w_e$ of the film.

We can now consider the case of strong skin effect in the metal grains and study the evolution of the optical properties of a semi-continuous metal film, as the volume fraction $p$ of the metal is increasing. When $p = 0$, the film is purely dielectric and $u_e = u_d$ and $w_e = w_d$, where $u_d$ and $w_d$ are the dielectric Ohmic parameters given by Eqs. (351). If we substitute $u_e = u_d$ and $w_e = w_d$ in Eqs. (354)–(356) and assume that the dielectric film has no losses and is optically thin (i.e., $dk\epsilon_d \ll 1$), we obtain the reflectance $\mathcal{R} = d^2(\epsilon_d - 1)^2 k^2/4$, the transmittance $\mathcal{T} = 1 - d^2(\epsilon_d - 1)^2 k^2/4$, and the absorbance $\alpha = 0$, well-known results for a thin dielectric film (see, for example, Jackson, 1998). The losses are also absent in the limit of full coverage, i.e., when the metal volume fraction $p = 1$. Indeed, substituting the Ohmic parameters $u_e = u_m$ and $w_e = w_m$ from Eqs. (350) in Eqs. (354)–(356) yields, $\mathcal{R} = 1$, $\mathcal{T} = 0$, and $\alpha = 0$. Note that in the limits $p = 0$ and $p = 1$, the optical properties of the film do not depend on the particle size $D$, because properties of the dielectric and continuous metal films should not depend on the shape of the metal grains.

Next, we consider the film at $p = p_c$ with $p_c = 1/2$ for a self-dual system. A semi-continuous metal film may be thought of as a mirror, which is broken into small pieces with typical size $D \ll \lambda$. At $p_c$, the exact equations (see Sections 3.1 and 3.3, and also Chapters 4 and 5 of Volume I), $u_e = \sqrt{u_d u_m}$ and $w_e = \sqrt{w_d w_m}$, which result from the exact duality relation, hold. Thus,

$$\frac{2\pi}{c} u_e(p_c) = \sqrt{\epsilon'_d}, \qquad \frac{2\pi}{c} w_e(p_c) = i\frac{Dk}{4}\sqrt{1 + \frac{2d}{D}}, \qquad (357)$$

from which it follows that $|w_e/u_e| \sim Dk \ll 1$, and hence, compared with $u_e$, $w_e$ can be neglected. Under this condition, we obtain

$$\mathcal{R}(p_c) = \frac{\epsilon'_d}{\left(1 + \sqrt{\epsilon'_d}\right)^2}, \qquad (358)$$

$$\mathcal{T}(p_c) = \frac{1}{\left(1 + \sqrt{\epsilon'_d}\right)^2}, \qquad (359)$$

$$\alpha(p_c) = \frac{2\sqrt{\epsilon'_d}}{\left(1 + \sqrt{\epsilon'_d}\right)^2}, \qquad (360)$$

where, $\epsilon'_d = 1 + 2\epsilon_d d/D$, as before. When metal grains are oblate enough that $\epsilon_d d/D \ll 1$ and $\epsilon'_d \to 1$, one obtains the universal result

$$\mathcal{R} = \mathcal{T} = 1/4, \qquad \alpha = 1/2, \qquad (361)$$

implying that there is effective absorption in semi-continuous metal films even for the case when neither dielectric nor metal grains absorb light energy. The effective absorption in a loss-free film means that the electromagnetic energy is stored in the system, and that the amplitudes of the local electromagnetic field can diverge. In practice, due to non-zero losses, the local field saturates in any semi-continuous metal film.

To determine the optical properties of semi-continuous films for arbitrary metal volume fraction $p$, the EMA can be used which then yields the following equations,

$$u_e^2 - \Delta p u_e(u_m - u_d) - u_d u_m = 0, \qquad (362)$$

$$w_e^2 - \Delta p w_e(w_m - w_d) - w_d w_m = 0, \qquad (363)$$

where $\Delta p = (p - p_c)/p_c$ ($p_c = 1/2$). Equation (363) indicates that, when the skin effect is strong and $w_m$ and $w_d$ are given by Eqs. (350) and (351), then $|w_e| \ll c$ for *all* metal volume fractions $p$, and therefore we can neglect $w_e$ in Eqs. (354) and (355). Moreover, compared with $u_m$, $u_d$ can also be neglected in the second term of Eq. (362). Thus, introducing the dimensionless Ohmic parameter $u'_e = (2\pi/c)u_e$ allows us to rewrite Eq. (362) as

$$u_e'^2 - 2i\frac{\lambda \Delta p}{\pi D} u'_e - \epsilon'_d = 0. \qquad (364)$$

At $p = p_c = 1/2$ (i.e., where $\Delta p = 0$), Eq. (364) predicts that, $u'_e(p_c) = \sqrt{\epsilon'_d}$, which coincides with the exact result, Eq. (357), and those given by Eqs. (358)–(360). For $p \neq p_c$, Eq. (364) predicts that

$$u'_e = i\frac{\lambda \Delta p}{\pi D} + \sqrt{\epsilon'_d - \left(\frac{\lambda \Delta p}{\pi D}\right)^2},$$   (365)

which becomes purely imaginary for $|\Delta p| > \pi D\sqrt{\epsilon'_d}/\lambda$. Then, $\alpha = 1 - |u'_e|^2/|1 + u'_e|^2 - 1/|1 + u'_e|^2 = 0$ (recall that $w_e$ was neglected). In the vicinity of $p_c$, namely, for $|\Delta p| < (\pi D/\lambda)\sqrt{\epsilon'_d}$, the effective Ohmic parameter $u'_e$ has a non-vanishing real part, and therefore

$$\alpha = \frac{2\sqrt{\epsilon'_d - [\lambda \Delta p/(\pi D)]^2}}{1 + \epsilon'_d + 2\sqrt{\epsilon'_d - [\lambda \Delta p/(\pi D)]^2}},$$   (366)

which has a well-defined maximum at $p_c$, with the width of the maximum being inversely proportional to the wavelength. These predictions were confirmed by extensive numerical simulations. They are also in agreement with the experimental data (see Sarychev and Shalaev, 2000, for detailed discussions). Note that the parameters $u_e$ and $w_e$ can be determined experimentally by measuring the amplitudes and phases of the transmitted and reflected waves using, for example, a waveguide technique (see, for example, Golosovsky et al., 1993 and references therein), or by measuring the film reflectance as a function of the fields $\mathbf{E}_1$ and $\mathbf{H}_1$.

## 3.8    Piecewise Linear Transport Processes

The last nonlinear transport process that we describe and analyze is not caused by strong morphological disorder and its interplay with a transport process, rather it has to do with the constitutive relation between the current and the potential gradient, augmented by a threshold in the potential gradient. Such nonlinear transport phenomenon are typically piecewise linear, or possibly nonlinear, and are characterized by at least one threshold. Several possible $I - V$ characteristics of such materials are shown in Figure 3.14. Because of the threshold, of course, even a piecewise linear transport is in fact a highly nonlinear process. In many cases, the regime below the threshold is degenerate in the sense that, nothing interesting happens if the driving potential applied to the material is below its threshold. The applications of this type of nonlinear transport process are numerous. For example, bipolar Zener diodes (which are commercially called *varistors*) switch from being a non-conducting link to a conducting one at an onset voltage threshold $v_c$. More generally, a network of such diodes can become conducting only if the voltage applied to it is larger than a critical value $V_c$. In brittle fracture, which will be studied in Chapters 6–8, no microcrack nucleation and propagation take place unless the external stress or strain applied to a solid material exceeds a critical value

FIGURE 3.14. Twelve types of physically realizable nonlinear $I - V$ characteristics, seven of which are also characterized by a threshold (after Sahimi, 1998).

which depends on the size of the sample. Bingham fluids are viscous and behave like Newtonian fluids if the shear stress applied to them is larger than a critical value, but do not flow if the stress is less than the threshold value. An example of such fluids, already described in Section 9.3 of Volume I, is foam. In order to mobilize foam and force it to flow, the applied pressure must exceed a critical value; otherwise it will not flow.

Let us consider a 2D or 3D resistor network in which every bond is characterized by the following current-voltage relation,

$$i = \begin{cases} g(v - v_c)^n, & v > v_c, \\ 0 & v \le v_c, \end{cases} \tag{367}$$

where $v_c$ is the critical voltage or threshold for the onset of transport. As in the case of strong and weak nonlinearities, we take $g$ to be a generalized bond conductance which, in general, can vary from bond to bond. On the other hand, in any physical situation involving a disordered material, one expects a distribution of the thresholds $v_c$, because due to a variety of factors, different parts of a material may become conductive beyond different thresholds. Therefore, one may make the simplification that, instead of assuming $g$ to be a statistically-distributed variable, $v_c$ is assumed to be randomly distributed which, for the sake of simplicity, is assumed to be distributed uniformly in $(0, 1)$. The conductivity $g$ is then the same for all bonds, and therefore its numerical value is irrelevant (we assume $g = 1$). The questions that we ask are:

(1) What is the critical voltage $V_c$ in order to have macroscopic transport in the network, and

(2) how do the macroscopic current $I$ and the effective conductivity $g_e$ of the network vary with the applied voltage? The piecewise linear process that we study here is *reversible*, i.e., if $I$ is lowered the conducting bonds become

insulating again. This is an important assumption since, if we assume that the process is irreversible, then converting even one insulating bond to a conducting one triggers an *avalanche effect*: The conversion of the first bond makes consecutive conversions easier. Such irreversible and nonlinear models have been used to model brittle fracture and electrical and dielectric breakdown of disordered materials, which will be discussed in Chapters 5–8.

It is clear that for any applied voltage $V$ less than a critical threshold $V_c$ no macroscopic current can flow. Therefore, it should also be clear that

$$V_c = min\left(\sum_i v_{ci}\right), \tag{368}$$

where $v_{ci}$ is the critical voltage of bond $i$, and the sum is taken over all paths between the two terminals of the network. Equation (368) immediately necessitates the concept of an *optimal path* between the two terminals of the network (see, for example, Cieplak *et al.*, 1994,1996; Porto *et al.*, 1997). Obviously, if the applied voltage is larger than a final or the last voltage threshold $V_l$, all bonds of the network will be conducting, and one will have the usual linear transport in which the current $I$ is simply proportional to $V$. Therefore, one generally has *three* regimes of interest:

(1) If $V < V_c$, then enough bonds have not become conducting to form a sample-spanning cluster, and therefore no macroscopic transport takes place. Hence, $I = 0$ and $g_e$=0.
(2) If $V_c < V < V_l$, then enough bonds have become conducting that make macroscopic transport possible, while some of the bonds are still not conducting. We expect $I$ to depend nonlinearly on $V - V_c$, because this is precisely the regime in which the effect of nonlinearity (random voltage thresholds) should manifest itself. As we show below, this is indeed the case (note that in linear transport above $p_c$, $I$ *always* varies linearly with $V$).
(3) If $V > V_l$, then every bond of the network is conducting, the normalized effective conductivity is $g_e = 1$, and $I$ depends linearly on $V$ again.

## 3.8.1 Computer Simulation

Computer simulation of this problem, even for $n = 1$, is difficult, and thus deserves to be discussed here. At the beginning of the simulations, one distributes the critical thresholds $v_c$ of the bonds and applies a large enough external voltage to the network, such that every bond becomes conducting (i.e., the voltage across it exceeds its critical voltage). The external voltage is then decreased gradually, and the nodal voltage distribution and hence the current distribution in the bonds are computed. As a result of lowering the applied voltage, some of the conducting bonds become insulating (since the voltage across them will be less than their critical voltage). The new voltage and current distributions are calculated, the newly-insulating bonds are identified, and so on.

### 3.8.2  Scaling Properties

Roux and Herrmann (1987) used accurate numerical simulations, and Gilabert *et al.* (1987) utilized an analogue resistor network, to find that in 2D and near $V_c$,

$$I \sim (V - V_c)^{\delta}, \tag{369}$$

with $\delta \simeq 2 \pm 0.08$. The power-law (369) is the only scaling property of piecewise linear transport that has been studied so far.

### 3.8.3  Effective-Medium Approximation

We now describe the predictions of an EMA for piecewise linear transport and compare its predictions with simulation results. We consider only the case $n = 1$ and present the final results; complete details are given by Sahimi (1993a). Suppose that $p$ is the fraction of the bonds that have become conducting. Then, in the non-conducting regime, i.e., before a sample-spanning conducting path has formed between two opposite faces of the network and $p < p_c = 2/Z$ (recall from Sections 3.1 and 3.2 that, since the problem is treated within an EMA, the percolation threshold is $p_c = 2/Z$), the applied voltage $V$ varies with $p$ according to

$$V = p - \frac{1}{2}p^2, \qquad p < 2/Z. \tag{370}$$

Equation (370) predicts how the applied voltage $V$ varies with $p$ *before* a sample-spanning conducting path is formed. At $p = p_c = 2/Z$ the first sample-spanning conducting path is formed and therefore

$$V_c = \frac{2}{Z} - \frac{2}{Z^2}, \tag{371}$$

which is obtained by substituting $p = p_c = 2/Z$ in Eq. (370). For $p > 2/Z$ we have a conducting system for which

$$V = \frac{Z - 2}{Z}p + \frac{2}{Z^2}, \qquad p \geq \frac{2}{Z}. \tag{372}$$

At $p = 1$ all the bonds are conducting, so that the corresponding last voltage for converting the last bond to a conducting bond is given by

$$V_l = \frac{Z - 2}{Z} + \frac{2}{Z^2}. \tag{373}$$

The corresponding equations for the effective conductivity $g_e$ of the network are as follows. Clearly, for $V < V_c$ we must have $g_e = 0$. For $V_c \leq V \leq V_l$ we have

$$g_e = \frac{Z^2}{(Z-2)^2}(V - \frac{2}{Z^2}) - \frac{2}{Z-2}, \qquad V_c \leq V \leq V_l. \tag{374}$$

Obviously, $g_e = 1$ for $V \geq V_l$.

We can also determine the macroscopic $I - V$ characteristic of the material. For $V < V_c$ there is no macroscopic transport and therefore, $I = 0$. For $V_c \leq V \leq V_l$ we have

$$I = \frac{Z^2}{2(Z-2)^2}(V - \frac{2}{Z^2})^2 - \frac{2}{Z-2}(V - \frac{2}{Z^2}) + \frac{2}{Z^2}, \qquad V_c \leq V \leq V_l. \quad (375)$$

For $V \geq V_l$, we have $g_e = 1$, and therefore the current $I$ is related to the applied voltage through a simple equation

$$I = V - \frac{1}{2}, \qquad (376)$$

independent of $Z$. Thus, the EMA predicts correctly the existence of the three transport regimes discussed above and, in particular, it predicts that for $V_c \leq V \leq V_l$, $I$ depends quadratically on $V - V_c$, where $V_c = 2/Z - 2/Z^2$.

Figure 3.15 presents the dependence of $g_e$ on the applied voltage $V$ in a square network. All the qualitative features of the transport process are correctly predicted by the EMA, except that the numerical simulations indicate smooth variations of $g_e$ with $V$, whereas the EMA predicts a sharp, discontinuous, transition at $V = V_l$. Figure 3.16 presents the variations of the macroscopic current $I$ with the applied voltage $V$ in the same system and, unlike $g_e$, both the numerical calculations and the EMA predict no sharp transition at $V = V_l$. However, the numerical value of the critical voltage $V_c$ does not agree with the prediction of the



FIGURE 3.15. Effective conductivity of the square network with piecewise linear resistors with a threshold, versus the applied voltage (after Sahimi, 1993a).

FIGURE 3.16. The $I - V$ characteristics of the square network of Figure 3.15 (after Sahimi, 1993a).

EMA. While computer simulations indicate that, $V_c \simeq 0.29$, the EMA predicts that, $V_c = 3/8 = 0.375$. Roux *et al.* (1987) used a transfer-matrix method described in Section 5.14.2 of Volume I and estimated that for a square network, *tilted at 45°*, one has, $V_c \simeq 0.23$ [in general, for the square network, $V_c$(tilted)=$V_c$(non-tilted)/$\sqrt{2}$]. This difference can be explained by the fact that, because the resistor network that Roux *et al.* (1987) used in their simulations was tilted, their network is different from a non-tilted one, since the distribution of currents in the bonds of their network is isotropic, whereas the same distribution is anisotropic in a non-tilted network. The difference is due to the fact that the bonds of a non-tilted network that are perpendicular to the direction of the applied voltage receive much less current than those that are aligned with it. As a result, formation of a sample-spanning conducting cluster is easier in a tilted network than in a non-tilted one, implying that the critical voltage $V_c$ for a tilted network should be smaller than that of a non-tilted one. Thus, such local anisotropies, which usually have no consequence for macroscopic properties of linear transport processes, are important in a nonlinear system, such as what is described here. Moreover, according to Eq. (375), in the nonlinear regime, the macroscopic current $I$ varies quadratically with $V - V_c$, which is in agreement with the simulations of Roux and Herrmann (1987), Eq. (369).

# Summary

Using the discrete models, we described and analyzed several types of nonlinear transport and optical properties of disordered materials. As our analyses indicate, the interplay of nonlinearity and the disordered morphology of a material gives rise to a rich set of phenomena that are absent in linear transport processes in the same material. In particular, strong heterogeneity, such as percolation-type disorder, enhances the nonlinear response of a material, and shrinks the range of the parameter space in which the material behaves linearly, and hence opens up the possibility of developing composite materials with highly unusual and useful properties.

# 4
# Nonlinear Rigidity and Elastic Moduli: The Continuum Approach

## 4.0   Introduction

In this chapter we consider nonlinear mechanical properties of heterogeneous materials. This class of problems has many applications that will be described throughout this chapter. However, to give the reader an interesting and somewhat unusual application of this class of phenomena, we consider the following problem. It has been observed (Gordon, 1978) that extensible biological tissues, such as skin, are difficult to tear, even though their specific work of fracture (see the discussions in Chapters 6 and 7) is not large compared to those of materials that tear easily. For example, the fracture toughness of animal membranes is around 1-10 kJm$^{-2}$, an order of magnitude smaller than aluminum foil which tears easily. Gordon reasoned that this difference is due to the markedly different shape of the stress-strain diagram of such materials. Figure 4.1 presents schematic stress-strain curves for extensible biological tissues, rubber, and the standard Hookean solid for which the diagram is a straight line. The small-strain portion of the J-shaped curve of the biological material is indicative of lack of shear connection in the material, i.e., absence of shear stiffness in what are anisotropic solids. This diagram provides an explanation as to why such materials are difficult to tear, because it is difficult to concentrate energy into the path of a putative crack. Note also the difference between the stress-strain diagrams for rubbers and the biological materials: For small strains, the rubber's curve is not J-shaped, which may also explain why we cannot replace human body arteries or veins by rubber tubes. We also remind the reader that when Nature does want fracture and tear to happen, as in, for example, amniotic membranes and eggshells, the stress-strain diagrams are Hookean linear elastic!

Studies of heterogeneous materials with nonlinear constitutive behavior go back to at least Taylor (1938) who studied the plasticity of polycrystals, and to the subsequent work by Bishop and Hill (1951a,b) and Drucker (1959) who investigated the behavior of ideally plastic polycrystals and composite materials. Over the past decade or so, numerical simulations of nonlinear materials with periodic microstructures have been carried out (see, for example, Christman *et al.*, 1989; Tvergaard, 1990; Bao *et al.*, 1991), as well as materials with more general microstructures (see, for example, Brokenborough *et al.*, 1991; Moulinec and Suquet, 1995). Such efforts will be briefly described in this chapter where we make comparison between the theoretical predictions and the numerical simulation results.

FIGURE 4.1. Schematic representation of different stress-strain relations.



The main advantage of such simulations is that they provide accurate description of the system under study, and yield useful insight into their properties. Their main disadvantage is that they require very intensive computations, especially when the material's microstructure is disordered.

In this chapter we describe and discuss recent advances in understanding the effective mechanical properties of disordered materials with constitutive nonlinearity. Although one may argue that numerical techniques, such as the finite-element methods, represent some form of discrete approach to this class of problems, to our knowledge very little work has been done using the discrete network models of the type that we have so far described and discussed for estimating various transport properties of disordered materials. Therefore, the main focus of this chapter is on the theoretical developments based on nonlinear continuum models of disordered materials. These theoretical approaches represent the mechanical analogues of those described in Chapter 2 for estimating the effective conductivity and dielectric constant of nonlinear materials. Thus, the methods that we describe in this chapter are based on rigorous variational principles which, in addition to possessing mathematical rigor, have the advantage of leading to bounds and relatively accurate estimates for the mechanical properties. As described and discussed in Chapter 2, such variational principles allow one to obtain estimates of the effective energy densities of nonlinear materials in terms of the corresponding information for linear composites with the same microstructure. A large portion of our analyses and discussions in this chapter is based on an excellent review by Ponte Castañeda and Suquet (1998).

## 4.1  Constitutive Relations and Potentials

Similar to Chapter 2, where we analyzed the effective nonlinear conductivity and dielectric constant of disordered materials, we also assume in the present chapter that the constitutive behavior of the individual phases of the material is governed by a potential, or strain-energy function, $w(\epsilon)$, in such a way that the (infinitesimal)

strain $\epsilon$ and stress $\sigma$ fields are related by

$$\sigma = \frac{\partial w}{\partial \epsilon}. \tag{1}$$

Although Eq. (1) is intended for nonlinear elastic behavior of materials in the limit of small strains, by interpreting $\epsilon$ and $\sigma$ as the Eulerian strain rate and Cauchy stress, it can also be used for modeling finite viscous deformations. Assuming then that $w$ is a convex function of $\epsilon$, Eq. (1) is inverted with aid of the Legendre transformation:

$$u(\sigma) = \sup_{\epsilon}\{\sigma : \tau - w(\epsilon)\}. \tag{2}$$

Equation (2) defines a convex stress-energy function $u$, such that

$$\epsilon = \frac{\partial u}{\partial \epsilon}. \tag{3}$$

The functions $w$ and $u$ are dual potentials and are related by the classical reciprocity relations. As in Chapter 2, the notation $u = w^*$ is used to express the relation between these two quantities.

For isotropic materials, general forms of $w$ and $u$ are given by

$$w(\epsilon) = \frac{9}{2}K\epsilon_m^2 + \varphi(\epsilon_{eq}), \tag{4}$$

and

$$u(\sigma) = \frac{1}{2K}\sigma_m^2 + \psi(\sigma_{eq}), \tag{5}$$

where $\varphi$ and $\psi$ are dual convex potentials, $\sigma_m$ and $\epsilon_m$ are the hydrostatic stress and strain given by

$$\sigma_m = \frac{1}{3}\text{tr}(\sigma) ; \quad \epsilon_m = \frac{1}{3}\text{tr}(\epsilon), \tag{6}$$

and $\sigma_{eq}$ and $\epsilon_{eq}$ are the Von Mises equivalent stress and strain,

$$\sigma_{eq} = \left(\frac{3}{2}\sigma_d : \sigma_d\right)^{1/2}, \quad \epsilon_{eq} = \left(\frac{2}{3}\epsilon_d : \epsilon_d\right)^{1/2}, \tag{7}$$

with $\sigma_d$ and $\epsilon_d$ being the stress and strain deviators (see also Chapter 7 of Volume I) given by, $\sigma_d = \sigma - \sigma_m\mathbf{U}$ and $\epsilon_d = \epsilon - \epsilon_m\mathbf{U}$. Thus, one can write

$$\sigma_m = 3K\epsilon_m, \quad \sigma_d = 2\mu(\epsilon_{eq})\epsilon_d, \tag{8}$$

with

$$\mu(\epsilon_{eq}) = \frac{1}{3}\frac{\sigma_{eq}}{\epsilon_{eq}} = \frac{1}{3}\frac{\varphi'(\epsilon_{eq})}{\epsilon_{eq}} = \frac{1}{3}\frac{\sigma_{eq}}{\psi'(\epsilon_{eq})}. \tag{9}$$

Therefore, each phase is assumed to be linear for purely hydrostatic loadings, characterized by a constant bulk modulus $K$ and nonlinear in shear, characterized by a strain-dependent shear modulus $\mu$.

Consider, as an example, high temperature creep of metals, which is commonly characterized in terms of a power-law constitutive relation. If we neglect elastic effects and assume incompressibility, then, the dissipation $\varphi$ and stress potential $\psi$ of the material are given by

$$\varphi(\epsilon_{eq}) = \frac{\sigma^0 \epsilon^0}{m+1} \left(\frac{\epsilon_{eq}}{\epsilon^0}\right)^{m+1}, \quad \psi(\sigma_{eq}) = \frac{\sigma^0 \epsilon^0}{n+1} \left(\frac{\sigma_{eq}}{\sigma^0}\right)^{n+1}, \tag{10}$$

where $\epsilon^0$ and $\sigma^0$ denote a reference strain rate and stress, respectively, $m$ and $n$ are two exponents such that $m = 1/n$, and $\epsilon_m = 0$. For example, for Newtonian viscous materials, $n = m = 1$, where $\eta = \sigma^0/3$ is the viscosity, while the Von Mises rigid, ideally plastic materials correspond to the limit $m \to 0$ ($n \to \infty$), where $\sigma^0$ now denotes the flow stress in tension. In the latter case, the stress potential becomes unbounded for stresses that exceed $\sigma^0$. It is then useful to introduce the strength domain $P$, defined by the set

$$P = \{\boldsymbol{\sigma} : \sigma_{eq} \leq \sigma^0\}. \tag{11}$$

The creep of crystalline materials can also be described within this framework. We consider a single crystal that undergoes creep on a set of $M$ preferred crystallographic slip systems, and is characterized by the second-order tensors $\boldsymbol{\mu}_i$, $i = 1, \ldots, M$, defined by

$$\boldsymbol{\mu}_i = \frac{1}{2}(\mathbf{n}_i \otimes \mathbf{m}_i + \mathbf{m}_i \otimes \mathbf{n}_i), \tag{12}$$

where $\mathbf{n}_i$ and $\mathbf{m}_i$ are the unit vectors normal to the slip plane and along the slip direction in the $i$th system, respectively, and $\otimes$ denotes the tensorial product of two vectors. If a stress $\boldsymbol{\sigma}$ is applied to the crystal, then, the resulting shear stress acting on the $i$th slip system is given by

$$\tau_i = \boldsymbol{\sigma} : \boldsymbol{\mu}_i, \tag{13}$$

while the strain rate $\boldsymbol{\epsilon}$ in the crystal is the superposition of the strain rates on each slip system,

$$\boldsymbol{\epsilon} = \sum_{i=1}^{M} \gamma_i \boldsymbol{\mu}_i, \tag{14}$$

where $\gamma_i$ is the shear strain rate acting on the $i$th system, which is given by

$$\gamma_i = \frac{\partial \psi_i}{\partial \tau_i}, \tag{15}$$

with the functions $\psi_i$ being convex. An equation commonly used for $\psi_i$ is

$$\psi_i(\tau) = \frac{\gamma^0 \tau_i^0}{n_i + 1} \left(\frac{|\tau|}{\tau_i^0}\right)^{n_i + 1}, \tag{16}$$

with $n_i \geq 1$ and $\tau_i^0$ being the creep exponent and reference stress of the $i$th slip system, respectively, and $\gamma^0$ is a reference strain rate. The constitutive relations

(14) and (15) can then be expressed in terms of the convex potential for the crystal:

$$u^{(c)}(\boldsymbol{\sigma}) = \sum_{i=1}^{M} \psi_i(\boldsymbol{\sigma} : \boldsymbol{\mu}_i), \tag{17}$$

such that

$$\boldsymbol{\epsilon} = \frac{\partial u^{(c)}}{\partial \boldsymbol{\sigma}}. \tag{18}$$

The limit $n_i \to \infty$ corresponds to a rigid, ideally plastic crystal, with a strength domain given by

$$P = \{\boldsymbol{\sigma}, \ \tau_i \leq \tau_i^0, \ i = 1, \ldots, M\}. \tag{19}$$

We can assume, more generally, that the potential $w$ can be expressed by

$$w(\boldsymbol{\epsilon}) = F(\mathcal{E}), \tag{20}$$

where $F$ is an appropriately-selected function, and $\mathcal{E}$ is a fourth-rank tensor which is defined by

$$\mathcal{E} = \frac{1}{2}\boldsymbol{\epsilon} \otimes \boldsymbol{\epsilon}, \tag{21}$$

and possesses the usual diagonal symmetry and positive-definitiveness property of an elasticity tensor. The function $F$ is then defined on the space of fourth-rank tensors $\mathcal{P}$ that have diagonal symmetry, so that the constitutive relation (1) can be written as

$$\boldsymbol{\sigma} = \mathcal{L}_s(\mathcal{E}) : \boldsymbol{\epsilon}, \tag{22}$$

with

$$\mathcal{L}_s(\mathcal{E}) = \frac{\partial F}{\partial \mathcal{P}}, \tag{23}$$

being the secant modulus tensor of the material, which also has diagonal symmetry. Given Eq. (20), the dual potential $u$ can be expressed as

$$u(\boldsymbol{\sigma}) = G(\mathbf{S}), \quad \mathbf{S} = \frac{1}{2}\boldsymbol{\sigma} \otimes \boldsymbol{\sigma}, \tag{24}$$

where $G$ is a function of fourth-rank tensors $\mathbf{S}$. In terms of the secant compliance tensor of the material, the constitutive relation (3) may be expressed in the following form

$$\boldsymbol{\epsilon} = \mathcal{M}_s(\mathbf{S}) : \boldsymbol{\sigma}, \quad \mathcal{M}_s(\mathcal{S}) = \frac{\partial G}{\partial \mathbf{S}}. \tag{25}$$

As an example, consider crystalline materials. First, note that

$$\tau_i^2 = 2\mathcal{M}_i :: \mathbf{S}, \quad \mathcal{M}_i = \boldsymbol{\mu}_i \otimes \boldsymbol{\mu}_i, \tag{26}$$

so that

$$u(\boldsymbol{\sigma}) = \sum_{i=1}^{M} g_i(2\mathcal{M}_i :: \mathbf{S}) = G(\mathbf{S}), \quad g_i(x) = \psi_i(\sqrt{x}), \tag{27}$$

and the compliance tensor is given by

$$\mathcal{M}_s(\mathbf{S}) = 2\sum_{i=1}^{M} \alpha_i \mathcal{M}_i, \quad \alpha_i = g_i'(\tau_i^2). \tag{28}$$

## 4.2 Formulation of the Problem

We now consider a representative volume element $\Omega$ of a heterogeneous material, such that the size of its heterogeneities is small compared to $\Omega$. The material consists of $N$ homogeneous phases $\Omega_i$, $i = 1, \ldots, N$, the distribution of which is defined by indicator functions $m_i(\mathbf{x})$, which are 1 when $\mathbf{x}$ belongs to the phase $i$, and zero otherwise. One can define two spatial averages, one over $\Omega$ and another one over $\Omega_i$, so that, for example

$$\langle \boldsymbol{\epsilon} \rangle_i = \frac{1}{|\Omega_i|} \int_{\Omega_i} \boldsymbol{\epsilon}(\mathbf{x}) d\mathbf{x}, \tag{29}$$

$$\langle \boldsymbol{\epsilon} \rangle = \frac{1}{|\Omega|} \int_{\Omega} \boldsymbol{\epsilon}(\mathbf{x}) d\mathbf{x} = \sum_{i=1}^{N} \phi_i \langle \boldsymbol{\epsilon} \rangle_i, \tag{30}$$

where $\phi_i$ is the volume fraction of phase $i$. All the phases are assumed to be homogeneous with potentials $w_i$ and $u_i$, and to be perfectly bonded at the interfaces. The total potentials $w$ and $u$ are then given by

$$w(\mathbf{x}, \boldsymbol{\epsilon}) = \sum_{i=1}^{N} m_i(\mathbf{x}) w_i(\boldsymbol{\epsilon}), \quad u(\mathbf{x}, \boldsymbol{\sigma}) = \sum_{i=1}^{N} m_i(\mathbf{x}) u_i(\boldsymbol{\sigma}). \tag{31}$$

As an example, consider a polycrystalline material, which we regard it as an aggregate of a large number of identical single crystals with different orientations, so that it can be treated as a composite, where phase $i$ is defined as the region occupied by all grains of a given orientation, relative to a reference crystal with potential $u^{(c)}$ given by (17). If $\mathbf{Q}_i$ denotes the rotation tensor that defines the orientation of phase $i$, the corresponding potential $u_i$ is given by

$$u_i(\boldsymbol{\sigma}) = u^{(c)}\left(\mathbf{Q}_i^{\mathrm{T}} \cdot \boldsymbol{\sigma} \cdot \mathbf{Q}_i\right) = \sum_{k=1}^{M} \psi_k\left[\tau_i^{(k)}\right], \tag{32}$$

where

$$\tau_i^{(k)} = \boldsymbol{\sigma} : \boldsymbol{\mu}_i^{(k)}, \quad \boldsymbol{\mu}_i^{(k)} = \mathbf{Q}_i^{\mathrm{T}} \cdot \boldsymbol{\mu}_k \cdot \mathbf{Q}_i. \tag{33}$$

The microscopic problem is one in which the local stress and strain fields within $\Omega$ solve a local problem that consists of the constitutive relation (1), the com-

patibility conditions satisfied by $\epsilon$, and the usual equilibrium equations satisfied by $\sigma$:

$$\sigma = \frac{\partial w}{\partial \epsilon}, \quad \epsilon = \frac{1}{2}\left[\nabla \mathbf{u} + (\nabla \mathbf{u})^{\mathrm{T}}\right], \quad \nabla \cdot \sigma = \mathbf{0}, \tag{34}$$

subject to one of the two classes of boundary conditions on $\partial \Omega$. One is in terms of affine displacements,

$$\mathbf{u}(\mathbf{x}) = \mathbf{E} \cdot \mathbf{x}, \tag{35}$$

while the second one is in terms of uniform traction,

$$\sigma(\mathbf{x}) \cdot \mathbf{n}(\mathbf{x}) = \Sigma \cdot \mathbf{n}(\mathbf{x}). \tag{36}$$

Here $\mathbf{E}$ and $\Sigma$ are the averages of the local strain and stress fields:

$$\mathbf{E} = \langle \epsilon \rangle = \sum_{i=1}^{N} \phi_i \langle \epsilon \rangle_i, \tag{37}$$

$$\Sigma = \langle \sigma \rangle = \sum_{i=1}^{N} \phi_i \langle \sigma \rangle_i, \tag{38}$$

and satisfy (Hill, 1963)

$$\Sigma : \mathbf{E} = \langle \sigma : \epsilon(\mathbf{u}) \rangle. \tag{39}$$

## 4.3    The Classical Variational Principles

As discussed in Chapter 2, and also Chapters 4 and 7 of Volume I, the solutions $\mathbf{u}$ and $\sigma$ of the local problem can be given two equivalent variational representations. One is in terms of the minimum potential energy which states that $\mathbf{u}$ is the solution of the problem

$$\inf_{V \in S_1(\mathbf{E})} \langle w[\epsilon(\mathbf{v})] \rangle, \tag{40}$$

where

$$S_1(\mathbf{E}) = \{\mathbf{v} = \mathbf{E} \cdot \mathbf{x} \text{ on } \partial \Omega\}, \tag{41}$$

while the second one is in terms of the minimum complementary energy, according to which $\tau$ is the solution of the problem

$$\inf_{\tau \in S_2(\Sigma)} \langle u(\sigma) \rangle, \tag{42}$$

with

$$S_2(\Sigma) = \{\tau, \nabla \cdot \tau = \mathbf{0}, \text{ in } \Omega, \langle \tau \rangle = \Sigma\}. \tag{43}$$

Then, since the infimum problem in (40) defines the average strain energy in the material, the effective strain-energy potential $\mathcal{H}_e$ is defined as

$$\mathcal{H}_e(\mathbf{E}) = \inf_{\mathbf{v} \in S_1(\mathbf{E})} \langle w[\boldsymbol{\epsilon}(\mathbf{v})] \rangle, \tag{44}$$

so that

$$\frac{\partial \mathcal{H}_e}{\partial \mathbf{E}} = \left\langle \frac{\partial w}{\partial \boldsymbol{\epsilon}}[\boldsymbol{\epsilon}(\mathbf{u})] : \boldsymbol{\epsilon}\left(\frac{\partial \mathbf{u}}{\partial \mathbf{E}}\right) \right\rangle = \left\langle \boldsymbol{\sigma} : \boldsymbol{\epsilon}\left(\frac{\partial \mathbf{u}}{\partial \mathbf{E}}\right) \right\rangle. \tag{45}$$

However, since $\partial \mathbf{u}/\partial \mathbf{E} = \mathbf{U} \cdot \mathbf{x}$, where $\mathbf{U}$ is the identity tensor in the space of fourth-rank tensors, it follows from Eq. (39) that

$$\frac{\partial \mathcal{H}_e}{\partial \mathbf{E}} = \langle \boldsymbol{\sigma} \rangle : \mathbf{U} = \boldsymbol{\Sigma}, \tag{46}$$

which defines the effective stress-strain relation for the material. Similarly, the effective stress-energy potential $\mathcal{H}_e^*$ is defined as

$$\mathcal{H}_e^*(\boldsymbol{\Sigma}) = \inf_{\boldsymbol{\tau} \in S_1(\boldsymbol{\Sigma})} \langle u(\boldsymbol{\tau}) \rangle, \tag{47}$$

in terms of which,

$$\mathbf{E} = \frac{\partial \mathcal{H}_e^*}{\partial \boldsymbol{\Sigma}}. \tag{48}$$

Both $\mathcal{H}_e$ and $\mathcal{H}_e^*$ are convex functions. Furthermore, it can be shown (Suquet, 1987; Willis, 1989a) that they are in fact the (Legendre) dual functions, such that

$$\mathcal{H}_e(\mathbf{E}) + \mathcal{H}_e^*(\boldsymbol{\Sigma}) = \langle w(\boldsymbol{\epsilon}) \rangle + \langle u(\boldsymbol{\sigma}) \rangle = \langle \boldsymbol{\sigma} : \boldsymbol{\epsilon} \rangle = \boldsymbol{\Sigma} : \mathbf{E}, \tag{49}$$

and that they correspond to the boundary conditions (35). Adopting the boundary condition (36) would lead to different pairs of dual potentials. However, under the assumption that the potentials $w$ and $u$ are strictly convex, the two types of boundary conditions are equivalent for the representative volume element, and are also equivalent to the periodic boundary conditions used in the theory of homogenization (see, for example, Sanchez-Palencia, 1980).

As an example, consider the limiting case of rigid, ideally plastic materials for which the potentials are convex, but not strictly. In this limit, which requires special treatment (Bouchitte and Suquet, 1991), $\mathcal{H}_e$ is a positively-homogeneous function of order one in $\mathbf{E}$, usually referred to as the *plastic dissipation function*. It may also be useful to introduce the *effective strength domain* of the material, defined as (Suquet, 1983)

$$P_e = \{\boldsymbol{\Sigma} \text{ such that there exists } \boldsymbol{\sigma}(\mathbf{x}) \text{ with } \langle \boldsymbol{\sigma} \rangle = \boldsymbol{\Sigma} \text{ and}$$

$$\nabla \cdot \boldsymbol{\sigma}(\mathbf{x}) = \mathbf{0}, \text{ with } \boldsymbol{\sigma}(\mathbf{x}) \in P_i, \text{ for } \mathbf{x} \text{ in phase } i\}. \tag{50}$$

Note that

$$\mathcal{H}_e(\mathbf{E}) = \sup_{\boldsymbol{\Sigma} \in P_e} \{\boldsymbol{\Sigma} : \mathbf{E}\}$$

and that

$$\mathcal{H}_e^*(\mathbf{\Sigma}) = \begin{cases} 0 & \text{if } \mathbf{\Sigma} \in P_e \\ +\infty & \text{otherwise} \end{cases}$$

### 4.3.1   One-Point Bounds

As already discussed in Chapter 2, the minimum energy principles can be utilized for deriving rigorous bounds for the effective potentials $\mathcal{H}_e$ and $\mathcal{H}_e^*$ for rigid, ideally plastic polycrystalline materials (Bishop and Hill, 1951a,b) and for materials with elastic, ideally plastic phases (Drucker, 1966). If we use uniform trial fields in the variational principles, the following rigorous bounds of the Voigt (1889) and Reuss (1929) type are obtained:

$$\mathcal{H}_e(\mathbf{E}) \leq \langle u \rangle (\mathbf{E}) = \sum_{i=1}^N \phi_i w_i(\mathbf{E}), \tag{51}$$

and

$$\mathcal{H}_e^*(\mathbf{\Sigma}) \leq \langle u \rangle (\mathbf{\Sigma}) = \sum_{i=1}^N \phi_i w_i(\mathbf{\Sigma}), \tag{52}$$

or, equivalently,

$$\left( \sum_{i=1}^N \phi_i u_i \right)^* (\mathbf{E}) \leq \mathcal{H}_e(\mathbf{E}) \leq \left( \sum_{i=1}^N \phi_i w_i \right) (\mathbf{E}), \tag{53}$$

where superscript $*$ denotes the convex dual function. In the context of polycrystalline materials, the bounds (51) and (52) are commonly referred to as the Taylor (1938) and Sachs (1928) bounds, respectively. For example, the Reuss and Voigt bounds for incompressible, isotropic power-law phases are given by

$$\frac{\langle (\sigma^0)^{-n} \rangle^{-m} \epsilon^0}{m+1} \left( \frac{E_{eq}}{\epsilon^0} \right)^{m+1} \leq \mathcal{H}_e(\mathbf{E}) \leq \frac{\langle \sigma^0 \rangle \epsilon^0}{m+1} \left( \frac{E_{eq}}{\epsilon^0} \right)^{m+1}. \tag{54}$$

Since the Voigt and Reuss bounds incorporate only limited information on the morphology of a material—the volume fractions of the phases—they are not very useful, particularly when the contrast between the phases is large. In fact, they can be shown to be exact only to first order in the contrast between the properties of the phases.

### 4.3.2   Two-Point Bounds: The Talbot–Willis Method

We have already described and discussed in Chapter 2, as well as Chapters 4 and 7 of Volume I, the variational procedure of Hashin and Shtrikman (Hashin and Shtrikman, 1962a,b, 1963). A generalization of the Hashin–Shtrikman variational principles, suitable for nonlinear materials, was developed by Talbot and Willis

(1985), following the earlier work of Willis (1983), which we now describe and discuss.

Let $w^0$ be the potential function of a linear, homogeneous reference material with uniform modulus tensor $\mathcal{L}^0$, such that

$$w^0(\boldsymbol{\epsilon}) = \frac{1}{2} \, \boldsymbol{\epsilon} : \mathcal{L}^0 : \boldsymbol{\epsilon}, \tag{55}$$

and assume that the difference potential $(w - w^0)$ is a concave function, so that the concave polar of this difference is defined as (see Ponte Castañeda and Suquet, 1998)

$$(w - w^0)_*(\mathbf{x}, \boldsymbol{\tau}) = \inf_{\boldsymbol{\epsilon}} \left\{ \boldsymbol{\tau} : \boldsymbol{\epsilon} - \left[ w(\mathbf{x}, \boldsymbol{\epsilon}) - w^0(\boldsymbol{\epsilon}) \right] \right\}.$$

The concavity of $(w - w^0)$ results in

$$w(\mathbf{x}, \boldsymbol{\epsilon}) - w^0(\boldsymbol{\epsilon}) = \inf_{\boldsymbol{\tau}} \left\{ \boldsymbol{\tau} : \boldsymbol{\epsilon} - (w - w^0)_*(\mathbf{x}, \boldsymbol{\tau}) \right\}, \tag{56}$$

Substituting (56) for $w$ in Eq. (47) and interchanging the order of the infima over $\boldsymbol{\epsilon}$ and $\boldsymbol{\tau}$, one arrives at

$$\mathcal{H}_e(\mathbf{E}) = \inf_{\boldsymbol{\tau}} \left\{ \inf_{\mathbf{v} \in S_1(\mathbf{E})} \left\{ \langle w^0[\boldsymbol{\epsilon}(\mathbf{v})] + \boldsymbol{\tau} : \boldsymbol{\epsilon}(\mathbf{v}) \rangle - \langle (w - w^0)_*(\mathbf{x}, \boldsymbol{\tau}) \rangle \right\} \right\}. \tag{57}$$

It then follows that minimizing the displacement field $\mathbf{u}$ is equivalent to finding the solution to the following boundary value problem:

$$\nabla \cdot \left[ \mathcal{L}^0 : \boldsymbol{\epsilon}(\mathbf{u}) \right] = -\nabla \cdot \boldsymbol{\tau}, \quad \mathbf{u} \in S_1(\mathbf{E}). \tag{58}$$

If one utilizes the Green function $\mathbf{G}^0$ associated with the system (58) in the domain $\Omega$, one obtains the following expressions for the strain tensor,

$$\boldsymbol{\epsilon} = \mathbf{E} - \boldsymbol{\Gamma}^0 * \boldsymbol{\tau}, \tag{59}$$

where, as before, $\mathbf{E}$ is the average strain over $\Omega$, and

$$\boldsymbol{\Gamma}^0 * \boldsymbol{\tau} = \int_{\Omega} \boldsymbol{\Gamma}^0(\mathbf{x}, \mathbf{x}') : \left[ \boldsymbol{\tau}(\mathbf{x}') - \langle \boldsymbol{\tau} \rangle \right] d\mathbf{x}', \tag{60}$$

with

$$\Gamma^0_{ijkl} = \left( \frac{\partial^2 G^0_{ik}}{\partial x_j \partial x'_l} \right)_{(ij),(kl)}.$$

Note that the concavity of $(w - w^0)$ is essential in attaining the equality in (57). Typically, however, $(w - w^0)$ is neither concave nor convex, as in the case of, for example, a power-law material. In such a case, the equality in (57) must be replaced by an inequality [either $\leq$ or $\geq$, depending on whether $(w - w^0)$ grows *weaker-than-affine* or *stronger-than-affine* at infinity, respectively].

However, as already pointed out in Chapter 2, as well as Chapters 4 and 7 of Volume I, it is very difficult, if not impossible, to determine the exact $\boldsymbol{\tau}$ that

satisfies (57). Because of this difficulty, an approximation of the following form (the so-called *piecewise constant polarization approximation*)

$$\boldsymbol{\tau}(\mathbf{x}) = \sum_{i=1}^{N} m_i(\mathbf{x}) \boldsymbol{\tau}_i \tag{61}$$

is usually used. Since $\phi_i = \langle m_i(\mathbf{x}) \rangle$ denotes the volume fraction of the phase $i$, and given the fact that the average of a tensor $\mathbf{T}$ over phase $i$ is given by, $\langle \mathbf{T} \rangle_i = \langle (m_i/\phi_i)\mathbf{T} \rangle$, it follows from (59) and (60) that

$$\langle \boldsymbol{\epsilon} \rangle_i = \mathbf{E} - \frac{1}{\phi_i} \sum_{j=1}^{N} \Gamma_{ij} : \boldsymbol{\tau}_j, \tag{62}$$

where

$$\Gamma_{ij} = \left\langle \int_{\Omega} m_i(\mathbf{x}) \left[ m_j(\mathbf{x}') - \phi_j \right] \Gamma^0(\mathbf{x}, \mathbf{x}') d\mathbf{x}' \right\rangle, \quad i, j = 1, \cdots, N \tag{63}$$

are tensors that depend only on the microstructure of the material and $\mathcal{L}^0$, and $\Gamma_{ij}$ are symmetric (Kohn and Milton, 1986) in $i$ and $j$ and are not all independent, since they satisfy the relations

$$\sum_{i=1}^{N} \Gamma_{ij} = \sum_{j=1}^{N} \Gamma_{ij} = 0.$$

After some algebra, one obtains

$$\mathcal{H}_e(\mathbf{E}) \leq \inf_{\boldsymbol{\tau}_l, \, l=1,\ldots,N} \left\{ w^0(\mathbf{E}) + \langle \boldsymbol{\tau} \rangle : \mathbf{E} - \sum_{i=1}^{N} \phi_i (w_i - w^0)_*(\boldsymbol{\tau}_i) \right.$$

$$\left. -\frac{1}{2} \sum_{i=1}^{N} \sum_{j=1}^{N} \boldsymbol{\tau}_i : \Gamma_{ij} : \boldsymbol{\tau}_j \right\}, \tag{64}$$

where $\langle \boldsymbol{\tau} \rangle = \sum_{i=1}^{N} \phi_i \boldsymbol{\tau}_i$. Then, optimizing over $\boldsymbol{\tau}_i$ (with $i = 1, \ldots, N$), one obtains the governing equations for the $\boldsymbol{\tau}_i$:

$$\frac{\partial}{\partial \boldsymbol{\tau}_i} (w_i - w^0)_*(\boldsymbol{\tau}_i) + \frac{1}{\phi_i} \sum_{j=1}^{N} \Gamma_{ij} : \boldsymbol{\tau}_j = \mathbf{E}, \quad i = 1, \ldots, N, \tag{65}$$

so that, from Eqs. (46), (64), and (65) one finally obtains [by replacing the inequality in (64) by an equality] an *approximate* stress-strain relation:

$$\boldsymbol{\Sigma} = \Gamma^0 : \mathbf{E} + \langle \boldsymbol{\tau} \rangle, \tag{66}$$

where the $\boldsymbol{\tau}_i$ are obtained from Eqs. (65).

The upper bound (64) for $\mathcal{H}_e(\mathbf{E})$, which was first given by Ponte Castañeda and Willis (1988), can be written in an alternative form (Willis, 1991) by noting,

through the use of (62), that the optimality conditions (65) can be rewritten in the form

$$\langle \boldsymbol{\epsilon} \rangle_i = \frac{\partial}{\partial \boldsymbol{\tau}_i} \left( w_i - w^0 \right)_* (\boldsymbol{\tau}_i)$$

which, when inverted, yield

$$\boldsymbol{\tau}_i = \frac{\partial}{\partial \langle \boldsymbol{\epsilon} \rangle_i} \left( w_i - w^0 \right)_* (\langle \boldsymbol{\epsilon} \rangle_i), \tag{67}$$

so that the (Legendre) dual variables $\langle \boldsymbol{\epsilon} \rangle_i$ satisfy the conditions

$$\langle \boldsymbol{\epsilon} \rangle_i + \frac{1}{\phi_i} \sum_{j=1}^{N} \Gamma_{ij} : \frac{\partial}{\partial \langle \boldsymbol{\epsilon} \rangle_j} (w_j - w^0)_* (\langle \boldsymbol{\epsilon} \rangle) = \mathbf{E}, \tag{68}$$

for $i = 1, \cdots, N$. Then, the bound (64) may be rewritten as

$$\mathcal{H}_e(\mathbf{E}) \geq w^0(\mathbf{E}) + \sum_{i=1}^{N} \phi_i [2\boldsymbol{\tau}_i : (\mathbf{E} - \langle \boldsymbol{\epsilon} \rangle_i) + (w_i - w^0)_{**}(\langle \boldsymbol{\epsilon} \rangle_i)], \tag{69}$$

where $\boldsymbol{\tau}_i$ are given in terms of $\langle \boldsymbol{\epsilon} \rangle_i$ by Eq. (67).

An upper bound for $\mathcal{H}_e$ is obtained from (69) for any choice of $w^0$; the sharpest bound is obtained by minimizing over $\mathcal{L}^0$. The resulting bound is finite only if $(w - w^0)$ has weaker-than-affine growth at infinity, which would be the case for, for example, power-law materials. The minimization with respect to $\mathcal{L}^0$ is complicated by the fact that computation of $(w_i - w^0)_{**}$ is difficult. Ponte Castañeda and Willis (1988) and Willis (1989a, b) obtained non-optimal bounds by utilizing values of $\mathcal{L}^0$ for which $(w_i - w^0)_{**} = (w_i - w^0)$. Willis (1991,1992) then showed that improved bounds, agreeing with those of the variational procedure of Ponte Castañeda (1991a), are obtained by eliminating this unnecessary restriction.

The estimate for $\mathcal{H}_e$ provided by Eqs. (64) and (65), or (68) and (69), after optimizing over the choice of $\mathcal{L}^0$, is explicit except for the microstructural parameters $\Gamma_{ij}$, which must be determined separately for each class of morphologies. Explicit expressions for these parameters were derived by Willis (1977,1978) and Ponte Castañeda and Willis (1995) for various classes of disordered morphologies with prescribed two-point correlation functions for the distribution of the phases, including particulate and granular materials (see below).

# 4.4   Variational Principles Based on a Linear Comparison Material

Variational methods for deriving improved bounds and estimates for the effective properties of nonlinear materials, utilizing the effective modulus tensor of suitably selected *linear-elastic comparison materials*, were introduced by Ponte Castañeda

(1991a) for materials with isotropic phases and by Suquet (1993a) for composites with power-law phases. Moreover, a hybrid of the Talbot–Willis and Ponte Castañeda procedures, using a linear thermoelastic comparison material, was proposed by Talbot and Willis (1992). These procedures can in fact be shown to be equivalent under appropriate hypotheses on the local potentials. An important advantage of the variational procedures that involve linear comparison materials is that, they can not only produce the nonlinear Hashin–Shtrikman-type bounds of the Talbot–Willis procedure directly from the corresponding linear Hashin–Shtrikman bounds, but also yield higher-order nonlinear bounds, such as Beran-type bounds, as well as other types of estimates. The application of this technique to deriving bounds and estimates for the effective nonlinear conductivity and dielectric constant of materials was described and discussed in Chapter 2. We now describe the analogous results for the effective nonlinear mechanical properties of materials.

### 4.4.1   Materials with Isotropic Phases

The potential $w$ of a nonlinear material with isotropic phases is written as

$$w(\mathbf{x}, \epsilon) = \frac{9}{2} K(\mathbf{x}) \epsilon_m^2 + f(\mathbf{x}, \epsilon_{eq}^2),$$

where

$$K(\mathbf{x}) = \sum_{i=1}^{N} m_i(\mathbf{x}) K_i, \quad f(\mathbf{x}, \epsilon_{eq}^2) = \sum_{i=1}^{N} m_i(\mathbf{x}) f_i(\epsilon_{eq}^2), \tag{70}$$

with the functions $f_i$, characterizing the deviatoric behavior of the material (see Chapter 7 of Volume I), being defined by the relations, $f_i(p) = \varphi_i(\epsilon_{eq})$ for $p = \epsilon_{eq}^2$. The functions $f_i$ are assumed to be concave functions of $p$, such that $f_i(p) = -\infty$ for $p < 0$, $f_i(0) = 0$, and $f_i \to \infty$ as $p \to \infty$. By definition, the concave dual function of $f_i$ is given by

$$f_i^*(q) = \inf_p \{pq - f_i(p)\} = \inf_{p>0} \{pq - f_i(p)\}.$$

It then follows from the concavity hypothesis that

$$f_i(p) = \inf_q \{pq - f_i^*(q)\} = \inf_{q>0} \{pq - f_i^*(q)\}. \tag{71}$$

Note that the above hypotheses on $f_i$ are consistent with weaker-than-quadratic growth for $w_i$ at infinity, in agreement with the physical requirements for plasticity and creep. For example, for power-law materials characterized by Eq. (10), $\varphi_i \sim \epsilon_{eq}^{1+m}$ ($0 \le m \le 1$), so that $f_i \sim p^{(1+m)/2}$ is a concave function in the interval $[0, \infty]$, even if $\varphi_i$ is itself convex.

We now introduce a linear comparison material with potential $w^0$, such that

$$w^0(\mathbf{x}, \epsilon) = \frac{9}{2} K(\mathbf{x}) \epsilon_m^2 + \frac{3}{2} \mu^0(\mathbf{x}) \epsilon_{eq}^2. \tag{72}$$

Then, using (71) with $q = 3\mu^0/2$, one finds that the potential of the nonlinear material $w$ is given by the *exact* equation,

$$w(\mathbf{x}, \boldsymbol{\epsilon}) = \inf_{\mu^0(\mathbf{x})>0} \left\{ w^0(\mathbf{x}, \boldsymbol{\epsilon}) + v(\mathbf{x}, \mu^0) \right\}, \tag{73}$$

where

$$v(\mathbf{x}, \mu^0) = \sum_{i=1}^{N} m_i(\mathbf{x}) v_i[\mu^0(\mathbf{x})], \quad \text{with } v_i(\mu^0) = -f_i^*(\tfrac{3}{2}\mu^0), \tag{74}$$

Note that (see also Chapter 2)

$$v_i(\mu^0) = \sup_{\boldsymbol{\epsilon}} \left\{ w_i(\boldsymbol{\epsilon}) - w_i^0(\boldsymbol{\epsilon}) \right\}, \tag{75}$$

so that

$$v(\mathbf{x}, \mu^0) = \sup_{\boldsymbol{\epsilon}} \left\{ w_i(\mathbf{x}, \boldsymbol{\epsilon}) - w^0(\mathbf{x}, \boldsymbol{\epsilon}) \right\}. \tag{76}$$

If one substitutes Eq. (73) into (44) for the effective potential $\mathcal{H}_e$, one obtains,

$$\mathcal{H}_e(\mathbf{E}) = \inf_{\mathbf{v} \in S_1(\mathbf{E})} \left\{ \inf_{\mu^0(\mathbf{x})} \left\{ \langle w^0[\mathbf{x}, \boldsymbol{\epsilon}(\mathbf{v})] \rangle - \langle v(\mathbf{x}, \mu^0) \rangle \right\} \right\},$$

from which one obtains, by interchanging the order of the infima over $\boldsymbol{\epsilon}$ and $\mu^0$,

$$\mathcal{H}_e(\mathbf{E}) = \inf_{\mu^0(\mathbf{x})} \left\{ \mathcal{H}_e^0(\mathbf{E}) + V(\mu^0) \right\}, \tag{77}$$

where $V(\mu^0) = \langle v[\mathbf{x}, \mu^0(\mathbf{x})] \rangle$, and $\mathcal{H}_e^0$ is the effective potential of the linear comparison material (Ponte Castañeda, 1992a):

$$\mathcal{H}_e^0(\mathbf{E}) = \inf_{\mathbf{v} \in S_1(\mathbf{E})} \langle w^0[\mathbf{x}, \boldsymbol{\epsilon}(\mathbf{v})] \rangle. \tag{78}$$

It must be emphasized that, under the concavity hypothesis on $f_i$, the variational representation (77) and the usual representation (44) are exactly equivalent.

One can also start from the complementary energy representation (47) for $\mathcal{H}_e^*$ to derive a corresponding dual version of the variational representation (77). In this case

$$\mathcal{H}_e^*(\boldsymbol{\Sigma}) = \sup_{\mu^0(\mathbf{x})} \left\{ (\mathcal{H}_e^*)^0(\boldsymbol{\Sigma}) - V(\mu^0) \right\}, \tag{79}$$

where

$$(\mathcal{H}_e^*)^0(\boldsymbol{\Sigma}) = \inf_{\boldsymbol{\sigma} \in S_2(\boldsymbol{\Sigma})} \langle u^0(\mathbf{x}, \boldsymbol{\sigma}) \rangle \tag{80}$$

is the effective stress potential of the linear comparison material.

### 4.4.2  Strongly Nonlinear Materials

We now consider similar variational principles for strongly nonlinear materials characterized by power-law constitutive equations (Suquet, 1993a). Such materials consist of power-law phases with strain-energy functions (10) with the same exponent $n$ and the reference strain $\epsilon^0$, but with different flow stresses $\sigma^0$. For such materials, the variational representation of the effective strain potentials is given by

$$\mathcal{H}_e(\mathbf{E}) = \inf_{\mathbf{v} \in S_1(\mathbf{E})} \left\{ \frac{1}{m+1} \frac{1}{(\epsilon^0)^m} \langle \sigma^0(\mathbf{x}) \epsilon_{eq}^{m+1}[\mathbf{v}(\mathbf{x})] \rangle \right\}.$$

One can then show that (Ponte Castañeda and Suquet, 1998)

$$\mathcal{H}_e(\mathbf{E}) =$$

$$\frac{1}{m+1} \frac{1}{(\epsilon^0)^m} \inf_{\mu^0(\mathbf{x})>0} \left\{ \mathcal{H}_e^0(\mathbf{E})^{(m+1)/2} \left\langle \left(\frac{3}{2}\mu^0\right)^{(m+1)/(m-1)} (\sigma^0)^{2/(1-m)} \right\rangle^{(1-m)/2} \right\},$$

(81)

and that

$$\mathcal{H}_e^*(\mathbf{\Sigma}) = \frac{\epsilon^0}{n+1} \sup_{\mu^0(\mathbf{x})>0} \left\{ (\mathcal{H}_e^*)^0(\mathbf{\Sigma})^{(n+1)/2} \left\langle \frac{(\sigma^0)^{2n/(n-1)}}{6(\mu^0)^{(n+1)/(n-1)}} \right\rangle^{(1-n)/2} \right\}.$$

(82)

### 4.4.3  Materials with Anisotropic Phases

To derive analogous results for nonlinear composite materials with anisotropic phases, we assume that the functions $F_i$, which define the strain potentials $w_i$ via Eq. (20), are concave on the space of positive, symmetric fourth-rank tensors $\mathcal{P}$, i.e., they satisfy

$$F_i[t\mathcal{P}_1 + (1-t)\mathcal{P}_2] \geq t F_i(\mathcal{P}_1) + (1-t)F_i(\mathcal{P}_2), \quad \forall \mathcal{P}_1 \text{ and } \mathcal{P}_2, \quad 0 \leq t \leq 1,$$

(83)

which implies weaker-than-quadratic growth for the potentials $w_i$ on the strain $\epsilon$, when $\mathcal{P}$ is set equal to $\mathcal{E}$, as defined by (21). The concave dual function of $F_i$ is defined by

$$F_i^*(\mathcal{L}) = \inf_{\mathcal{P}} \{\mathcal{L} :: \mathcal{P} - F_i(\mathcal{P})\},$$

based on which one defines $F(\mathbf{x}, \mathcal{P})$ as

$$F(\mathbf{x}, \mathcal{P}) = \sum_{i=1}^{N} m_i(\mathbf{x}) F_i(\mathcal{P}).$$

Because of definition of $\mathcal{E}$ by Eq. (21), one has

$$F\{\mathbf{x}, \mathcal{E}[\mathbf{v}(\mathbf{x})]\} = \inf_{\mathcal{L}^0(\mathbf{x})} \left\{ \mathcal{L}^0(\mathbf{x}) :: \mathcal{E}[\mathbf{v}(\mathbf{x})] - F^*[\mathbf{x}, \mathcal{L}^0(\mathbf{x})] \right\}.$$

Therefore, from definition (44) one obtains

$$\mathcal{H}_e(\mathbf{E}) = \inf_{\mathbf{v} \in S_1(\mathbf{E})} \langle F[\mathcal{E}(\mathbf{v})] \rangle = \inf_{\mathbf{v} \in S_1(\mathbf{E})} \inf_{\mathcal{L}^0(\mathbf{x}) > 0} \left\{ \langle \mathcal{L}^0 :: \mathcal{E}(\mathbf{v}) \rangle - \langle F^*[\mathbf{x}, \mathcal{L}^0(\mathbf{x})] \rangle \right\}.$$

(84)

Then, introducing a linear comparison material with a local potential,

$$w^0[\mathbf{x}, \boldsymbol{\epsilon}(\mathbf{v})] = \mathcal{L}^0 :: \mathcal{E}(\mathbf{v}) = \frac{1}{2} \boldsymbol{\epsilon}(\mathbf{v}) : \mathcal{L}^0(\mathbf{x}) : \boldsymbol{\epsilon}(\mathbf{v}),$$

(85)

and interchanging the infima in (84), one obtains the following exact variational representation for the effective potential,

$$\mathcal{H}_e(\mathbf{E}) = \inf_{\mathcal{L}^0 > 0} \left\{ \mathcal{H}_e^0(\mathbf{E}) + V(\mathcal{L}^0) \right\},$$

(86)

where $\mathcal{H}_e^0$ is the effective potential of the linear comparison material defined by the local potential (85), and $V(\mathcal{L}^0) = \langle v[\mathbf{x}, \mathcal{L}^0(\mathbf{x})] \rangle$, given by

$$v[\mathbf{x}, \mathcal{L}^0(\mathbf{x})] = -F^*[\mathbf{x}, \mathcal{L}^0(\mathbf{x})] = \sup_{\mathcal{P}} [F(\mathbf{x}, \mathcal{P}) - \mathcal{L}^0(\mathbf{x}) :: \mathcal{P}].$$

(87)

Equation (86) expresses the nonlinear effective properties of the material in terms of two functions which are, (1) $\mathcal{H}_e^0$, the elastic energy of a *fictitious linear heterogeneous solid*, called the linear comparison material, that consists of phases with stiffness $\mathcal{L}^0(\mathbf{x})$, and (2) $v(\mathbf{x}, \cdot)$, the role of which is to measure the difference between the non-quadratic potential $w(\mathbf{x}, \cdot)$ and the quadratic energy of the linear comparison solid. The linear comparison solid is selected from amongst all the possible comparison materials by solving the optimization problem (86).

Equation (86), which is exact, is strictly equivalent to the variational representation of $\mathcal{H}_e$ given by (44). However, determining the exact solution of (86) is not possible. The difficulty lies in the precise determination of the energy $\mathcal{H}_e^0$ for a linear comparison solid consisting of infinitely many different phases, about which very little is known. For this reason, except for very simple microstructures, such as laminates considered in Section 2.3, the optimal solution of (86) is not known, and only sub-optimal solutions can be determined. We now consider application of these principles to a few classes of materials.

### 4.4.3.1  Polycrystalline Materials

If the individual phases of a material are single crystals, the functions $G_i$ [see Eqs. (24) and (25)] are given by

$$G_i(\mathbf{S}) = \sum_{k=1}^{M} g_{(k)} \left[ 2\mathcal{M}_i^{(k)} :: \mathbf{S} \right],$$

where, as before

$$\mathcal{M}_i^{(k)} = \boldsymbol{\mu}_i^{(k)} \otimes \boldsymbol{\mu}_i^{(k)}.$$

It can then be shown that

$$G_i^*(\mathcal{M}_i) = \inf_{\alpha_i^{(k)} > 0} \sum_{k=1}^{M} \left(g_i^{(k)}\right)^* \left(\alpha_i^{(k)}\right),$$

with

$$\mathcal{M}_i^{(k)} = \begin{cases} 2\sum_{k=1}^{N} \alpha_i^{(k)} \mathcal{M}_i^{(k)}, & \alpha_i^{(k)} > 0, \\ +\infty, & \text{otherwise.} \end{cases} \tag{88}$$

Therefore, the corresponding local stress potential for the linear comparison polycrystalline material is given by

$$u_i^0(\boldsymbol{\sigma}) = \frac{1}{2}\boldsymbol{\sigma} : \mathcal{M}_i : \boldsymbol{\sigma} = \sum_{k=1}^{M} \alpha_i^{(k)} \left|\tau_i^{(k)}\right|^2, \quad \left|\tau_i^{(k)}\right|^2 = 2\mathcal{M}_i :: \boldsymbol{\sigma}. \tag{89}$$

Hence, for polycrystalline materials one obtains

$$\mathcal{H}_e^*(\boldsymbol{\Sigma}) = \sup_{\alpha_i^{(k)} > 0} \left\{(\mathcal{H}_e^*)^0(\boldsymbol{\Sigma}) - V[\alpha_i^{(k)}]\right\}, \quad i = 1, \cdots, N, \quad k = 1, \cdots, M \tag{90}$$

where $(\mathcal{H}_e^*)^0$ is the effective potential of the linear comparison polycrystalline material with grain potentials (89), and

$$V[\alpha_i^{(k)}] = \sum_{i=1}^{N} \sum_{k=1}^{M} \phi_i \left\langle [g_i^{(k)}]^* \left[\alpha_i^{(k)}\right]\right\rangle_i,$$

which was first derived by deBotton and Ponte Castañeda (1995). The functions $\alpha_i^{(k)}(\mathbf{x})$, $k = 1, \ldots, M$, are defined over the region in space that is occupied by the crystals with fixed orientation $i$.

### 4.4.3.2    Strongly Nonlinear Materials

If the individual phases of a composite are power-law materials with the same exponent $m$ (with $0 \leq m \leq 1$), the composite itself is also a power-law material. That is, the local potentials and the effective macroscopic potential are given by

$$w_i(\lambda\boldsymbol{\epsilon}) = \lambda^{m+1} w_i(\boldsymbol{\epsilon}), \quad \mathcal{H}_e(\lambda\mathbf{E}) = \lambda^{m+1}\mathcal{H}_e(\mathbf{E}), \quad \forall \lambda \geq 0,$$

i.e., they are homogeneous function of order $m + 1$. The function $F$ that defines the strain potential $w$ is itself a power-law function of degree $\frac{1}{2}(m + 1)$, and its dual is a power-law function of degree $(m + 1)/(m - 1)$. If we let $\mathcal{L}^0(\mathbf{x}) = t\hat{\mathcal{L}}^0(\mathbf{x})$ for any $t > 0$, and note that $\mathcal{H}_e^0$ and $V = -\langle F^*\rangle > 0$ are homogeneous functions of orders 1 and $(m + 1)/(m - 1)$ in $\mathcal{L}^0$, respectively, it follows from the variational statement (86) that

$$\mathcal{H}_e(\mathbf{E}) = \inf_{\hat{\mathcal{L}}^0 > 0} \inf_{t > 0} \left\{t\hat{\mathcal{H}}_e^0(\mathbf{E}) + t^{(m+1)/(m-1)} V(\hat{\mathcal{L}}^0)\right\}.$$

where $\hat{\mathcal{H}}_e^0$ is the same as $\mathcal{H}_e^0$ in (86), but with $\mathcal{L}^0$ replaced by $\hat{\mathcal{L}}^0$. Evaluating the minimum over $t$ yields an exact representation for $\mathcal{H}_e$:

$$\mathcal{H}_e(\mathbf{E}) = \frac{2}{m+1} \inf_{\hat{\mathcal{L}}^0 > 0} \left\{ \mathcal{H}_e^0(\mathbf{E})^{(m+1)/2} \left[ \frac{1+m}{1-m} V(\mathcal{L}^0) \right]^{(1-m)/2} \right\}, \qquad (91)$$

where the hat notation has been deleted for simplicity. The analogous representation for $\mathcal{H}_e^*$ is given by

$$\mathcal{H}_e^*(\mathbf{\Sigma}) = \frac{2}{n+1} \sup_{\mathcal{L}^0 > 0} \left\{ (\mathcal{H}_e^*)^0(\mathbf{\Sigma})^{(n+1)/2} \left[ \frac{n+1}{n-1} V(\mathcal{L}^0) \right]^{(1-n)/2} \right\}. \qquad (92)$$

### 4.4.3.3   Materials with Isotropic and Strongly Nonlinear Phases

In this case, it is sufficient to consider isotropic linear comparison materials. If such materials are governed by Eq. (10), then, the functions $f$, $f^*$, $g$ and $g^*$ are given by

$$f(x) = \frac{\sigma^0 \epsilon^0}{m+1} \left[ \frac{|x|}{(\epsilon^0)^2} \right]^{(m+1)/2}, \quad f^*(y) = \frac{m-1}{m+1} \left[ \frac{\sigma^0}{2(\epsilon^0)^m} \right]^{2/(1-m)} |y|^{(m+1)/(m-1)},$$

$$g(x) = \frac{\sigma^0 \epsilon^0}{n+1} \left[ \frac{|x|}{(\sigma^0)^2} \right]^{(n+1)/2}, \quad g^*(y) = \frac{n-1}{n+1} \left[ \frac{2(\sigma^0)^n}{\epsilon^0} \right]^{2/(n-1)} |y|^{(n+1)/(n-1)},$$

$$(93)$$

and

$$V(\mathcal{L}^0) = -\left\langle f^* \left( \frac{3}{2} \mu^0 \right) \right\rangle$$

$$= \frac{1-m}{1+m} \left[ \frac{1}{2(\epsilon^0)^m} \right]^{2/(1-m)} \left\langle \left( \frac{3}{2} \mu^0 \right)^{(m+1)/(m-1)} (\sigma^0)^{2/(1-m)} \right\rangle.$$

### 4.4.3.4   Strongly Nonlinear Polycrystalline Materials

The corresponding result for power-law polycrystalline materials with potentials (16) is obtained directly from Eq. (92), with the result being

$$\mathcal{H}_e^*(\mathbf{\Sigma}) =$$

$$\frac{\gamma^0}{n+1} \sup_{\alpha_i^{(k)} > 0} \left\{ (\mathcal{H}_e^*)^0(\mathbf{\Sigma})^{(n+1)/2} \left[ \sum_{i=1}^{N} \sum_{k=1}^{M} \phi_i \left\langle [\alpha_i^{(k)}]^{(n+1)/(n-1)} \left[ (\tau^0)_i^{(k)} \right]^{2n/(n-1)} \right\rangle_i \right]^{(1-n)/2} \right\},$$

$$(94)$$

for $i = 1, \cdots, N$ and $k = 1, \cdots, M$.

#### 4.4.3.5    Ideally Plastic Materials

In the ideally plastic limit, $m \to 0$, the variational representations (91) and (92) reduce to

$$\mathcal{H}_e(\mathbf{E}) = 2 \inf_{\mathcal{L}^0} \left\{ \mathcal{H}_e^0(\mathbf{E}) V(\mathcal{L}^0) \right\}^{1/2}, \tag{95}$$

and

$$\mathcal{H}_e^*(\mathbf{\Sigma}) = \begin{cases} 0 & \text{if } (\mathcal{H}_e^*)^0(\mathbf{E}) \leq V(\mathcal{L}^0) \ \forall \mathcal{L}^0 = (\mathcal{M}^0)^{-1} > 0, \\ +\infty & \text{otherwise.} \end{cases} \tag{96}$$

## 4.5    Bounds with Piecewise Constant Elastic Moduli

The exact computation of $\mathcal{H}_e$ requires the determination of the effective potential of a linear material with infinitely many different phases, an extremely difficult problem, which may be simplified by restricting the optimization over $\mathcal{L}(\mathbf{x})$ to the set of piecewise constant moduli,

$$\mathcal{L}^0(\mathbf{x}) = \sum_{i=1}^{N} m_i(\mathbf{x}) \mathcal{L}_i^0, \tag{97}$$

where the tensors $\mathcal{L}_i^0$ are assumed constant. In this manner an upper bound for $\mathcal{H}_e$, given by

$$\mathcal{H}_e(\mathbf{E}) \leq \inf_{\mathcal{L}_i^0 > 0, \ i=1,\cdots,n} \left\{ \mathcal{H}_e^0(\mathbf{E}) + \sum_{i=1}^{N} \phi_i v_i(\mathcal{L}_i^0) \right\}, \tag{98}$$

is obtained in which $\mathcal{H}_e^0$ is the effective potential [see Eq. (78)] of a linear composite with the same microstructure as the nonlinear material with the domains $\Omega_i$ occupied by linear phases with stiffness $\mathcal{L}_i^0$. The comparison material has an effective stiffness $\mathcal{L}_e^0$, such that

$$\mathcal{H}_e^0(\mathbf{E}) = \frac{1}{2} \mathbf{E} : \mathcal{L}_e^0 : \mathbf{E}, \tag{99}$$

and the functions $v_i$ are defined by

$$v_i(\mathcal{L}_i^0) = -F_i^*(\mathcal{L}_i^0). \tag{100}$$

The bound (98) is a generalization for materials with anisotropic phases of a corresponding bound for composites with isotropic phases, introduced by Ponte Castañeda (1991a).

A bound equivalent to (98) can be derived by considering the stress potential $\mathcal{H}_e^*$ and its variational representation. Thus, utilizing the piecewise constant compliances $\mathcal{M}_i^0$, one obtains

$$\mathcal{H}_e^*(\mathbf{\Sigma}) \geq \sup_{\mathcal{M}_i^0 > 0, \ i=1,\cdots,N} \left\{ (\mathcal{H}_e^*)^0(\mathbf{\Sigma}) - \sum_{i=1}^{N} \phi_i v_i \left[ (\mathcal{M}_i^0)^{-1} \right] \right\}, \tag{101}$$

with $(\mathcal{H}_e^*)^0$ now being the effective stress potential associated with the same linear comparison material as for $\mathcal{H}_e^0$ given above, i.e., one with the same microstructure as the nonlinear material, but with the domains $\Omega_i$ occupied by linear phases with compliances $\mathcal{M}_i^0$. From Eq. (101) one obtains

$$\mathbf{E} = \frac{\partial \mathcal{H}_e^*}{\partial \boldsymbol{\Sigma}}(\boldsymbol{\Sigma}) = \mathcal{M}_e^0(\mathcal{M}_i^0) : \boldsymbol{\Sigma}. \tag{102}$$

It also follows from (100) that

$$\mathcal{H}_e^*(\boldsymbol{\Sigma}) \geq \inf_{\boldsymbol{\tau} \in S_2(\boldsymbol{\Sigma})} \left\{ \sum_{i=1}^N \phi_i G_i(S_i^*) = \sum_{i=1}^N \phi_i G_i(\langle S \rangle_i) \right\}, \tag{103}$$

where $S_i^* = \langle S \rangle_i = \frac{1}{2}\langle \boldsymbol{\sigma} \otimes \boldsymbol{\sigma} \rangle_i$ is the second moment of the stress field in phase $i$ of the linear comparison material. The compliances $\mathcal{M}_i^0$ of the comparison material are determined as the solution of the optimization problem (100), which can alternatively be written in terms of the solution of the following nonlinear problem for the variables $S_i^*$:

$$\mathcal{M}_i^0 = \frac{\partial G_i}{\partial S}(S_i^*), \quad S_i^* = \frac{1}{2\phi_i} \boldsymbol{\Sigma} : \frac{\partial \mathcal{M}_e^0}{\partial \mathcal{M}_i^0} : \boldsymbol{\Sigma}. \tag{104}$$

## 4.5.1  Materials with Isotropic Phases

If the nonlinear phases are isotropic, then the constituent phases of the linear comparison material can also be selected to be isotropic. The effective bulk modulus is then equal to the bulk modulus $K_i$ of the nonlinear constituent phase $i$, and therefore the only modulus that must be determined is the shear modulus $\mu_i^0$ of each phase. Thus, the bound (98) reduces to

$$\mathcal{H}_e(\boldsymbol{\Sigma}) \geq \inf_{\mu_i^0 > 0, \ i=1,\cdots,N} \left\{ \frac{1}{2}\mathbf{E} : \mathcal{L}_e^0(\mu_i^0) : \mathbf{E} + \sum_{i=1}^N \phi_i v_i(\mu_i^0) \right\}, \tag{105}$$

where the functions $v_i$ are defined by Eq. (75). The upper bound (105), as well as the analogous lower bound,

$$\mathcal{H}_e^*(\boldsymbol{\Sigma}) \geq \sup_{\mu_i^0 > 0, \ i=1,\cdots,N} \left\{ \frac{1}{2}\boldsymbol{\Sigma} : \mathcal{M}_e^0(\mu_i^0) : \boldsymbol{\Sigma} - \sum_{i=1}^N \phi_i v_i(\mu_i^0) \right\}, \tag{106}$$

were first derived by Ponte Castañeda (1991a). Based on the associated optimality conditions (deBotton and Ponte Castañeda, 1992,1993), it can be shown that

$$\mathcal{H}_e \leq \frac{9}{2} \sum_{i=1}^N \phi_i K_i \left[ \epsilon_i^{(m)} \right]^2 + \sum_{i=1}^N \phi_i \varphi_i \left( \epsilon_i^{eq} \right)^2, \tag{107}$$

$$\mathcal{H}_e^* \geq \frac{1}{2} \sum_{i=1}^N \frac{\phi_i}{K_i} \left[ \sigma_i^{(m)} \right]^2 + \sum_{i=1}^N \phi_i \psi_i \left( \sigma_i^{eq} \right)^2, \tag{108}$$

where

$$\epsilon_i^{(m)} = \left( \frac{1}{9\phi_i} \mathbf{E} : \frac{\partial \mathcal{L}_e^0}{\partial K_i^0} : \mathbf{E} \right)^{1/2}, \quad \sigma_i^{(m)} = \left( \frac{1}{\phi_i} \mathbf{\Sigma} : \frac{\partial \mathcal{M}_e^0}{\partial (1/K_i^0)} : \mathbf{\Sigma} \right)^{1/2}, \quad (109)$$

and

$$\epsilon_i^{eq} = \left( \frac{1}{3\phi_i} \mathbf{E} : \frac{\partial \mathcal{L}_e^0}{\partial \mu_i^0} : \mathbf{E} \right)^{1/2}, \quad \sigma_i^{eq} = \left( \frac{3}{\phi_i} \mathbf{\Sigma} : \frac{\partial \mathcal{M}_e^0}{\partial (1/\mu_i^0)} : \mathbf{\Sigma} \right)^{1/2}. \quad (110)$$

These simplified bounds were first given by Suquet (1995,1997).

For power-law materials, one obtains,

$$\mathcal{H}_e(\mathbf{E}) \leq \frac{1}{m+1} \frac{1}{(\epsilon^0)^m} \inf_{\mu_i^0 > 0} \left\{ \mathcal{H}_e^0(\mathbf{E})^{(m+1)/2} \left[ \sum_{i=1}^{N} \phi_i \left( \frac{3}{2} \mu_i^0 \right)^{(m+1)/(m-1)} (\sigma_i^0)^{2/(1-m)} \right]^{(1-m)/2} \right\}, \quad (111)$$

and

$$\mathcal{H}_e^*(\mathbf{\Sigma}) \geq \frac{\epsilon^0}{n+1} \sup_{\mu_i^0 > 0} \left\{ (\mathcal{H}_e^*)^0(\mathbf{\Sigma})^{(n+1)/2} \left[ \sum_{i=1}^{N} \phi_i (6\mu_i^0)^{(n+1)/(1-n)} (\sigma_i^0)^{2n/(n-1)} \right]^{(1-n)/2} \right\}, \quad (112)$$

which were also derived by Suquet (1993a).

## 4.5.2  Polycrystalline Materials

In this case, one restricts the optimization in (101) to compliance tensors that yield finite values for the functions $v_i$, which then leads to (deBotton and Ponte Castañeda, 1995)

$$\mathcal{H}_e^*(\mathbf{\Sigma}) \geq \sup_{\alpha_i^{(k)} > 0} \left\{ \frac{1}{2} \mathbf{\Sigma} : \mathcal{M}_e^0(\boldsymbol{\alpha}) : \mathbf{\Sigma} - \sum_{i=1}^{N} \sum_{k=1}^{M} \phi_i \left[ g_i^{(k)} \right]^* \left[ \alpha_i^{(k)} \right] \right\}, \quad (113)$$

where the suprema should be performed for $i = 1, \cdots, N$ and $k = 1, \cdots, M$. Here, $\boldsymbol{\alpha}$ denotes the entire set of positive slip compliances $\alpha_i^{(k)}$, and $\mathcal{M}_e$ is the effective compliance tensor of the linear comparison polycrystalline material with grain compliances $\mathcal{M}_i$, as given by Eq. (88) in terms of the slip compliances $\boldsymbol{\alpha}$. One may also approximate the effective stress-strain relation of the polycrystalline material by the relation (102) of the linear comparison material, with the optimal $\mathcal{M}_i^0$ replaced by the optimal $\alpha_i^{(k)}$. In that case, the nonlinear optimality relations are given by

$$\mathcal{M}_i = 2 \sum_{k=1}^{M} \alpha_i^{(k)} \mathcal{M}_i^{(k)}, \quad \alpha_i^{(k)} = \frac{\partial g_i^{(k)}}{\partial \tau} \left( 2\mathcal{M}_i^{(k)} :: \bar{\sigma}_i \right),$$

$$\bar{\sigma}_i = \frac{1}{2\phi_i} \mathbf{\Sigma} : \frac{\partial \mathcal{M}_e}{\partial \mathcal{M}_i} : \mathbf{\Sigma}. \quad (114)$$

These nonlinear equations can be expressed more explicitly in terms of the slip compliances $\alpha_i^{(k)}$ and the corresponding second moment of the resolved shears,

$$\bar{\tau}_i^{(k)} = \left[ 2\mathcal{M}_i^{(k)} :: \bar{\sigma}_i \right]^{1/2} .$$

with the result being

$$\alpha_i^{(k)} = \frac{\partial g_i^{(k)}}{\partial \tau} \left[ \bar{\tau}_i^{(k)} \right]^2 = \frac{1}{2\bar{\tau}_i^{(k)}} \frac{\partial \psi_i^{(k)}}{\partial \tau} \left[ \bar{\tau}_i^{(k)} \right]^* ,$$

$$\bar{\tau}_i^{(k)} = \left( \frac{1}{2\phi_i} \boldsymbol{\Sigma} : \frac{\partial \mathcal{M}_e}{\partial \alpha_i^{(k)}} (\boldsymbol{\alpha}) : \boldsymbol{\Sigma} \right)^{1/2} , \tag{115}$$

in terms of which the bound is rewritten as

$$\mathcal{H}_e^*(\boldsymbol{\Sigma}) \geq \sum_{i=1}^{N} \sum_{k=1}^{M} \phi_i \psi_i^{(k)} \left[ \bar{\tau}_i^{(k)} \right] , \tag{116}$$

For power-law polycrystalline materials, one obtains the following result,

$$\mathcal{H}_e^* \geq \frac{\gamma^0}{n+1} \sup_{\alpha_i^{(k)} > 0} \left\{ (\mathcal{H}_e^*)^0(\boldsymbol{\Sigma})^{(n+1)/2} \left( \sum_{i=1}^{N} \sum_{k=1}^{M} \phi_i \left[ \alpha_i^{(k)} \right]^{(n+1)/(n-1)} \left[ (\tau^0)_i^{(k)} \right]^{2n/(n-1)} \right)^{(1-n)/2} \right\} , \tag{117}$$

where the suprema must be carried out over $i = 1, \cdots, N$ and $k = 1, \cdots, M$.

It should be emphasized that any estimate for the effective modulus tensor of a linear elastic material can be used to generate, by the variational procedures described above, a corresponding estimate for a nonlinear material with the same microstructure. This is in contrast to several other schemes which are closely connected with specific types of estimates. For example, the Talbot–Willis method described above provides only estimates of the Hashin–Shtrikman-type. Moreover, similar to the case of the effective conductivity and dielectric constant of nonlinear materials discussed in Chapter 2, if the estimate for the effective modulus tensor of the linear elastic material is an upper bound to $\mathcal{L}_e$, then an upper bound is obtained for $\mathcal{H}_e$. If, on the other hand, the linear estimate is a lower bound, then, the variational method *cannot*, in general, be used for deriving a lower bound for the nonlinear material. However, if accurate estimates (but not necessarily bounds) are available for a specific type of linear material, such as those provided by the effective-medium approximation, then the above variational methods can be utilized for generating the corresponding estimates for a nonlinear material with the same microstructure. The resulting estimates for $\mathcal{H}_e$ would tend to err on the high side, because of the nature of the approximations intrinsic to the variational method. In addition, these variational methods can be used for deriving higher-order ($\geq 2$) bounds, such as Beran-type bounds (Ponte Castañeda, 1992a), as well as other types of estimates, such as the generalized self-consistent estimates of Suquet (1993b). Let us also mention that Smyshlyaev and Fleck (1995; see also

Fleck and Hutchinson, 1997) proposed extensions of the variational methods in the context of strain gradient plasticity.

We remind the reader that the above variational methods use the concavity hypothesis on the function $F$ associated with the local strain potential $w$. Except for some pathological cases, this mild hypothesis is satisfied by the standard models of plasticity and creep (Willis, 1992; Ponte Castañeda and Willis, 1993). When this hypothesis is satisfied, Ponte Castañeda (1992c) showed, in the context of materials with isotropic phases, that the Talbot–Willis variational method (see Section 4.3) can be directly derived, via variational principles, from the Hashin–Shtrikman variational principles for linear materials. An alternative, simpler way of analyzing materials for which the concavity hypothesis is violated was proposed by Ponte Castañeda (1996b, 1997); see also Kohn and Little (1997) and Bhattacharya and Kohn (1997) in the context of polycrystalline materials.

## 4.6   Second-Order Exact Results

We now describe and discuss exact results for the effective mechanical properties of weakly heterogeneous nonlinear materials, and also estimates for arbitrary contrast of the phases. This analysis represents an extension of a similar theory for linear elasticity, for which it is well-known that the effective moduli tensor of a weakly heterogeneous material can be determined exactly to second order in the contrast (see Chapter 7 of Volume I). The present theory also represents an extension of the analogous theoretical developments for the effective nonlinear conductivity and dielectric constant of heterogeneous materials that were described and discussed in Section 2.6. The analysis that follows also establishes that the above variational estimates are exact only to first order in the phase contrast, when estimates that are exact to second order are used to evaluate the mechanical properties of the linear comparison material.

### 4.6.1   Weak-Contrast Expansion

It is assumed that the contrast between the properties of the phases is small. To incorporate this assumption into the analysis, the potential $w$ is assumed to depend on a small parameter $t$ that characterizes the contrast between the properties of the material and those of a homogeneous nonlinear reference material with energy function $w^0(\epsilon)$, such that

$$w(\mathbf{x}, \epsilon, t) = w^0(\epsilon) + t\delta w(\mathbf{x}, \epsilon). \tag{118}$$

The effective potential also depends on the parameter $t$:

$$\mathcal{H}_e(\mathbf{E}, t) = \langle w[\mathbf{x}, \epsilon(\mathbf{u}_t), t]\rangle, \tag{119}$$

where $\mathbf{u}_t$ and $\epsilon(\mathbf{u}_t)$ are the local displacement and the associated strain fields induced by appropriate boundary conditions that generate an average strain $\mathbf{E}$ in $\Omega$. Furthermore, it is assumed that $\mathcal{H}_e(., t)$ and $\mathbf{u}_t$ are continuously differentiable

functions of $t$. Since $t$ is small, one can write down a perturbation series expansion of $\mathcal{H}_e$ about $t = 0$, given formally by

$$\mathcal{H}_e(\mathbf{E}, t) = \mathcal{H}_e(\mathbf{E}, 0) + t \frac{\partial \mathcal{H}_e}{\partial t}(\mathbf{E}, 0) + \frac{1}{2} t^2 \frac{\partial^2 \mathcal{H}_e}{\partial t^2}(\mathbf{E}, 0) + O(t^3). \quad (120)$$

The problem to be solved for $\mathbf{u}_t$ is given by

$$\nabla \cdot \left\{ \frac{\partial w}{\partial \boldsymbol{\epsilon}}[\mathbf{x}, \boldsymbol{\epsilon}(\mathbf{u}_t), t] \right\} = \mathbf{0}, \quad \mathbf{u}_t \in S_1(\mathbf{E}). \quad (121)$$

If we differentiate (121), we find that $\dot{\mathbf{u}}_t = \partial \mathbf{u}_t / \partial t$ is the solution of the following system of equations

$$\nabla \cdot [\mathcal{L}_t : \boldsymbol{\epsilon}(\dot{\mathbf{u}}_t)] + \nabla \cdot \boldsymbol{\tau}_t = \mathbf{0}, \quad \dot{\mathbf{u}}_t \in S_1(0). \quad (122)$$

where

$$\mathcal{L}_t = \frac{\partial^2 w}{\partial \boldsymbol{\epsilon} \partial \boldsymbol{\epsilon}}[\mathbf{x}, \boldsymbol{\epsilon}(\mathbf{u}_t), t], \quad \boldsymbol{\tau}_t = \frac{\partial^2 w}{\partial t \partial \boldsymbol{\epsilon}}[\mathbf{x}, \boldsymbol{\epsilon}(\mathbf{u}_t), t] = \frac{\partial}{\partial \boldsymbol{\epsilon}}(\delta w)[\mathbf{x}, \boldsymbol{\epsilon}(\mathbf{u}_t), t].$$

Therefore,

$$\frac{\partial \mathcal{H}_e}{\partial t}(\mathbf{E}, t) = \left\langle \frac{\partial w}{\partial \boldsymbol{\epsilon}}[\mathbf{x}, \boldsymbol{\epsilon}(\mathbf{u}_t), t] : \boldsymbol{\epsilon}(\dot{\mathbf{u}}_t) \right\rangle + \left\langle \frac{\partial w}{\partial t}[\mathbf{x}, \boldsymbol{\epsilon}(\mathbf{u}_t), t] \right\rangle. \quad (123)$$

The first term of (123) vanishes due to Eq. (39) (the so-called Hill's lemma), and therefore,

$$\frac{\partial \mathcal{H}_e}{\partial t}(\mathbf{E}, t) = \langle \delta w[\mathbf{x}, \boldsymbol{\epsilon}(\mathbf{u}_t)] \rangle. \quad (124)$$

Using Eq. (121), one obtains

$$\frac{\partial^2 \mathcal{H}_e}{\partial t^2}(\mathbf{E}, t) = \left\langle \frac{\partial}{\partial \boldsymbol{\epsilon}}(\delta w)[\mathbf{x}, \boldsymbol{\epsilon}(\mathbf{u}_t)] : \boldsymbol{\epsilon}(\mathbf{u}_t) \right\rangle = -\langle \boldsymbol{\epsilon}(\dot{\mathbf{u}}_t) : \mathcal{L}_t : \boldsymbol{\epsilon}(\dot{\mathbf{u}}_t) \rangle. \quad (125)$$

Because the material is homogeneous for $t = 0$, $\mathbf{u}_0 = \mathbf{E} \cdot \mathbf{x}$, and therefore

$$\mathcal{H}_e(\mathbf{E}, 0) = w^0(\mathbf{E}),$$
$$\frac{\partial \mathcal{H}_e}{\partial t}(\mathbf{E}, 0) = \langle \delta w \rangle(\mathbf{E}), \quad (126)$$
$$\frac{\partial^2 \mathcal{H}_e}{\partial t^2}(\mathbf{E}, 0) = -\langle \boldsymbol{\epsilon}(\dot{\mathbf{u}}_0) : \mathcal{L}^0 : \boldsymbol{\epsilon}(\dot{\mathbf{u}}_0) \rangle,$$

where

$$\mathcal{L}^0 = \frac{\partial^2 w^0}{\partial \boldsymbol{\epsilon} \partial \boldsymbol{\epsilon}}(\mathbf{E}). \quad (127)$$

Here, $\dot{\mathbf{u}}_0$ is the solution of the linear elasticity problem,

$$\nabla \cdot [\mathcal{L}^0 : \boldsymbol{\epsilon}(\dot{\mathbf{u}}_0)] + \nabla \cdot \boldsymbol{\tau} = \mathbf{0}, \quad \dot{\mathbf{u}}_0 \in S_1(0), \quad (128)$$

with

$$\boldsymbol{\tau}(\mathbf{x}) = \frac{\partial}{\partial \boldsymbol{\epsilon}}(\delta w)(\mathbf{x}, \mathbf{E}).$$

Since the modulus tensor $\mathcal{L}^0$ is constant, the problem posed by (128) is a linear elasticity problem for a homogeneous material with a distribution of body forces determined by $\boldsymbol{\tau}$.

If the material consists of $N$ homogeneous phases, then $\boldsymbol{\tau}$ is piecewise constant, i.e., it is a constant in each phase, with

$$\boldsymbol{\tau}(\mathbf{x}) = \sum_{i=1}^{N} m_i(\mathbf{x}) \, \boldsymbol{\tau}_i, \quad \boldsymbol{\tau}_i = \frac{\partial}{\partial \boldsymbol{\epsilon}}(\delta w_i)(\mathbf{E}),$$

in terms of which one has

$$\frac{\partial^2 \mathcal{H}_e}{\partial t^2}(\mathbf{E}, 0) = -\sum_{i=1}^{N} \sum_{j=1}^{N} \boldsymbol{\tau}_i : \Gamma_{ij} : \boldsymbol{\tau}_j, \tag{129}$$

where the microstructural tensors $\Gamma_{ij}$ are defined by Eq. (63). Therefore (Suquet and Ponte Castañeda, 1993),

$$\mathcal{H}_e(\mathbf{E}, t) = \langle w \rangle(\mathbf{E}) - \frac{1}{2}t^2 \sum_{i=1}^{N} \sum_{j=1}^{N} \boldsymbol{\tau}_i : \Gamma_{ij} : \boldsymbol{\tau}_j + O(t^3). \tag{130}$$

As an example, consider a material with $N$ isotropic phases, with strain potentials defined by Eq. (118) and

$$w^0(\boldsymbol{\epsilon}) = \frac{9}{2}K^0\epsilon_m^2 + f^0(\epsilon_{eq}^2),$$

$$\delta w_i(\boldsymbol{\epsilon}) = \frac{9}{2}\delta K_i\epsilon_m^2 + \delta f_i(\epsilon_{eq}^2).$$

Then, it is straightforward to show that (Ponte Castañeda and Suquet, 1995),

$$\mathcal{H}_e(\mathbf{E}, t) = \frac{9}{2}K^0 E_m^2 + f^0(E_{eq}^2) + t \sum_{i=1}^{N} \phi_i \left[\frac{9}{2}\delta K_i E_m^2 + \delta f_i(E_{eq}^2)\right]$$

$$-\frac{1}{2}t^2 \sum_{i=1}^{N} \sum_{j=1}^{N} \left(9\delta K_i \delta K_j \mathbf{U} : \Gamma_{ij} : \mathbf{U} E_m^2 + 4\delta\mu_i \delta\mu_j \mathbf{E}_d : \Gamma_{ij} : \mathbf{E}_d\right) + O(t^3).$$

$$\tag{131}$$

where $\mathbf{E}_d = \mathbf{E} - E_m \mathbf{U}$ is the average strain deviator.

## 4.6.2 Strong-Contrast Expansion

Another method for estimating the effective mechanical properties of nonlinear materials was proposed by Ponte Castañeda (1996a). His method uses a linear heterogeneous comparison material and the associated tangent modulus tensors of the constituent phases. This choice of comparison material ensures that the resulting nonlinear estimates are exact to second order in the contrast, and thus are in agreement with the small-contrast asymptotic results of the last section.

Similar to the case of the effective conductivity and dielectric constant of nonlinear heterogeneous materials that was discussed in Chapter 2, this method is based on a Taylor expansion for the phase potentials $w_i$. Thus, introducing reference strains $\mathbf{E}^{(i)}$, the Taylor expansion for $w_i$ about $\mathbf{E}^{(i)}$ is given by

$$w_i(\boldsymbol{\epsilon}) = w_i[\mathbf{E}^{(i)}] + \boldsymbol{\rho}_i : [\boldsymbol{\epsilon} - \mathbf{E}^{(i)}] + \frac{1}{2}[\boldsymbol{\epsilon} - \mathbf{E}^{(i)}] : \mathcal{L}_i : [\boldsymbol{\epsilon} - \mathbf{E}^{(i)}], \quad (132)$$

where $\boldsymbol{\rho}_i$ and $\mathcal{L}_i$ are, respectively, an internal stress and a tangent modulus tensor, with components

$$(\rho_k)_{ij} = \frac{\partial w_k}{\epsilon_{ij}}[\mathbf{E}^{(k)}], \quad (\mathcal{L}_m)_{ijkl} = \frac{\partial^2 w_m}{\partial \epsilon_{ij} \partial \epsilon_{kl}}[\tilde{\mathbf{E}}^{(m)}]. \quad (133)$$

$\mathcal{L}_i$ depends on the strain $\tilde{\mathbf{E}}^{(i)} = \lambda^{(i)}\mathbf{E}^{(i)} + [1 - \lambda^{(i)}]\boldsymbol{\epsilon}$, where $\lambda^{(i)}$ depends on $\boldsymbol{\epsilon}$ and is such that $0 < \lambda^{(i)} < 1$.

In terms of the average $\mathbf{E}$ and fluctuating $\boldsymbol{\epsilon}'$ components of $\boldsymbol{\epsilon} = \mathbf{E} + \boldsymbol{\epsilon}'$, Eq. (132) is rewritten as

$$w_i(\mathbf{E} + \boldsymbol{\epsilon}') = v_i + \boldsymbol{\tau}_i : \boldsymbol{\epsilon}' + \frac{1}{2}\boldsymbol{\epsilon}' : \mathcal{L}_i : \boldsymbol{\epsilon}', \quad (134)$$

where

$$v_i = w_i[\mathbf{E}^{(i)}] + \boldsymbol{\rho}_i : [\mathbf{E} - \mathbf{E}^{(i)}] + \frac{1}{2}[\mathbf{E} - \mathbf{E}^{(i)}] : \mathcal{L}_i : [\mathbf{E} - \mathbf{E}^{(i)}], \quad (135)$$

$$\boldsymbol{\tau}_i = \boldsymbol{\rho}_i + \mathcal{L}_i : [\mathbf{E} - \mathbf{E}^{(i)}]. \quad (136)$$

Making the approximation that the strains $\tilde{\mathbf{E}}^{(i)}$ are constant in each phase, the effective potential $\mathcal{H}_e$ of the nonlinear material is then estimated as

$$\mathcal{H}_e(\mathbf{E}) \simeq \tilde{\mathcal{H}}_e(\mathbf{E}) = \sum_{i=1}^{N} \phi_i v_i + P, \quad (137)$$

where

$$P = \inf_{\mathbf{v}' \in S_1(0)} \left\langle \frac{1}{2}\boldsymbol{\epsilon}(\mathbf{v}') : \mathcal{L} : \boldsymbol{\epsilon}(\mathbf{v}') + \boldsymbol{\tau} : \boldsymbol{\epsilon}(\mathbf{v}') \right\rangle. \quad (138)$$

$\boldsymbol{\tau}$ and $\mathcal{L}(\mathbf{x})$ are defined by equations similar to (61). The advantage of approximation (137), relative to the exact result (44), is that it requires only the solution of a linear problem for an $N$-phase *thermoelastic material*, as defined by the Euler–Lagrange equations of the variational problem $P$ in (138):

$$\nabla \cdot [\mathcal{L} : \boldsymbol{\epsilon}(\mathbf{u}')] = -\nabla \cdot \boldsymbol{\tau}, \quad \mathbf{u}' \in S_1(0). \quad (139)$$

Estimates for $N$-phase linear-thermoelastic materials can, in general, be obtained by appropriate extension of the corresponding methods for $N$-phase linear-elastic composites (see, for example, Willis, 1981). Similar, but not equivalent, representations for the effective mechanical properties of nonlinear materials, which also utilize heterogeneous thermoelastic reference materials, were proposed by Molinari *et al.* (1987) and Talbot and Willis (1992).

Equation (137) provides an estimate for $\mathcal{H}_e$ for *any* choice of $\mathbf{E}^{(i)}$ and $\tilde{\mathbf{E}}^{(i)}$, if we supply it with an estimate for $P$. A plausible approximation for the $\mathbf{E}^{(i)}$ is to set them equal to the averages of the strain field over the phases $i$. However, because the exact strain field is not known, the approximate field $\epsilon$, as determined by (139), is used, so that

$$\mathbf{E}^{(i)} = \langle\epsilon\rangle_i, \tag{140}$$

where, $\langle\epsilon\rangle_i = \mathbf{E} + \langle\epsilon'\rangle_i$. Equation (140) is a reasonable choice because the strain $\epsilon$ in phase $i$ is expected to fluctuate about its average in phase $i$ in such a way that large deviations would only be expected in regions of relatively small measure. The following identity, obtained from Eq. (139),

$$\langle\epsilon'\rangle_i = \frac{1}{\phi_i}\frac{\partial P}{\partial\tau_i}, \tag{141}$$

which can be used to obtain $\langle\epsilon\rangle_i$ directly from $P$, via

$$\langle\epsilon\rangle_i = \mathbf{E} + \frac{1}{\phi_i}\frac{\partial P}{\partial\tau_i}, \tag{142}$$

is also useful, since the reference strains $\mathbf{E}^{(i)}$ may also be computed from $P$ by means of Eqs. (140) and (142). It can also be shown that Eq. (140) provides the *optimal* choice for $\mathbf{E}^{(i)}$ in the sense that, estimate (137) for $\mathcal{H}_e$ is stationary with respect to the $\mathbf{E}^{(i)}$.

One important consequence of stationarity of Eq. (140) is that the overall stress-strain relation (46) for the material may be approximated as

$$\mathbf{\Sigma} = \sum_{i=1}^{N}\phi_i\left\{\boldsymbol{\rho}_i + \frac{1}{2}\langle(\epsilon - \langle\epsilon\rangle_i) : \mathcal{N}_i : (\epsilon - \langle\epsilon\rangle_i)\rangle_i : \frac{\partial\tilde{\mathbf{E}}^{(i)}}{\partial\mathbf{E}}\right\}, \tag{143}$$

where $\langle\epsilon\rangle_i$ are determined by Eq. (142), and

$$(\mathcal{N}_m)_{ijklpq} = \frac{\partial^3 w_m}{\partial\epsilon_{ij}\partial\epsilon_{kl}\partial\epsilon_{pq}}[\tilde{\mathbf{E}}^{(m)}],$$

which can be derived by taking the derivative of Eq. (137) with respect to $\mathbf{E}$, with $\mathbf{E}^{(i)}$ held fixed (because of stationarity), and enforcing (140).

Equation (140) also allows simplification of estimate (137) for $\mathcal{H}_e$. Note that the Euler–Lagrange equations, Eqs. (139), of problem (138) for $P$ imply that

$$\langle\tau : \epsilon(\mathbf{u}')\rangle = -\langle\epsilon(\mathbf{u}') : \mathcal{L} : \epsilon(\mathbf{u}')\rangle, \tag{144}$$

which, together with Eq. (140) and the definition (136) of $\tau_i$, are used to rewrite the estimate (137) in the following simpler form,

$$\mathcal{H}_e \simeq \tilde{\mathcal{H}}_e(\mathbf{E}) = \sum_{i=1}^{N}\phi_i\left\{w_i\langle\epsilon\rangle_i + \frac{1}{2}\frac{\partial w_i}{\partial\epsilon}(\langle\epsilon\rangle_i) : (\mathbf{E} - \langle\epsilon\rangle_i)\right\}, \tag{145}$$

with $\langle\epsilon\rangle_i$ being determined by Eqs. (142).

The choice of $\tilde{\mathbf{E}}^{(i)}$ in definition (133) of $\mathcal{L}_i$ is not as straightforward, and, in particular, stationarity of $\tilde{\mathcal{H}}_e$ with respect to $\tilde{\mathbf{E}}^{(i)}$ cannot be implemented. For this reason, Ponte Castañeda (1996a) proposed the following physically motivated equation for $\tilde{\mathbf{E}}^{(i)}$:

$$\tilde{\mathbf{E}}^{(i)} = \langle \boldsymbol{\epsilon} \rangle_i = \mathbf{E}^{(i)}, \tag{146}$$

an interesting consequence of which is that it implies that

$$\frac{\partial^2 \tilde{\mathcal{H}}_e}{\partial \mathbf{E}^{(i)} \partial \mathbf{E}^{(i)}} = 0. \tag{147}$$

It is now not difficult to show that

$$\tilde{\mathcal{H}}_e(\mathbf{E}) = \sum_{i=1}^{N} \phi_i w_i(\mathbf{E}) - \frac{1}{2} t^2 \langle \boldsymbol{\epsilon}(\dot{\mathbf{u}}_0) : \mathcal{L}^0 : \boldsymbol{\epsilon}(\dot{\mathbf{u}}_0) \rangle + O(t^3), \tag{148}$$

in agreement with the small-contrast expansion (120) together with (126).

As an example, consider two-phase materials, for which a well-known result due to Levin (1967) allows further simplification of the thermoelastic problem $P$, and hence of the corresponding estimate for $\mathcal{H}_e$. The result for $P$, which depends only on the effective modulus tensor $\mathcal{L}_e$ of a two-phase, linear elastic material with phase modulus tensors $\mathcal{L}_1$ and $\mathcal{L}_2$, is given by

$$P = \frac{1}{2} (\Delta \boldsymbol{\tau}) : (\Delta \mathcal{L})^{-1} : (\mathcal{L}_e - \langle \mathcal{L} \rangle) : (\Delta \mathcal{L})^{-1} : (\Delta \boldsymbol{\tau}), \tag{149}$$

where $\Delta \mathcal{L} = \mathcal{L}_1 - \mathcal{L}_2$ and $\Delta \boldsymbol{\tau} = \boldsymbol{\tau}_1 - \boldsymbol{\tau}_2$. It then follows from Eqs. (142) and (146) that

$$\mathbf{E}^{(i)} = \tilde{\mathbf{E}}^{(i)} = \langle \boldsymbol{\epsilon} \rangle_i = \mathbf{E} + (\mathcal{A}_i - \mathcal{U}) : (\Delta \mathcal{L})^{-1} : (\Delta \boldsymbol{\tau}), \tag{150}$$

where $\mathcal{A}_i$ denote the strain-concentration tensors (Hill, 1965a) for the linear elastic material problem, such that

$$\phi_1 \mathcal{A}_1 + \phi_2 \mathcal{A}_2 = \mathcal{U}, \quad \mathcal{L}_e = \phi_1 \mathcal{L}_1 : \mathcal{A}_1 + \phi_2 \mathcal{L}_2 : \mathcal{A}_2, \tag{151}$$

which can be solved for the tensors $\mathcal{A}_i$ in terms of the $\mathcal{L}_i$ and $\mathcal{L}_e$.

It must be emphasized that any estimate of any type for $\mathcal{L}_e$ can be used for generating the corresponding estimates for $\mathcal{H}_e$, that the second-order term in the above expansion depends only on the two-point statistics of the material and completely specifies its effective properties (to second order in the contrast), and that by comparison with this exact result, it becomes clear that the variational estimates described above are exact only to first order in the contrast. In addition, the second-order theory does produce estimates that are exact to second order in the contrast. However, the approximations involved in the second-order theory are such that it is *not* possible to control the sign of the error, so that the resulting estimates, unlike the earlier variational estimates, cannot be guaranteed to be bounds to the effective properties. Another important limitation of the second-order theory is the existence of a duality gap, i.e., it can be shown that, $\tilde{\mathcal{H}}_e^* \neq (\tilde{\mathcal{H}}_e)^*$. As a practical matter, in plasticity and creep, as in conductivity and dielectric constant (see Chapter 2),

the second-order estimates based on the estimate $\mathcal{H}_e$ are more accurate than the analogous estimates for $\mathcal{H}_e^*$.

## 4.7    Applications of Second-Order Exact Results

We now consider some applications of the above theoretical results to modeling of mechanical properties of porous materials, and composites with a superrigid phase, as well as more general two-phase power-law and perfectly plastic composites. In each case, we first describe the various bonds that can be obtained from the above general formulation, and then discuss the application of the second-order theory.

### 4.7.1    Porous Materials

We consider porous materials with isotropic matrix phases, so that the strain and stress potentials of the matrix phase are given by Eqs. (4) and (5). Designating the matrix as phase 1, the bound (107) becomes

$$\mathcal{H}_e(\mathbf{E}) \leq \frac{9}{2}\phi_1 K_1 \left[\epsilon_1^{(m)}\right]^2 + \phi_1\varphi_1 \left(\epsilon_1^{eq}\right)^2 \tag{152}$$

where

$$\epsilon_1^{eq} = \sqrt{\frac{1}{3\phi_1}\mathbf{E} : \left[\frac{\partial \mathcal{L}_e^0}{\partial \mu_1^0}(\mu_1^0, K_1)\right] : \mathbf{E}}, \tag{153}$$

$$\epsilon_1^{(m)} = \sqrt{\frac{1}{9\phi_1}\mathbf{E} : \left[\frac{\partial \mathcal{L}_e^0}{\partial K_1}(\mu_1^0, K_1)\right] : \mathbf{E}}. \tag{154}$$

Equation (153) must be solved for $\epsilon_1^{eq}$ with

$$\mu_1^0 = \frac{1}{3\epsilon_1^{eq}}\frac{\partial \varphi_1}{\partial \epsilon_1^{eq}}\left(\epsilon_1^{eq}\right).$$

If we utilize any upper bound on, or estimate for, the effective modulus tensor of a linear porous material with an isotropic matrix, then, the bound (152) would lead to a corresponding upper bound or estimate for the effective strain potential of the nonlinear porous material. If the matrix phase is incompressible ($K_1 \to \infty$), so that the effective modulus and compliance tensors of the linear comparison porous material can be written as

$$\mathcal{L}_e^0 = \mu_1^0\hat{\mathcal{L}}, \quad \mathcal{M}_e^0 = \left(\mu_1^0\right)^{-1}\hat{\mathcal{M}},$$

where $\hat{\mathcal{L}}$ and $\hat{\mathcal{M}}$ are two microstructural tensors that are independent of $\mu_1^0$, then, the estimate (152) for $\mathcal{H}_e$ and the corresponding estimate for $\mathcal{H}_e^*$ reduce to

$$\mathcal{H}_e(\mathbf{E}) \leq \phi_1\varphi_1\left(\epsilon_1^{eq}\right), \quad \mathcal{H}_e^*(\mathbf{\Sigma}) \geq \phi_1\psi_1\left(\sigma_1^{eq}\right), \tag{155}$$

where

$$\epsilon_1^{eq} = \sqrt{\frac{1}{3\phi_1} \mathbf{E} : \hat{\mathcal{L}} : \mathbf{E}}, \quad \sigma_1^{eq} = \sqrt{\frac{3}{\phi_1} \boldsymbol{\Sigma} : \hat{\mathcal{M}} : \boldsymbol{\Sigma}}. \tag{156}$$

### 4.7.1.1  Two-Point Bounds

When the distribution of the pore phase is statistically isotropic, the linear Hashin–Shtrikman bound (see Chapter 7 of Volume I) leads to a corresponding upper (lower) bound for $\mathcal{H}_e$ ($\mathcal{H}_e^*$) with

$$\epsilon_1^{eq} = \sqrt{\frac{4}{\phi_2} E_m^2 + \left(1 + \frac{2}{3}\phi_2\right)^{-1} E_{eq}^2}, \tag{157}$$

$$\sigma_1^{eq} = \frac{1}{\phi_2} \sqrt{\frac{9}{4}\Sigma_m^2 + \left(1 + \frac{2}{3}\phi_2\right) \Sigma_{eq}^2}, \tag{158}$$

which was first derived Ponte Castañeda (1991a) and Suquet (1992). It was also derived as an *ad hoc* estimate (not a bound) by Qiu and Weng (1992) by estimating the stress in the matrix from the energy in the porous material. If, on the other hand, the voids' shapes and distribution are cylindrical with circular cross section, one obtains a Hashin–Shtrikman-type bound given by (152) with (Suquet, 1992)

$$\sigma_1^{eq} = \sqrt{\frac{1}{1 - \phi_2}\left[\Sigma_{eq}^2 + \frac{3}{2}\phi_2(\Sigma_{11}^2 + \Sigma_{22}^2) + 3\phi_2(\Sigma_{13}^2 + \Sigma_{23}^2 + \Sigma_{12}^2)\right]}, \tag{159}$$

where the axis of symmetry has been taken to be aligned with the $x_3$ direction.

Another important case is when one of the aspect ratios of the voids approaches zero, leading to cracks, in which case, $\phi_2 \to 0$. When the cracks are penny shaped, aligned, and distributed isotropically, one obtains the Hashin–Shtrikman-type bound (152) with

$$\sigma_1^{eq} = \sqrt{\Sigma_{eq}^2 + \frac{3\rho}{\pi}\left[\left(1 - \frac{32}{15}\frac{\rho}{\pi}\right)^{-1}\Sigma_{33}^2 + \frac{4}{3}\left(1 - \frac{4}{15}\frac{\rho}{\pi}\right)^{-1}(\Sigma_{13}^2 + \Sigma_{23}^2)\right]}, \tag{160}$$

where $\rho = \frac{4}{3}\pi n_2 a^3$ is the crack density corresponding to $n_2$ cracks of mean radius $a$ per unit volume. The corresponding results for flat distributions of cracks, i.e., when the crack interactions are weak, which are obtained by linearizing (with respect to $\alpha_2$) Eq. (160), were first given by Suquet (1992) and Talbot and Willis (1992). When the cracks are randomly oriented and distributed isotropically, the following upper bound is obtained:

$$\epsilon_1^{eq} = \sqrt{\frac{3\pi}{\rho} E_m^2 + \left(1 - \frac{12}{25}\frac{\rho}{\pi}\right)\left(1 + \frac{8}{25}\frac{\rho}{\pi}\right)^{-1} E_{eq}^2}, \tag{161}$$

which was derived by Ponte Castañeda and Willis (1995) in the linear context.

If we use the self-consistent or effective-medium approximation estimate of Hill (1965b) and Budiansky (1965), described in Chater 7 of Volume I, it follows (Ponte Castañeda, 1991a) that expression (152) provides an estimate for $\mathcal{H}_e$ with

$$\epsilon_1^{eq} = \sqrt{\frac{1}{\phi_1}\left(\frac{1-2\phi_2}{1-\phi_2/3}\right)\left(\frac{4}{\phi_2}E_m^2 + E_{eq}^2\right)}. \tag{162}$$

### 4.7.1.2  Three-Point Bounds

One can also obtain third-order, Beran-type bounds for this class of materials. If one uses the Milton (1982) simplified form of the third-order bounds for linear elastic materials (see also Chapter 7 of Volume I), it is straightforward to derive a third-order upper bound for porous materials with statistically isotropic microstructures, given by (152), with

$$\epsilon_1^{eq} = \sqrt{\frac{4\zeta_1}{\phi_2}E_m^2 + \eta_1\left(\eta_1 + \frac{2}{3}\phi_2\right)^{-1}E_{eq}^2}, \tag{163}$$

where $\eta_1$ and $\zeta_1$ are two microstructural parameters defined and described in Chapters 4 and 7 of Volume I. Bounds of this type for nonlinear materials were first proposed by Ponte Castañeda (1992a, 1997). Note that when $\zeta_1 = \eta_1 = 1$, the bound (152), together with Eq. (163), reduce to the Hashin–Shtrikman upper bound [together with Eqs. (157) and (158)], but the bound (152) is generally tighter than the Hashin–Shtrikman bound for $\zeta_1 \neq 1$ and $\eta_1 \neq 1$.

One may also utilize the second-order theory of Ponte Castañeda (1996a) in order to derive certain results for porous materials with an incompressible matrix and statistically isotropic microstructures (or isotropic distributions of spherical pores), provided that the pores are also incompressible, so that the material as a whole is incompressible ($E_m = 0$). This would be the case if the pores are saturated with an incompressible fluid. Assuming isotropy of the matrix, as characterized by the function $f_1$ in

$$w_i(\epsilon) = \frac{9}{2}K_i\epsilon_m^2 + f_i(\epsilon_{eq}^2), \tag{164}$$

where $f_i$ characterizes the shear modulus of phase $i$, and letting $K_1 \to \infty$, the second-order estimate (145) for such fluid-saturated porous materials is written as

$$\tilde{\mathcal{H}}_e(\mathbf{E}) = \phi_1 f_1\left[(1+\phi_2\omega)^2 E_{eq}^2\right] - \phi_1\phi_2\omega(1+\phi_2\omega)(f_1)'\left[(1+\phi_2\omega)^2 E_{eq}^2\right]E_{eq}^2, \tag{165}$$

where $\langle\epsilon\rangle_1 = (1+\phi_2\omega)\mathbf{E}$ is obtained from Eq. (150) for the average strain in the matrix phase. In a similar manner, a self-consistent or an effective-medium approximation estimate can also be obtained (Ponte Castañeda and Suquet, 1998).

To see the application of these results, consider, as an example, the Hashin–Shtrikman-type variational bounds [i.e., using Eqs. (157) and (158) in (155) and (156)] and second-order estimates (165) for statistically-isotropic porous materials. The behavior of the incompressible matrix is characterized by the power-law relation (10), so that for purely deviatoric loading conditions ($E_m = 0$), the effective

potential $\mathcal{H}_e$ can be written in the form

$$\mathcal{H}_e(\mathbf{E}) = \frac{\sigma_e^0(\theta)\epsilon^0}{m+1}\left(\frac{E_{eq}}{\epsilon^0}\right)^{m+1}, \tag{166}$$

where $\theta$ depends on the determinant of the strain, with $\theta = 0$ corresponding to axisymmetric deformation and $\theta = \pi/6$ to simple shear. In general, one finds that the second-order estimates lie below the variational bounds. Moreover, although the variational bounds are independent of the type of loading, the corresponding second-order estimates are different for such cases as uniaxial tension and simple shear, with the shear results always lying below the tensile results. In addition, the difference between the shear and tensile results becomes progressively larger, as the level of nonlinearity increases, with the second-order estimates remaining close to the variational estimates for tension, but predicting sharper drops in the load-carrying capacity of the porous material in shear. As first pointed out by Drucker (1959), the sharper drop for large values of $n$ (tending to perfectly plastic behavior) is possible under shear loading because of the availability of localized deformation modes (i.e., slip bands) passing through the pores. There is also experimental evidence for this type of behavior (Spitzig *et al.*, 1988). On the other hand, for the axisymmetric deformation mode, the plastic deformation is diffused through the matrix (Duva and Hutchinson, 1984), and the differences between the variational and second-order estimates are relatively small (Ponte Castañeda, 1996a). It must, however, be emphasized that the second-order procedure can capture more accurately the anisotropy of the localized deformation fields by means of the use of the anisotropic tangent modulus tensors (Ponte Castañeda, 1992a).

## 4.7.2  Rigidly Reinforced Materials

Let us now discuss composite materials with isotropic nonlinear matrix phases, reinforced by a rigid phase. The phase strain and stress potentials are assumed to be given by Eqs. (4) and (5). Designating the matrix as phase 1, Eq. (108) becomes

$$\mathcal{H}_e^*(\boldsymbol{\Sigma}) \geq \frac{\phi_1}{2K_1}\left[\sigma_1^{(m)}\right]^2 + \phi_1\psi_1(\sigma_1^{eq}), \tag{167}$$

where

$$\sigma_1^{eq} = \sqrt{\frac{3}{\phi_1}\boldsymbol{\Sigma} : \left[\frac{\partial\mathcal{M}_e^0}{\partial(1/\mu_1^0)}(\mu_1^0, K_1)\right] : \boldsymbol{\Sigma}}, \tag{168}$$

$$\sigma_1^{(m)} = \sqrt{\frac{1}{\phi_1}\boldsymbol{\Sigma} : \left[\frac{\partial\mathcal{M}_e^0}{\partial(1/K_1)}(\mu_1^0, K_1)\right] : \boldsymbol{\Sigma}}. \tag{169}$$

Equation (168) must be solved for $\sigma_1^{eq}$ with

$$\frac{1}{\mu_1^0} = \frac{3}{\sigma_1^{eq}}\frac{\partial\psi_1}{\partial\sigma^{eq}}(\sigma_1^{eq}).$$

Use of any lower bound or any estimate for the effective compliance tensor of a rigidly-reinforced material with an isotropic matrix then leads to corresponding lower bounds and estimates for the effective stress potential of the corresponding nonlinear, rigidly-reinforced composites. When the matrix phase is also incompressible (i.e., when $K_1 \to \infty$), the resulting material is also incompressible and the corresponding estimates for $\mathcal{H}_e$ and $\mathcal{H}_e^*$ can be written in a form similar to (155) and (156), with $E_m = 0$, in terms of appropriate microstructural tensors $\hat{\mathcal{L}}$ and $\hat{\mathcal{M}} = (\hat{\mathcal{L}})^{-1}$. For example, the Hashin–Shtrikman estimates can be interpreted as appropriate variational estimates for particulate microstructures, and thus the corresponding nonlinear results can be thought of as appropriate variational estimates for particulate microstructures. Thus, (155) and (156), with the inequality replaced by an approximate equality, yield estimates for $\mathcal{H}_e$ and/or $\mathcal{H}_e^*$. In particular, for spherical particles that are distributed with statistically-isotropic symmetry, the following estimate should be used:

$$\epsilon_1^{eq} = \frac{1}{\phi_1}\sqrt{1 + \frac{3}{2}\phi_2}\, E_{eq}, \tag{170}$$

This "lower estimate" was proposed by Ponte Castañeda (1991b, 1992a) for isotropic, rigidly-reinforced composites and generalized by Talbot and Willis (1992) for anisotropic materials. Talbot and Willis (1992) and Li *et al.* (1993) also presented predictions for aligned spheroidal inclusions. Gărăjeu and Suquet (1997) also discussed an application to rigidly-reinforced materials.

### 4.7.2.1    Two-Point Bounds

As discussed in Chapter 7 of Volume I, in the case of statistically-isotropic morphologies, the Hashin–Shtrikman upper bounds for linear elastic materials with arbitrary microstructures are unbounded, and therefore the corresponding upper bounds for $\mathcal{H}_e$ are also unbounded. Physically, this is due to the fact that statistical isotropy does not exclude the possibility of formation of a sample-spanning percolation cluster of rigid materials. However, for particulate microstructures (which, at least for small enough volume fractions of the inclusions, exclude the possibility of formation of rigid percolation clusters), one can obtain finite upper bounds for the effective modulus tensor of rigidly-reinforced materials. Linear Hashin–Shtrikman bounds of this type were derived by Hervé, Stolz, and Zaoui (HSZ) (1991) for coated-spheres models, and for more general morphologies by Bornert *et al.* (1996). In the coated-spheres model (see also Chapters 3, 4 and 7 of Volume I) the material consists of composite spheres that are composed of a spherical core of elastic stiffness tensor $\mathbf{C}_2$ and radius $a$, surrounded by a concentric shell of elastic stiffness tensor $\mathbf{C}_1$ with an outer radius $b > a$. The ratio $a/b$ is fixed, and the volume fraction $\phi_2$ of inclusions in $d$ dimensions is given by $\phi_2 = (a/b)^d$. The composite spheres fill the space, implying that there is a sphere size distribution that extends to infinitesimally-small spheres.

Bornert (1996) pointed out that in fact the (lower) bounds for coated-spheres model can be interpreted as rigorous bounds for materials with the larger class

of particulate microstructures considered by Ponte Castañeda and Willis (1995). When both the shapes of the rigid inclusions and their distribution are spherical, the upper bound can be explicitly computed from the corresponding linear bound of Hashin (1962) and HSZ (which are identical in this case), and is given by

$$
\epsilon_1^{eq} = \sqrt{\frac{1}{\phi_1}\left\{1 + \phi_2\left[\frac{2}{5}(1-\phi_2) - \frac{21\phi_2(1-\phi_2^{2/3})^2}{10(1-\phi_2^{7/3})}\right]^{-1}\right\}} E_{eq}, \qquad (171)
$$

which was first derived by Suquet (1993a).

For fiber-reinforced materials with cylindrical inclusions that have circular cross sections, the following result (Li *et al.*, 1993) is obtained:

$$
\sigma_1^{eq} = \sqrt{\frac{3}{1+\phi_2}\left[\Sigma_{12}^2 + \frac{1}{4}(\Sigma_{11}-\Sigma_{22})^2 + \Sigma_{13} + \Sigma_{23}\right]}, \qquad (172)
$$

where the axis of symmetry was assumed to be along the $x_3$ direction. Such materials are inextensible along the fiber direction and can only support shear in the transverse and longitudinal directions. Another important case is one in which one of the aspect ratios of the rigid inclusions approaches zero, leading to disk-like inclusions (in this limit, $\phi_2 \to 0$). If the disks have circular cross sections, and are aligned and distributed isotropically, the following result is obtained:

$$
\sigma_1^{eq} = \left(\Sigma_{eq}^2 - \frac{4\rho}{\pi}\left\{\frac{1}{3}\left(1 + \frac{4}{5}\frac{\rho}{\pi}\right)^{-1}\left[\Sigma_{33} - \frac{1}{2}(\Sigma_{11}+\Sigma_{22})\right]^2\right.\right.
$$
$$
\left.\left. + 2\left(1 + \frac{24}{15}\frac{\rho}{\pi}\right)^{-1}\left[\Sigma_{12}^2 + \frac{1}{2}(\Sigma_{11}-\Sigma_{22})^2\right]^2\right\}\right)^{1/2}. \qquad (173)
$$

where, as before, $\rho = \frac{4}{3}\pi n_2 a^3$ is the disk density corresponding to $n_2$ disks (per unit volume) of mean radius $a$. The corresponding results for flat distributions of disks (i.e., when the disk interactions are weak) are obtained by linearizing Eq. (173), and were first given by Talbot and Willis (1992) and Li *et al.* (1993). If the disks are randomly oriented and distributed isotropically, one obtains (Ponte Castañeda and Willis, 1995)

$$
\epsilon_1^{eq} = \sqrt{\left(1 + \frac{12}{15}\frac{\rho}{\pi}\right)\left(1 - \frac{8}{15}\frac{\rho}{\pi}\right)^{-1}} E_{eq}. \qquad (174)
$$

### 4.7.2.2   Three-Point Bounds and Estimates

Utilizing the third-order bounds for linear, elastic materials (see Chapter 7 of Volume I), it is straightforward to derive the following third-order estimates for rigidly-reinforced materials with statistically isotropic microstructures:

$$
\epsilon_1^{eq} = \frac{1}{\phi_1}\sqrt{1 + \frac{3}{2}\left(\frac{11\zeta_1 + 5\eta_1}{21\eta_1 - 5\zeta_1}\right)\phi_2} \, E_{eq}, \qquad (175)
$$

which, in the limit, $\zeta_1 = \eta_1 = 1$, reduces to the Hashin–Shtrikman estimate. Because for these values of $\zeta_1$ and $\eta_1$, the Beran upper and lower bounds coincide, Eq. (175) is a rigorous upper bound for nonlinear materials with $\zeta_1 = \eta_1 = 1$.

One may also obtain the self-consistent or effective-medium approximation estimates for statistically-isotropic microstructures by utilizing the estimates of Hill (1965b) and Budiansky (1965) (see also Chapter 7 of Volume I). The result is then given by (Ponte Castañeda, 1991)

$$\epsilon_1^{eq} = \sqrt{\frac{1}{\phi_1}\left(1 - \frac{5}{2}\phi_2\right)^{-1}} \, E_{eq}. \tag{176}$$

The results presented so far represent rigorous bounds for the effective mechanical properties of rigidly-reinforced materials. The corresponding second-order Hashin–Shtrikman estimate, Eq. (170), and the self-consistent estimate, Eq. (176), for statistically-isotropic microstructures (or isotropic distributions of spherical voids) were derived by Ponte Castañeda (1996a) and Ponte Castañeda and Nebozhyn (1997) for materials with an isotropic, incompressible matrix phase, as characterized by the function $f_1$ in Eq. (164).

### 4.7.3  Completely Plastic Materials

Another class of nonlinear composites for which explicit analytical results are available consists of two-phase, rigid, perfectly plastic materials with isotropic constituents. In certain limits of this class, the associated nonlinear equations for the comparison moduli or reference strain in the phases can even be solved exactly. For example, consider a two-phase material with isotropic, ductile phases governed by the Von Mises criterion,

$$\sigma^{eq}(\mathbf{x}) \leq \sigma_i^0, \quad \text{in phase } i. \tag{177}$$

Then, the variational representation (95) can be utilized for deriving explicit results for some cases of practical interest, which are now briefly discussed.

If the material is isotropic, the dissipation potential $\mathcal{H}_e$ depends only on the second and third invariants of the strain and, due to homogeneity, can be written as

$$\mathcal{H}_e(\mathbf{E}) = \sigma_e^0(\theta) E_{eq}, \tag{178}$$

where $\theta$ depends on the determinant of the normalized deviatoric strain. Use of a piecewise constant shear modulus $\mu^0(\mathbf{x})$ in (95) then leads to the following upper bound for $\sigma_e^0$:

$$\frac{\sigma_e^0}{\sigma_2^0} \leq \inf_{\mu_1^0/\mu_2^0 \geq 0} \left\{ \sqrt{\left(\frac{\mu_e^0}{\mu_2^0}\right)\left[\phi_1\left(\frac{\sigma_1^0}{\sigma_2^0}\right)^2 \frac{\mu_2^0}{\mu_1^0} + \phi_2\right]} \right\}, \tag{179}$$

which is independent of $\theta$ and therefore of the third invariant. Rigorous upper bounds for the effective flow stress $\sigma_e^0$ of isotropic, two-phase materials can then be obtained by incorporating upper bounds for the effective shear modulus $\mu_e$ of the linear comparison material in (179). For example, assuming that $\sigma_1^0 \geq \sigma_2^0$, the

Hashin–Shtrikman upper bound for a $d$-dimensional material is given by (Ponte Castañeda and deBotton, 1992; Suquet, 1993a; Olson, 1994)

$$\frac{\sigma_e^0}{\sigma_2^0} = \frac{(d+2)\phi_2}{d+2\phi_2} + \frac{d\phi_1}{d+2\phi_2}\sqrt{\left(\frac{\sigma_1^0}{\sigma_2^0}\right)^2 + \frac{2}{d}\phi_2\left[\left(\frac{\sigma_1^0}{\sigma_2^0}\right)^2 - 1\right]}. \qquad (180)$$

Similarly, the Hashin–Shtrikman estimates for spherical inclusions, distributed with statistical isotropy, can also be derived. In this case, estimates for the effective flow stress of the material can be obtained by using the appropriate estimates for $\mu_e$ for this class of microstructures. If the estimate for $\mu_e$ is accurate for arbitrary contrast $\mu_1/\mu_2$, then, the resulting expression for $\sigma_e^0$ is likely to be an upper bound for the same class of microstructures. For example, the Hashin–Shtrikman lower bound is appropriate for describing the effective shear modulus of dispersions of spherical inclusions (phase 2) in a matrix (phase 1) at moderate volume fractions of inclusions which, as mentioned earlier in this chapter (see also Chapter 7 of Volume I), is a rigorous upper bound for materials with the microstructural parameters $\zeta_1 = \eta_1 = 1$. When used in (179), the optimization procedure can be carried out analytically. Assuming that $\sigma_2^0 \geq \sigma_1^0$, the estimate for the overall flow stress resulting from this calculation is given by (Ponte Castañeda and deBotton, 1992)

$$\frac{\sigma_e^0}{\sigma_1^0} = \frac{(d+2)\phi_2}{d+2\phi_2}\frac{\sigma_2^0}{\sigma_1^0} + \frac{d\phi_1}{d+2\phi_2}\sqrt{1 + \frac{2\phi_2}{d}\left[1 - \left(1 - \frac{\sigma_2^0}{\sigma_1^0}\right)^2\right]}, \qquad (181)$$

with

$$\frac{\sigma_1^0}{\sigma_2^0} \geq \frac{2}{d+2}\sqrt{1 + \frac{1}{2}d\phi_2},$$

and

$$\frac{\sigma_e^0}{\sigma_1^0} = \sqrt{1 + \frac{1}{2}d\phi_2}, \qquad (182)$$

where

$$\frac{\sigma_1^0}{\sigma_2^0} \leq \frac{2}{d+2}\sqrt{1 + \frac{1}{2}d\phi_2}.$$

These results, which may be interpreted as *approximate estimates* for materials with particulate microstructures, are upper bounds for composites with morphologies for which the Hashin–Shtrikman lower bound for $\mu_e$ is exact (for example, sequentially-laminated composites; see Chapter 2). Note that the estimate (182) predicts that the strengthening effect of the inclusions (when they are stronger than the matrix) saturates after a certain finite increase in the strength of the inclusions. This is a consequence of the non-hardening character of the matrix phase, which would be expected to carry all the deformation, for sufficiently strong (but still non-rigid) inclusions.

For unidirectional materials with transverse isotropy (or for fiber-reinforced composites with circular fibers of phase 2 dispersed isotropically), the expression for the effective yield function reduces to

$$
P_e(\mathbf{\Sigma}) \geq \max_{y \geq 0} \left\{ \left[ \phi_2 + \phi_1 \left( \frac{\sigma_1^0}{\sigma_2^0} \right)^2 y \right]^{-1} \mathbf{\Sigma} \cdot \hat{\mathcal{M}}_e^0 (y) \cdot \mathbf{\Sigma} - \left( \sigma_2^0 \right)^2 \right\}, \quad (183)
$$

where $y = \mu_1^0 / \mu_2^0$, and the tensor $\hat{\mathcal{M}}_e^0 = \mu_1^0 \mathcal{M}_e^0$ is the (normalized) effective compliance of the fiber-reinforced linear comparison material with incompressible and isotropic phases. In general, this result requires numerical computation, but for transverse and longitudinal shear, the result simplifies to expressions similar in form to (180)–(182) with $d = 2$. Similarly, for (axisymmetric) uniaxial tension, one obtains

$$
\sigma_e^0 = \phi_1 \sigma_1^0 + \phi_2 \sigma_2^0, \quad (184)
$$

in agreement with the Voigt estimate. These results are due to Ponte Castañeda and deBotton (1992) and Moulinec and Suquet (1995); see also deBotton (1995).

In a similar way, one may obtain second-order estimates for the effective mechanical properties of this class of nonlinear materials. For example, for two-phase, rigid, perfectly plastic materials with statistically-isotropic microstructures [or with isotropic distributions of spherical inclusions (phase 2) in a matrix (phase 1)], the second-order estimates (145) for $\mathcal{H}_e$ can be simplified. The result, for simple shear loading conditions, is given by

$$
\frac{\sigma_e^0}{\sigma_1^0} = \begin{cases} 1 - \dfrac{1}{2}(1 + \phi_2)\left(1 - \dfrac{\sigma_2^0}{\sigma_1^0}\right), & \text{if} \quad \dfrac{\sigma_2^0}{\sigma_1^0} < 1, \\[3mm] 1, & \text{if} \quad \dfrac{\sigma_2^0}{\sigma_1^0} \geq 1. \end{cases} \quad (185)
$$

An identical result is obtained for fiber-reinforced microstructures with transverse isotropy loaded in transverse shear. We should point out that the small-contrast expansion described in Section 4.6.1 diverges for simple shear loading, whereas, as indicated by Eq. (185), the corresponding second-order estimate does not.

Finite-element computations carried out by Suquet (1993a), for particle-reinforced materials with inclusion volume fraction $\phi_2 = 0.15$, indicate that, although the two types of nonlinear estimates obtained from the linear Hashin–Shtrikman lower bound exhibit the same general trends, the second-order estimates are in closer agreement with the numerical results. Moreover, the variational estimates lie above the numerical results, consistent with the fact that the variational estimates are expected to overestimate the effective yield strength of the composite at this value of $\phi_2$. The nonlinear estimate obtained from the linear Hashin–Shtrikman upper bound lies below the microstructure-independent Voigt (one-point) upper bound (see Section 4.3.1), and is such that the second-order estimate lies below the variational estimate, which is known to be a rigorous bound for all statistically-isotropic microstructures.

One may also compare the results of numerical simulations by Moulinec and Suquet (1995) for the effective yield strength of fiber-reinforced materials with the corresponding predictions (183) obtained from the variational method. These authors considered cylindrical fibers (phase 2) with circular cross section and aligned with the $x_3$ axis, distributed randomly in a matrix (phase 1). The overall stresses considered by these authors consisted of the superposition of uniaxial tension and transverse shear,

$$\mathbf{\Sigma} = \Sigma_{11}(\mathbf{e}_1 \otimes \mathbf{e}_1 - \mathbf{e}_2 \otimes \mathbf{e}_2) + \Sigma_{33}\mathbf{e}_3 \otimes \mathbf{e}_3.$$

Various contrast ratios for the strengths of the two phases were investigated: $\sigma_2^0/\sigma_1^0 = 0.5, 1.1, 2, 3, 5$, and 10. For $\sigma_2^0/\sigma_1^0 = 2$, 11 different realizations were used, while for the other ratios, the computations were performed on a single realization, representative of the average of the predictions over the entire set of configurations for $\sigma_2^0/\sigma_1^0 = 2$, a configuration that approaches transverse isotropy, with its overall strain/stress response being close to the mean response of all the realizations, both under multiaxial loading and uniaxial tension. The results are shown in Figure 4.2. The agreement between the numerical simulation results and the variational estimates (183) is good. In particular, the variational estimates (183) capture rather well the flat sectors on the yield surfaces.

For the cases that involve sufficiently strong fibers, the shape of the observed extremal surfaces was found to be bimodal in character. Bimodal surfaces were used by Hashin (1980), Dvorak and Bahei-El-Din (1987), and de Buhan and Taliercio (1991) for describing the initial yield or the flow surface of unidirectional composites. The numerical and variational results are consistent with these models and with experimental observations (Dvorak *et al.*, 1988). The numerical calculations



FIGURE 4.2. Effective yield strength $\Sigma_{11}$ of composites with cylindrical fibers aligned in the $x_3$-direction (perpendicular to the plane of this page) with volume $\phi_2$. The curves are, from left to right, for $\sigma_2^0/\sigma_1^0 = 0.5, 1.1, 2, 3, 5$, and 10. Symbols represent the results of numerical simulations for randomly isotropic configurations (averaged over 11 realizations), while the curves show the predictions of the variational method in which the Hashin–Shtrikman lower bound for the linear comparison material has been used (after Moulinec and Suquet, 1995).

also suggest closed-form expression for the bimodal surface:

$$|\Sigma_{33}| \le \phi_1 \sigma_1^0 \left[ 1 - \left( \frac{\Sigma_{11}}{K_1} \right)^2 \right]^{1/2} + \phi_2 \sigma_2^0 \left[ 1 - \left( \frac{\Sigma_{11}}{K_e} \right)^2 \right]^{1/2}, \qquad (186)$$

where $K_1 = \sigma_1^0/\sqrt{3}$ is the in-plane shear strength of phase 1, and $K_e$ is the in-plane shear strength of the composite, which can either be fitted to the numerical simulations (Moulinec and Suquet, 1995), or be taken from the prediction of the variational procedure used with the Hashin–Shtrikman lower bound (Ponte Castañeda and deBotton, 1992): $K_e = (1/\sqrt{3})\sigma_e^0$, with $\sigma_e^0$ being given by (181) and (182) with $d = 2$. In the second case, the agreement with the predictions of the variational procedure for the full yield surface was found to be quite good.

## 4.8    Other Theoretical Methods

In addition to what was discussed above, several other theoretical methods have been proposed over the past 30 years for predicting the overall effective mechanical properties of nonlinear materials. Two noteworthy of such methods are the (classical) secant method developed by Chu and Hashin (1971), Berveiller and Zaoui (1979), and Tandon and Weng (1988), and the incremental method originally proposed by Hill (1965a) in conjunction with the self-consistent or the effective-medium approximation method. Briefly, the secant method consists of writing down the constitutive relation in phase $i$ with the secant tensor of phase $i$, evaluated at the average strain $\langle \epsilon \rangle_i$. In the incremental method, one writes down the constitutive law of phase $i$ in the form $\dot{\sigma} = \mathcal{L}_i^{(t)}(\langle \epsilon \rangle_i) : \dot{\epsilon}_i$, where $\mathcal{L}_i^{(t)}$ is now the tensor of instantaneous or tangent moduli of the phase, given by the second derivative of the energy $w_i$ with respect to the strain.

Two-phase, incompressible, power-law materials with the same exponent provide an important test for comparing the different models. Particulate power-law materials were considered by Ponte Castañeda and Willis (1988) in the context of the Talbot–Willis procedure, by Ponte Castañeda (l991a) and Suquet (1993a) in the context of the variational method with a linear comparison material, and by Ponte Castañeda (1996a) in the context of the second-order procedure. Granular microstructures were also considered by these groups, as well as by Gilormini (1995), who compared the different methods using the self-consistent method for estimating the effective properties of the linear comparison material. He pointed out that the predictions of the incremental and classical secant method can violate the rigorous variational upper bound for isotropic materials. Michel (1996) proposed a nonlinear extension of the self-consistent method for power-law materials.

Consider, as an example, two-phase materials with particulate microstructure. Both phases are characterized by Eq. (10) with the same exponent $m$ but different stresses $\sigma_i^0$. Suppose that the material consists of inclusions (phase 2) that are distributed randomly in a softer matrix (phase 1). If the volume fraction $\phi_2$ of the inclusions is not too large, the Hashin–Shtrikman lower bound provides accurate

estimates for the effective linear properties of the comparison material with the same microstructure as that of the nonlinear material. The material itself is a power-law composite with the same exponent as the individual phases, and is, in addition, incompressible. Under the assumption of statistical isotropy, the effective potential is a function of the second and third invariant of the average strain $\mathbf{E}$ and, by homogeneity, is given by Eq. (166). The variational bounds, derived above for power-law materials provide bounds for $\sigma_e^0$ that are independent of the parameter $\theta$ of Eq. (166), whereas the estimates provided by the second-order theory do depend on this parameter. It can then be shown (Ponte Castañeda and Suquet, 1998) that the incremental and secant procedures lead to the stiffest predictions, whereas the variational and second-order methods provide more compliant predictions. In particular, since, as already noted in Chapter 2 (see also Chapters 4 and 7 of Volume I), the linear Hashin–Shtrikman lower bound is attained by certain particulate microstructures, the variational estimates are actually upper bounds for the nonlinear composites with the same type of microstructure. Therefore, both the corresponding secant and incremental estimates violate this bound, whereas the second-order estimates do not. In fact, the incremental estimates violate even the Hashin–Shtrikman upper bound for statistically-isotropic microstructures, at sufficiently large values of the exponent $n$. This is somewhat unexpected, as this type of bound is known to correspond to the opposite type of microstructure, with the stronger material occupying the matrix phase.

A similar observation was made by Gilormini (1995) in the context of the self-consistent estimate (instead of the Hashin–Shtrikman lower bound). These results indicate that the tendency of the incremental model to approach the Voigt (one-point) bound (see Section 4.3.1) when $m \rightarrow 0$ is not due to the approximate nature of the self-consistent method, but is because of the shortcomings of the incremental method itself. Let us emphasize again that of the four nonlinear homogenization procedures described above, only the second-order theory yields estimates that are exact to second order in the contrast between the properties of the phases. The other three (variational, secant, and incremental) provide estimates that are exact only to first order in the contrast.

Finally, Gibiansky and Torquato (1998b) derived approximations for the effective energy of $d$-dimensional nonlinear, isotropic, elastic dispersions. These approximations are similar to those described in Sections 2.2.2.1 and 2.2.2.2, derived by Gibiansky and Torquato (1998a), for the effective conductivity of the materials with the same morphology. In addition, Gibiansky and Torquato (1998b) derived cross-property relations that link the effective energy of nonlinear materials with their effective conductivity.

## 4.9  Critique of the Variational Procedure

A valid criticism of the variational procedures is that they rely, from the very beginning, on the assumption that the mechanical behavior of the constituent phases can be described by a potential, which is not the case for many nonlinear (usually

elasto-plastic) materials. A partial response to this criticism was provided by Ponte Castañeda and Suquet (1998) who argued that, at least for certain loading conditions of practical interest, it is possible to use a deformation theory of plasticity, instead of a flow theory, to describe the mechanical properties of the constituent materials. This substitution is rigorous only when the loading is radial and monotonic at every point **x** in the volume element $\Omega$, but it may also be appropriate for small deviations from proportionality (Budiansky, 1959). The assumption of proportionality is rarely met and deviations from radial paths are likely to be the rule. Nevertheless, numerical simulations of the transverse response of nonlinear matrices, reinforced by aligned continuous fibers, suggest that, even though local deviations from this assumption are actually observed and found to affect the local stress and strain fields, they seem to have little influence on the overall stress-strain response of the material under monotonic loading, implying that using a deformation theory for the constituents can be a good approximation for materials that are subjected to a monotonic radial loading, such as uniaxial tension. Strictly speaking, although this model is not applicable to general loadings, its predictions for those loadings to which it is applicable are much more accurate (Suquet, 1997) than those of theories that allow for more general loadings, such as the incremental method or the transformation field analysis (Dvorak, 1992).

However, use of a deformation theory for non-monotonic loadings is not appropriate. Instead, one must use a flow theory for which a variational method cannot be utilized. The variational method can still yield useful insight into how to construct approximate effective constitutive relations, expressed in terms of two thermodynamic potentials, the free energy for reversible effects and the dissipation potential for irreversible phenomena (see, for example, Rice, 1970; Mandel, 1972; Germain *et al.*, 1983).

## Summary

Several continuum approaches to estimating the effective nonlinear mechanical properties of multiphase materials were described and discussed. One method, due to Talbot and Willis, is based on a nonlinear extension of the Hashin–Shtrikman variational principles, while the second method, developed by Ponte Castañeda (for nonlinear isotropic materials) and Suquet (for power-law composites) utilizes new variational principles that involve a linear comparison material with the same microstructure as that of the nonlinear composite. These methods provide at least one type of rigorous bounds (i.e., upper or lower bounds). The Talbot–Willis procedure yields the bounds of the Hashin–Shtrikman type, while the Ponte Castañeda–Suquet method provides bounds and estimates of *any* type, given the corresponding bounds and estimates for the linear comparison material. In both cases, the resulting bounds and estimates are exact to first order in the contrast between the properties of the phases. A third method, also developed by Ponte Castañeda, yields estimates that are exact to second order in the contrast. The resulting estimates are not, however, bounds of any type.

Despite this considerable progress, much remains to be done, especially since it appears that the constitutive laws that characterize the behavior of many materials are rather complex. In addition, true second-order bounds, i.e., those that are exact to second order in the phase contrast, remain to be derived. When the deformations are finite, a material may undergo microstructural evolution. An example is deformations that are present in metal-forming processes. Little is known about modeling and predicting the mechanical properties of such evolving materials. Finally, no discrete model of the types that have been described throughout this book has been developed for studying the mechanical properties of nonlinear heterogeneous materials. This research field is wide open.

# Part II

# Fracture and Breakdown of Heterogeneous Materials

# 5
# Electrical and Dielectric Breakdown:
# The Discrete Approach

## 5.0  Introduction

Beginning with this chapter, and in the next three, we study and analyze failure and fracture of heterogeneous materials. In the present chapter, electrical and dielectric breakdown of composite materials, which constitute a set of complex, nonlinear, and non-local transport processes, are described. Their nonlinearity stems from the existence of a threshold: Below and far from the threshold nothing particularly complex happens. The laws of linear (or constitutively nonlinear) transport hold, and the electrical properties of the materials are described by the models that were described in the previous chapters and in Volume I. However, at the threshold, the materials' behavior and their transport properties abruptly change and become very complex. Note that, unlike the percolation threshold, the threshold in electrical or dielectric breakdown is not geometrical but dynamical although, as discussed below, the interplay between the heterogeneities and the dynamical threshold gives rise to a rich set of phenomena that are completely absent in the linear transport regime in the same system.

Dielectric breakdown in gases, liquids, and solids is a complex problem and has been studied for a long time. Many breakdown phenomena in gases are relatively well-understood (see, for example, Meek and Craggs, 1978), while some, such as atmospheric lightning, are more difficult to analyze, because the density, conductivity, and humidity of air are distributed inhomogeneously. Another well-known example, in addition to lightning, is surface discharges, also known as Lichtenberg figures. These phenomena are beyond the scope of our book and will not be considered.

In dielectric breakdown in solids, the material is initially non-conducting when an electric field is applied across the sample. If the field exceeds a certain threshold, the material breaks down and becomes conducting. The microscopic mechanisms of dielectric breakdown in solid materials are much more complex than those in gases since, in addition to dielectric effects, mechanical and chemical effects can also intervene and make the problem more difficult. From a practical view point, dielectric breakdown is an important phenomenon, since it limits the application of dielectrics as insulators. For this reason, dielectric breakdown in solids has received much attention over the past several decades, and has been especially studied intensively over the past decade. A well-known example of such phenomena is formation and growth of electrical trees (as in, for example, discharge treeing in polymers).

FIGURE 5.1. Schematic representation of a tree growing between two electrodes on two parallel planes (after Hill and Dissado, 1983a).

The trees themselves may not cause breakdown unless they grow so large that they span the thickness of the material. A diagrammatic representation of this phenomena is shown in Figure 5.1. We will come back to this phenomenon shortly.

Another important example is dielectric breakdown in metal-loaded dielectrics, which are disordered materials consisting of a mixture of conducting and non-conducting components. For example, solid-fuel rocket propellant is a mixture of aluminum and perchlorate particles in a polymer binder (Kent and Rat, 1985). It has been reported that the breakdown field of this material decreases significantly by the presence of the aluminum particles, and is also a strong function of the volume fraction of the constituent particles. Dielectric breakdown of such composite solids is dominated by space charge effects due to the large electric fields near any sharp metal tips occurring in the composite, and thus the composite is unusually sensitive to breakdown. Recall that about a decade ago the solid fuel of a United States Air Force rocket experienced dielectric breakdown, with the fuel becoming electrically conductive, setting the rocket on fire.

Electrical breakdown occurs when the current through a conducting medium causes an irreversible resistance change in the medium. In this phenomenon, the material is initially conducting. The failure occurs when the current density flowing in the material exceeds a threshold value at and beyond which the material becomes insulating. Unlike dielectric breakdown, the mechanism of electrical failure is well-understood; it is merely Joule effect which causes degradation of metallic interconnects (or the metal lines) which, due to electromigration phenomenon, lose their conducting properties. Note that in this phenomenon the material behaves precisely like a fuse, which is broken when the applied voltage exceeds a certain limit. Electrical breakdown is a major obstacle to development of nano-size devices. Experimental realizations of electrical and dielectric breakdown in

metal-insulator films, with a view to explain them in terms of the statistical physics of disordered media, were reported by Yagil *et al.* (1992, 1993) among others. Hill and Dissado (1983b) analyzed the older experimental data. We will come back to these experiments later in this chapter.

Another important phenomenon that belongs to this class of problems is electromigration failure in polycrystalline metal films (see, for example, Huntington, 1975; Ho and Kwok, 1989). If a high current density passes through a thin metal film, collisions between the conduction electrons and the metal ions result in drifting of the ions and their electromigration. If there is a divergence in the flux of the ions at some points, voids nucleate, grow and overlap with each other until conduction ceases and the film suffers electrical breakdown (see, for example, Rodbell *et al.*, 1987). This phenomenon is particularly important in integrated circuits, where the continuing miniaturization of the circuits exposes the conducting thin metal films to increasingly large current densities. Under such conditions, electromigration failure decreases the circuit lifetime which is unacceptable from an economical view point.

Throughout this book, both in Volume I and in the present Volume, we have grouped the models for any phenomenon of interest to us into two classes—the continuum models and the discrete models. In this chapter, we deviate from this general approach because the continuum models of electrical and dielectric breakdown of heterogeneous solid materials are well-documented (see, for example, Whitehead, 1951; O'Dwyer, 1973; see also Niklasson, 1989a; Dissado and Fothergill, 1992; Ohring, 1998, for more recent references); hence, the best we could do would be summarizing these works, an unwise action. In addition, as will be discussed in this chapter, many phenomena associated with electrical and dielectric breakdown have a vector analogue in brittle fracture of solids, for which many continuum models have been developed that will be described and discussed in detail in Chapter 7. Thus, we restrict our discussion of the continuum models to a few recent efforts that utilized extensive numerical solution of the discretized continuum equations in order to study the breakdown phenomena in strongly-disordered solids. On the other hand, over the past several decades several discrete models of breakdown of heterogeneous materials have been developed. These models are either stochastic or completely deterministic. Their general features for modeling both the electrical and dielectric breakdown are the same, and in fact, with appropriate modifications, a model for one of the phenomena can be used for studying the other one. In this chapter, we describe these models in detail, discuss their predictions, and, whenever possible, compare the predictions to the relevant experimental data.

## 5.1    Continuum Models of Dielectric Breakdown

Typical of continuum models of dielectric breakdown are those of Garboczi (1988), who studied the problem analytically, and of Gyure and Beale (1989, 1992) who carried out a numerical study of the problem. What follows is a brief description of each model.

### 5.1.1   Griffith-like Criterion and the Analogy with Brittle Fracture

Garboczi (1988) extended the analysis of Griffith (1920) for brittle fracture (see Chapters 6 and 7) to dielectric breakdown, and derived the criterion for nucleation and development of a single conducting "crack" in an isotropic dielectric (insulating) material. The problem that one solves is one of an elliptical inclusion with dielectric constant $\epsilon'$ placed in an isotropic linear dielectric material with dielectric constant $\epsilon$. A far-field electric field $\mathbf{E}^0$ is then applied to the material, and the Laplace equation, $\nabla^2 V = 0$ is solved for the distribution of the voltage $V$ in the material, subject to the boundary conditions that far from the inclusion the electric field $\mathbf{E} \to \mathbf{E}^0$, and that the normal component of the displacement field $\mathbf{D} = \epsilon \mathbf{E}$ is continuous at the inclusion boundary. In the limit $\epsilon' \to \infty$ and fixed $\epsilon$, the latter boundary condition becomes $V = 0$ at the boundary of the inclusion.

This problem is easily solved by using elliptical cylindrical coordinates $(u, \theta, z)$ (see, for example, Jackson, 1998), where we assume that all the quantities in the $z$-direction are uniform. Then, the transformation between the $(x, y)$ and $(u, \theta)$ coordinates is given by

$$x = c \cosh u \cos \theta, \quad y = c \sinh u \sin \theta, \tag{1}$$

valid for $0 \le u < \infty$ and $0 \le \theta \le 2\pi$. The inclusion's surface is defined by $u = \beta$, where $\beta$ is a constant. If $\beta \to 0$, then the inclusion degenerates into a "crack" of length $2c$ with its tip at $x = \pm c$. The solution of the problem is given by

$$V = -cE^0 \cosh u \cos \theta + \frac{cE^0}{2C} \exp(\beta - u)(\epsilon' - \epsilon) \sinh(2\beta) \cos \theta, \quad u > \beta, \tag{2}$$

$$V = -\frac{c\epsilon E^0}{C}(\cosh \beta + \sinh \beta) \cosh u \cos \theta, \quad u < \beta, \tag{3}$$

with

$$C = \epsilon \cosh \beta + \epsilon' \sinh \beta.$$

From this solution, the components of the electric field, namely, $E_u = -\tau^{-1}\partial V/\partial u$, and $E_\theta = -\tau^{-1}\partial V/\partial \theta$, are computed, where $\tau = c(\sinh^2 u + \sin^2 \theta)^{1/2}$. One then defines a *field multiplication factor*, $E_u(\beta, 0)/E^0 = E_x(\beta, 0)/E^0$, which is given by

$$\frac{E_x(\beta, 0)}{E^0} = 1 + \frac{a(\epsilon' - \epsilon)}{a\epsilon + b\epsilon'}, \tag{4}$$

where $a$ and $b$ are the semi-major and semi-minor axes of the elliptical inclusion, respectively.

The critical question to be answered is: What is the difference $\Delta\mathcal{H}$ in the electrostatic energy between a material with and without the inclusion? If the sources of the applied field are fixed, then for the elliptical inclusion embedded in an infinite medium, one has, $\Delta\mathbf{H} = -\frac{1}{2}p_x E^0$, where $p_x$ is the $x$-component of the dipole

moment **p** of the inclusion. It is straightforward to show that

$$p_x = \frac{2\pi c^2 E^0 \epsilon}{C} (\epsilon' - \epsilon)(\cosh \beta + \sinh \beta) \sinh\left(\frac{1}{2}\beta\right). \tag{5}$$

Therefore, in the limit of a conducting ($\epsilon' \to \infty$ with $\epsilon$ held fixed) crack ($\beta \to 0$), one obtains

$$\Delta \mathcal{H} = -\frac{1}{2}\pi \epsilon c^2 (E^0)^2, \tag{6}$$

which is negative, indicating that the presence of the conducting crack *lowers* the energy of the system. Had we made the same computations but for a fixed potential (the common situation in practice), we would have obtained the same $\Delta\mathcal{H}$, but with the opposite sign.

Now, suppose that $\mathcal{H}_b$ is the breakdown energy required to create a unit area of conducting crack (per unit length in the $z$-direction). Then, the surface energy of the crack is $4\mathcal{H}_b c > 0$. Hence, the total energy difference between a cracked and uncracked material is given by

$$\Delta\mathcal{H} = -\frac{1}{2}\pi\epsilon c^2 (E^0)^2 + 4\mathcal{H}_b c. \tag{7}$$

The linear term of Eq. (7) will dominate if $c$ is small, implying that it is energetically unfavorable to have the conducting crack nucleate or propagate. The reverse is true for large enough $c$. The equilibrium point is thus found from $d\Delta\mathcal{H}/dc = 0$, yielding

$$E_c^0 = \sqrt{\frac{4\mathcal{H}_b}{\pi c \epsilon}}, \tag{8}$$

for the critical value of the applied far-field. Equation (8), which was first derived by Horowitz (1927), is the analogue of the Griffith's prediction for brittle fracture, which will be described in detail in Chapter 7. It is easy to show that the point represented by $E_c^0$ is a point of unstable equilibrium, and therefore for any applied field $E_0 > E_c^0$ dielectric breakdown will occur spontaneously.

Similar to brittle fracture of materials, of great interest is the region around the tip of the conducting crack where the most intense electric fields are located, and where the dielectric breakdown actually takes place. For simplicity, consider the limits $\epsilon' \to \infty$ and $\beta \to 0$, and consider the $x = c$ crack tip. One can then use a new coordinate system consisting of $r$, the distance from the crack tip, and $\Phi$, the angle from the $x$-axis. Then, in the limit $(u, \theta) \to 0$, we obtain

$$E_r = E^0 \sqrt{\frac{c}{2r}} \, \cos\left(\frac{1}{2}\Phi\right),$$

$$E_\Phi = -E^0 \sqrt{\frac{c}{2r}} \, \sin\left(\frac{1}{2}\Phi\right), \tag{9}$$

$$V = -E^0 \sqrt{2cr} \, \cos\left(\frac{1}{2}\Phi\right).$$

In analogy with brittle fracture, which is associated with a quantity referred to as the *stress-intensity factor* (see Chapters 6 and 7), we define an *electric field-intensity factor $K_I$* or, more simply, field-intensity factor,

$$K_I = \sqrt{\pi c}\, E^0,$$

in terms of which one has

$$E_r = \frac{K_I}{\sqrt{2\pi r}} \cos\left(\frac{1}{2}\Phi\right),$$

$$E_\Phi = -\frac{K_I}{\sqrt{2\pi r}} \sin\left(\frac{1}{2}\Phi\right), \tag{10}$$

$$V = -K_I \sqrt{\frac{2r}{\pi}} \cos\left(\frac{1}{2}\Phi\right).$$

Physically, $K_I$ is the amplitude of the $r^{-1/2}$ electric field singularity at the tip of the conducting crack. One may also define the *electrostatic energy release rate* $\mathcal{H}^R$ by

$$\mathcal{H}^R = \frac{d[\frac{1}{2}\pi\epsilon c^2 (E^0)^2]}{dc} = \pi\epsilon c(E^0)^2, \tag{11}$$

where $\mathcal{H}^R dc$ is the amount of electrostatic energy released when the crack extends by $dc$, with its critical value being, $\mathcal{H}_c^R = 4E_c^0$. Moreover,

$$E_c^0 = \frac{K_{Ic}}{\sqrt{\pi c}}, \tag{12}$$

where $K_{Ic}$ represents the critical value of $K_I$.

Finally, Rice (1968) developed a line integral, usually called the $J$-integral, which is independent of the contour. This quantity was originally developed for fracture of material, and its usefulness becomes evident when the contour encloses the tip of the fracture. Thus, $J$ yields $\mathcal{H}^R$, the elastic energy release rate. The $J$-integral for the elasticity problem is defined by

$$J = \oint \left[ -(\boldsymbol{\sigma} \cdot \mathbf{n}) \cdot \left(\frac{\partial \mathbf{u}}{\partial x}\right) ds + \mathcal{H}_e dy \right], \tag{13}$$

where $\mathbf{u}$ is the displacement vector, $\boldsymbol{\sigma}$ is the stress tensor, $\mathbf{n}$ is the unit vector normal to the contour, and $\mathcal{H}_e$ is the elastic energy density. Since the analogue of the stress tensor is the displacement field $\mathbf{D}$, then, the $J$-integral for the electrostatic problem is given by

$$J = \oint [-(\mathbf{D} \cdot \mathbf{n}) E_x ds + \mathcal{H} dy], \tag{14}$$

where $\mathcal{H}$ is the electrostatic energy. Garboczi (1988) showed that, similar to mechanical fracture, the $J$ integral for the electrostatic energy is independent of the path. Equation (14) was also suggested by Hoeing (1984).

The above discussions should make it clear that, many results that have been derived for brittle fracture of materials, based on the continuum models and described in Chapter 7, can be directly translated into analogous results for dielectric breakdown of materials.

### 5.1.2   *Computer Simulation*

Gyure and Beale (1989,1992) developed two-dimensional (2D) and 3D models of breakdown of metal-loaded dielectric materials. Their model consisted of a random array of perfectly conducting cylinders (in 2D) or spheres (in 3D), embedded in a uniform dielectric. The same type of boundary conditions that were used in Garboczi's work (described above), i.e., continuity of the normal component of the displacement field at the inclusions' boundaries, and the far-field condition, $\mathbf{E} \to \mathbf{E}^0$, were also utilized by Gyure and Beale. In their model, the numerical solution of the Laplace equation was obtained by the boundary element method (Kim and Karrila, 1991) described in Section 7.8.2 of Volume I. After determining the solution of the Laplace equation—the voltage distribution in the composite material—those regions of the system that are vulnerable to breakdown are identified by using the fact that the largest electric fields lie along lines joining the centers of the (cylindrical or spherical) inclusions which are closely spaced, with center-to-center lines that are nearly parallel to the applied field. It is then assumed that local breakdown occurs only between the pair of inclusions that has the largest electric field between them, and that, as a result of the breakdown, an electrical connection between the two inclusions is established, so that the two conductors attain the same electrical potential. This assumption is based on the experimental observation that such local breakdowns occur by vaporization of a portion of the metallic particles followed by resolidification as a single conductor. The voltage distribution of the new (defected) system was then calculated, the next region to suffer breakdown was identified, and so on. Various quantities of interest, such as the breakdown field, the geometry of the breakdown path, and the dielectric constant of the medium, as a function of the packing fraction were calculated by Gyure and Beale (1989,1992). These properties are further discussed below, where we describe the discrete models.

## 5.2   Discrete Models of Electrical Breakdown

We first describe and discuss discrete or lattice models of electrical breakdown of materials with percolation-type disorder. As we have been emphasizing throughout this book, the reason for considering percolation-type heterogeneities is that, they represent strong disorder and therefore any theory that provides reasonable predictions for a material with percolation disorder should be at least as accurate for other less extreme types of disorder. We will, however, discuss the effect of other types of disorder.

Thus, the problem that we wish to study is the following. We are given a disordered material, represented by a lattice in which the conductance of every bond is selected from a probability density function $f(g)$. In this state, the material is completely conducting (it contains no insulating region). We now select at random a fraction $1 - p$ of the bonds and convert them to insulators; that is, the fraction of the conducting bonds is $p$. So long as $p \gg p_c$, where $p_c$ is the percolation threshold of the lattice, the material will still be conducting, albeit with a smaller effective conductivity than when $p = 1$. We now apply a voltage $V$ across the material. If $V$ is small enough, then there would be no change in the conductivity state of the material. We now increase $V$ by an amount large enough that the first microscopic failed region (or the first failed bond in the lattice model) appears in the material. Then, the material may behave according to one of the two scenarios.

(1) As soon as the first failed region appears, the entire material may fail rapidly by an avalanche of local failed regions, *without any need for increasing the applied voltage V.*
(2) The state of the material may be such that the macroscopic failure of the material is more gradual, as the disorder distributes the current in an "equitable" way, rather than concentrating it in a few weak regions. In this case, after the first failed region appears, nothing further happens, unless we increase the applied voltage so that new failed regions can emerge.

Corresponding to any applied voltage, there exists a current that flows through the material. Since in practice macroscopic failure of the material is what one is interested in, we consider the behavior of the macroscopic current and its influence on the material. If this current exceeds a threshold $I_f$, then, the material as whole is converted to an insulator and fails. Two important questions that must be addressed by any model are as follows.

(1) How does $I_f$ depend on $p$?
(2) How does the breakdown process take place? In other words, how does the first sample-spanning path of the failed regions (or bonds in the lattice models) appear for the first time?

   Let us analyze the problem in detail for two limiting cases, namely, the dilute limit when $p \simeq 1$ (very few insulating regions), and the opposite limit, $p \simeq p_c$ (most of the sample being insulating).

## 5.2.1   The Dilute Limit

Consider first the dilute limit. In a completely conducting material (no insulating regions), the current lines are more or less parallel to each other and perpendicular to the electrode surface. Suppose now that there is only one insulating defect in the material which, for simplicity, is assumed to be spherical (or circular in 2D). In the lattice model, the corresponding defect consists of a few insulating bonds that form a cell with a regular shape, placed at the lattice's center. Then, the current lines around the defect are "deformed," leading to a current enhancement. If $i_d$

and $i_u$ are, respectively, the current densities around the defect and far from it in the unperturbed state, then, one can write

$$i_d = i_u(1 + \mathcal{E}), \tag{15}$$

where $\mathcal{E}$ is the *enhancement factor*, the magnitude of which depends on the material's morphology. For example, for an elliptical defect with major and minor axes $2a$ and $2b$, $\mathcal{E} = a/b$. The total current flowing through the material is then, $I = Si_u = Si_d/(1 + \mathcal{E})$, where $S$ is the surface area of the electrode. The first failure happens when $i_d = i_w$, where $i_w$ is the threshold current density for the failure of the sample *without* the defect. Therefore,

$$I_f = \frac{Si_w}{1 + \mathcal{E}}, \tag{16}$$

implying that the current enhancement *decreases* the failure current $I_f$. Typically, the current for the complete first failure is also the current for failure of the sample, since as soon as the regions in the vicinity of the defect fail, the current density around the new defect is further enhanced, leading to a rapid failure of the entire material. Clearly, the most damaging defects are those that are perpendicular to the current lines, and are in the form of long cylinders or rods. The probability of developing a defect depends on its shape.

In the context of the lattice models, the simplest and smallest defect is one insulating bond which is positioned parallel to the direction of the current lines and is far from the lattice's boundaries (see Li and Duxbury, 1987, for the effect of the defects that are near the boundaries of the lattice). If no defect is present in the lattice, then, $I_f = Li_w$, where $L$ is the linear size of the lattice. For a defect of size one (i.e., one bond), it is not difficult to show that, $\mathcal{E} = \pi/4$, and therefore in this case,

$$I_f = \frac{\pi}{4}Li_w. \tag{17}$$

## 5.2.2   *The Effect of Sample Size*

The most damaging defect consists of $N$ neighboring insulating bonds that are in the same plane which is perpendicular to the current lines. Thus, in 2D the most damaging defect is a line of $N$ of such insulting bonds, while in 3D it is a set of such bonds with roughly the shape of a disk. Since in 3D the current that is diverted by the $N$ bonds should be distributed over the perimeter of the defect, which is proportional to $\sqrt{N}$, one obtains

$$i_d = \begin{cases} i_w(1 + a_2 N), & \text{2D}, \\ i_w(1 + a_3\sqrt{N}), & \text{3D}. \end{cases} \tag{18}$$

The next issue to be addressed is the relation between $N$ and $L$, the linear size of the lattice. Since the probability that a bond has failed is proportional to $(1 - p)$, then, $P_N$, the probability that $N$ bonds are insulating, is given by

$$P_N \sim (1 - p)^N L^d, \tag{19}$$

where $L^d$ represents the volume of the system. The most probable, most damaging defect is formed when $P_N \sim 1$, and therefore the critical number $N_c$ for the formation of such a defect is given by

$$N_c \sim -\frac{d}{\ln(1-p)} \ln L. \tag{20}$$

Therefore, the corresponding current density $i_d$ is given by

$$i_d = \begin{cases} i_w \left[ 1 + a_2 \dfrac{-2 \ln L}{\ln(1-p)} \right], & \text{2D,} \\[2ex] i_w \left\{ 1 + a_3 \left[ \dfrac{-3 \ln L}{\ln(1-p)} \right]^{1/2} \right\}, & \text{3D.} \end{cases} \tag{21}$$

Because the total current in the system is $i L^{d-1}$, the failure current is obtained by setting $i_d = i_w$, resulting in

$$I_f = \begin{cases} \dfrac{i_w L}{1 + 2a_2 \ln L / \ln(1-p)}, & \text{2D,} \\[2ex] \dfrac{i_w L^2}{1 + 3a_3 [\ln L / \ln(1-p)]^{1/2}}, & \text{3D.} \end{cases} \tag{22}$$

The most interesting aspect of Eq. (22) is its prediction for the size-dependence of the failure current. According to this equation, the failure current per bond, $i_f = I_f / L^{d-1}$, *decreases* with the linear size of the sample in a complex way (in practice, $L$ is the ratio of the linear size of the actual sample and the typical size of the insulating defects). If $\ln(1-p)$ is not too large, then

$$I_f \sim \begin{cases} (\ln L)^{-1}, & \text{2D,} \\ (\ln L)^{-1/2}, & \text{3D.} \end{cases} \tag{23}$$

Thus, for a fixed size of the insulating defect, *the larger the sample, the smaller the failure current.*

## 5.2.3  Electrical Failure in Strongly Disordered Materials

In the limit, $p \simeq p_c$, where the material is strongly heterogeneous, the distribution of the current in the materials is controlled by the links or the red bonds of the percolation lattice model (see Chapter 2 of Volume I) that connect two multiply-connected conducting clusters. These are the bounds that, if cut, would break the sample-spanning clusters into two parts. They break down and become insulating by only a small current. Therefore, it is reasonable to assume that as $p \to p_c$, the critical current $I_f$ *vanishes*. The number of the links is proportional to $\xi_p^{d-1}$, where $\xi_p$ is the correlation length of percolation and $d$ is the Euclidean dimensionality of the system. Since near $p_c$, $\xi_p \sim |p - p_c|^{-\nu}$, if $\ell$ is the thickness of the links,

then we must have

$$I_f \sim i_w \frac{\ell}{\xi_p^{d-1}} \sim (p - p_c)^{(d-1)\nu}. \tag{24}$$

On the other hand, $I_f = g_e V_f$, where $g_e$ is the effective conductivity of the sample, and $V_f$ is the failure voltage. Since near $p_c$ one has, $g_e \sim (p - p_c)^\mu$, where $\mu$ is the critical exponent of the effective conductivity near $p_c$, we obtain

$$V_f \sim (p - p_c)^{(d-1)\nu - \mu}. \tag{25}$$

Equation (25) can also be derived by the following more detailed analysis. If, for length scales $L \ll \xi_p$, we cut one red bond, it splits the sample-spanning conducting cluster (and hence the backbone) of the material into two pieces, and therefore the total critical current for breakdown is $I \sim O(1)$ (because all the current must go through this red bond), and thus the failure (breakdown) current density (current per length of the sample) is $I_f = I/L \sim 1/L$. Therefore, the failure voltage $V_f$ is given by $V_f \sim I_f/G_e$, where $G_e$ is the effective *conductance* of the sample. As $G_e \sim L^{d-2}L^{-\mu/\nu}$ for a $d$-dimensional system (note that the factor $L^{d-2}$ is included to convert the effective conductivity to the effective conductance), we obtain $V_f \sim L^{\mu/\nu - (d-1)}$. For $L \gg \xi_p$, we replace $L$ by $\xi_p \sim (p - p_c)^{-\nu}$ and obtain Eq. (25). Equation (25) indicates that there is a qualitative difference between 2D and 3D materials. In 2D where $\mu \simeq 1.3$ and $\nu = 4/3$, $(d-1)\nu - \mu > 0$, and therefore $V_f$ vanishes as $p_c$ is approached, in agreement with the results of computer simulations (see below). On the other hand, in 3D where $\mu \simeq 2.0$ and $\nu \simeq 0.88$, $(d-1)\nu - \mu < 0$, and therefore $V_f$ *diverges* as $p_c$ is approached. Therefore, a thin (2D) conducting film (attached to a substrate) suffers electrical breakdown quite differently than a bulk (3D) material.

If, instead of a lattice model, we utilize a continuum one, then, the exponent that characterizes the power-law (25) will be different from its lattice counterpart $(d - 1)\nu - \mu$. For example, this exponent for the Swiss-cheese model in which spherical or circular holes are distributed in an otherwise uniform conducting matrix, is given by $\nu + d - 1 + \delta$, where $\delta = 1$ and $3/2$ for $d = 2$ and $3$, respectively, so that the voltage $V_f$ for a continuum near its percolation threshold is *smaller* than the corresponding value for a discrete system.

Consider now the effect of the sample size which, in the context of the lattice model, leads us to the size of the most damaging defect which is an inclusion of size $\xi_p$, in the direction parallel to the macroscopic voltage, and $\ell_\perp$ perpendicular to it. The total probability $P$ of having a defect of size $\ell$ is, $P = P_\ell (L/\xi_p)^d$, where $P_\ell$ is the probability density of defect clusters of size $\ell$. Percolation theory predicts (Stauffer and Aharony, 1992; Sahimi, 1994a) that

$$P_\ell \sim \xi_p^{-1} \exp\left(-\frac{\ell}{\xi_p}\right), \tag{26}$$

and therefore, $P = \exp(-\ell/\xi_p)(L^d/\xi_p^{d+1})$. The linear size $\ell_\perp$ is that value of $\ell$ for which $P \sim 1$, which then yields

$$\ell_\perp \sim \xi_p \ln L. \tag{27}$$

Since the current that flows through the side link of the defect is proportional to $(\ell_\perp)^{d-1} I$, one obtains

$$I_f \sim \frac{(p - p_c)^{(d-1)\nu}}{(\ln L)^{d-1}}, \qquad (28)$$

implying that, finite size of the sample generates a (weak) correction to Eq. (24). Thermal effects also modify Eq. (24) which will be described in Section 5.2.7.

Duxbury and Li (1990) proposed that one may combine the above results for the dilute limit and the region near $p_c$, Eqs. (22) and (24), into a single unified equation, given by

$$I_f = I_w \frac{\left(\dfrac{p - p_c}{1 - p_c}\right)^\phi}{1 + c \left[\dfrac{\ln(L/\xi_p)}{\ln(1 - p)}\right]^\psi}, \qquad (29)$$

where $c$ is a constant, and $\psi$ is an exponent, the precise value of which is not known, but can be bounded by

$$\frac{1}{2(d - 1)} < \psi < 1. \qquad (30)$$

Thus, in general, there exist three regimes.

(1) For $p = 1$, one has $I_f = I_w$, as expected.
(2) For $p \simeq 1$, the numerator of Eq. (29) is essentially a constant of order unity, and one recovers Eq. (22).
(3) For $p \simeq p_c$, the denominator of Eq. (29) is of the order of 1, and one recovers Eq. (24) with $\phi = (d - 1)\nu$.

### 5.2.4   Computer Simulation

One of the first computer simulations of a discrete model of electrical breakdown problem was carried out by de Arcangelis *et al.* (1985). In their model, a fraction $p$ of the bonds are conducting, while the rest, with fraction $(1 - p)$, are insulating. A voltage is then applied to the lattice. Once the current in one bond reaches the failure value, the failed bond is removed (its conductance is set to zero), a voltage is applied again, and the next bond to fail is looked for. This procedure is repeated until the system fails macroscopically and its effective conductivity vanishes. A slightly more general version of this model was studied by Duxbury *et al.* (1995) in which each bond of a lattice, with probability $p$, is a conductor with conductance $g_1$ and failure current threshold $i_1$, while the rest of the bonds, with a fraction $(1 - p)$, have a conductance $g_2$ and a threshold $i_2$. Söderberg (1987) and Stephens and Sahimi (1987) suggested another model in which each bond burns out and becomes insulating if the dissipated Joule heat in it exceeds a threshold value.

In general, as more bonds fail, the necessary applied voltage for failing a bond decreases. de Arcangelis *et al.* (1985) determined two voltages: One, $V_i$, is the

FIGURE 5.2. Size dependence of the failure current $I_f$ in the square network. The curves from top to bottom are, respectively, for $p = 0.6, 0.7, 0.8$ and $0.9$ (after Duxbury *et al.*, 1987).



voltage at which the first bond fails, and a second one, $V_l$, causes the last bond, and hence the sample, to fail. The two voltages exhibit very different behaviors as $p$, the fraction of the conducting bonds in the original lattice, was varied. $V_i$ first decreases up to $p \simeq 0.7$, and then increases again. On the other hand, $V_l$ increases monotonically with $p$ until, in the vicinity of $p_c$, it becomes roughly equal to $V_i$. Duxbury *et al.* (1987) employed the same model and analyzed the dependence of the failure current $I_f$ on the sample size $L$. Figure 5.2 presents the results, where $L/I_f$ is plotted versus $\ln L$. The linear dependence of $L/I_f$ on $\ln L$, for several values of $p$, is in agreement with Eq. (29). In addition, when $I_f$ was determined as a function of $p$, it was found to follow Eq. (24) [or Eq. (29)], although when their data are fitted to this equation, the exponent $\phi$ is about 1, rather than the theoretical prediction (for $d = 2$), $\phi = \nu = 4/3$.

de Arcangelis and Herrmann (1989) studied a model of electrical breakdown in which each conducting bond is characterized by a voltage threshold, such that if the voltage along the bond exceeds the threshold, the bond breaks down and becomes an insulator. This model can be thought of as the scalar analog of brittle fracture of materials, in which a microscopic portion of a material behaves elastically until the stress or the force that it suffers exceeds a threshold, in which case the material breaks. The thresholds in the model of de Arcangelis and Herrmann (1989) were distributed according to a probability distribution function. Interesting scaling properties, in addition to what we described above, were discovered for the model. For example, the total current $I$ that passes through the network, as the conducting bonds burn out, scales with the linear size $L$ of the network as

$$I \sim L^\zeta h(N_b/L^{D_f}) \tag{31}$$

where $N_b$ is the number of burnt-out bonds, and $h(x)$ is a universal scaling function. Numerical simulations in 2D indicated that $\zeta \simeq 0.85$ and $D_f \simeq 1.7$. Note that $D_f$ represents the fractal dimension of the set of *all* the burnt-out bonds. If one considers only those burnt-out bonds that form a sample-spanning cluster, then one finds that, $D_f \simeq 1.1$, indicating that the cluster is almost like a straight line. Moreover, de Arcangelis and Herrmann (1989) found that the distribution of the local currents in the network just before it fails macroscopically is multifractal, so that each of its moments is characterized by a distinct exponent (which is similar

to the distribution of currents in random resistor networks studied in Chapter 3, and also Chapter 5 of Volume I), whereas the same distribution obeys constant-gap scaling (i.e., there is a constant difference between the exponents so that from one exponents all other exponents are computed) before the catastrophic failure sets in (i.e., the point beyond which the network burns out very quickly and becomes insulating). Since, as pointed out above, many properties of such models of electrical breakdown have analogues in the problem of brittle fracture, we postpone a more detailed discussion of these properties to Chapter 8 where we describe and discuss the discrete models of brittle fracture and other types of mechanical breakdown. Simulation of large 3D models of this type was carried out by Batrouni and Hansen (1998) who found that their results follow Eq. (31).

### 5.2.5 Distribution of the Failure Currents

Equations (22), (24), and (29) predict the value of the *most probable* failure current. In practice, this quantity is *not* a self-averaged property. That is, nominally identical samples have different failure currents. Therefore, there is a distribution of such currents, which also depends on the linear size of the sample. Duxbury *et al.* (1987) determined this distribution by computing $P_L(N)$, the probability that in a sample of linear size $L$, no defect of insulating configuration with a size larger than $N$ (bonds) is formed. In order to accomplish this, the lattice is divided into smaller elementary cubes (or squares in 2D) of linear size $L_c$. Due to the statistical independence of the elementary cubes, the probability that no defect of size larger than $N$ forms is $[P_L(N)]^n$, where $n$ is the number of the elementary cubes of linear size $L_c$. Since the distribution functions must have the same form on the lattice and its elementary cubes or sublattices, one must have

$$[P_L(N)]^n = P_L(a_n N + b_n), \qquad (32)$$

where $a_n$ and $b_n$ are scaling functions that remain finite as $n \to \infty$. Two general solutions can now be derived.

(1) $a_n = 0$, in which case one has

$$P_L(N) = \exp[-x_1 \exp(-x_1 N)], \qquad (33)$$

where $x_1 > 0$ and $x_2 > 0$ are two parameters to be determined.
(2) $b_n = 0$, in which case one obtains

$$P_L(N) = \exp(-r/N^m), \qquad (34)$$

with $r > 0$ and $m > 0$. To determine the constants $x_1$ and $x_2$, we note that the probability that a defect of size $N$ is formed is given by $dP_L/dN$, and the maximum of this probability is obtained when $N = N_c$, where $N_c$ is given by Eq. (20). Consequently, one finds that, $x_1 = cL^d$ and $x_2 \sim -\ln(1-p)$, where $c$ is a constant which depends on the dimensionality $d$ of the system. Thus, the numerical value of $x_1$ is large, ensuring that $P_L(0) \simeq 0$. Combining these results with Eqs. (18) and (22), the cumulative probability of failure,

$F_L(I_f) = 1 - P_L(N)$ is then given by

$$F_L(I_f) = 1 - \exp\left\{-cL^d \exp\left[-dc\left(\frac{I_w/I - 1}{I_w/I_f - 1}\right)^{d-1}\ln L\right]\right\}. \quad (35)$$

Distribution (35), which was derived by Chakrabarti and Benguigui (1997), is a double exponential distribution and is normally referred to as the *Gumbel* distribution (Gumble, 1958). We remind the reader that $I_w$ is the current for the failure of the pure sample (without any insulating region, or the limit $p = 1$). If $L$ is large enough, then the current $I_f$ that appears in Eq. (35) is indeed the most probable failure current. Note also that $F_L(\infty) \to 1$ only when $L$ is large enough. How large is large enough cannot be answered very precisely, because the constant $c$ depends on the dimensionality $d$. Duxbury *et al.* (1987) derived the probability $F_L$ in terms of the failure voltage $V_f$. Their equation is given by

$$F_L(V_f) = 1 - \exp\left[-cL^d \exp\left(-\frac{kL^{d-1}}{V_f}\right)\right], \quad (36)$$

where $k$ is a constant.

On the other hand, Eq. (34) does not lead straightforwardly to a corresponding cumulative probability of failure. However, it is often stated that, the cumulative probability distribution $F_L$ that corresponds to Eq. (34) is the Weibull distribution, given by

$$F_L(I_f) = 1 - \exp\left[-rL^d\left(\frac{I}{I_f}\right)^m\right]. \quad (37)$$

If the parameter $m$ is large enough, then $I_f$ that appears in Eq. (37) is indeed the failure current. Two points are now worth mentioning.

(1) Distributions (35)–(37) are valid if the material is far from its percolation threshold. It has been proposed that, near $p_c$, the following cumulative failure distribution should be valid,

$$F(I_f) = 1 - \exp\left\{-c'L^d \exp\left[-\frac{k'(p - p_c)^\nu}{I_f^{1/(d-1)}}\right]\right\}, \quad (38)$$

where $c'$ and $k'$ are two constants, the precise values of which are not known. The distribution (38) is similar to (35) (in the sense of being double exponential) although, unlike (35), it has never been checked against the results of computer simulations or experimental data.

(2) It is difficult to test the validity of the Gumbel distribution against the Weibull distribution by simply fitting the data to them. However, if one defines a quantity $A$ by

$$A = -\ln\left\{-\frac{\ln[1 - F_L(V_f)]}{L^d}\right\}, \quad (39)$$

then, the corresponding quantity, for example, for the distribution (37) (when written for the failure voltage $V_f$) is given by

$$A_W = a_1 \ln \left( \frac{1}{V_f} \right) + b_1, \tag{40}$$

and thus a plot of $A_W$ versus $\ln(1/V_f)$ must be linear. On the other hand, for the Gumbel distribution, Eq. (35) or (36) [or (38)], one has

$$A_G = a_2 \left( \frac{1}{V_f} \right) + b_2, \tag{41}$$

which predicts linear variation of $A_G$ with $1/V_f$. In this way, one can clearly determine which cumulative distribution provides a better fit of the data. Duxbury *et al.* (1987) found, using this method, that the Gumbel distribution provides a more accurate fit of their numerical data. In Chapter 8 we will utilize this method in order to test the accuracy of analogous distributions for the failure stress of brittle materials.

## 5.2.6   *The Effect of Failure Thresholds*

In practice, different parts of a material may exhibit different resistance to electrical breakdown. Therefore, a more realistic model may be one in which one characterizes the conducting bonds by a threshold in the voltage or current, beyond which it breaks down and becomes insulating. The thresholds can be selected from a probability density function, which then introduces into the model a heterogeneity that is different from percolation disorder. Kahng *et al.* (1988) considered such a model in which each bond is characterized by a failure voltage uniformly distributed over the range $v_- = 1 - \frac{1}{2}w$ to $v_+ = 1 + \frac{1}{2}w$, where $0 < w \le 2$. All the bonds have the same resistance. The limit $w = 0$ represents a system without any disorder, while the limit $w = 2$ corresponds to a uniform distribution in (0,2). A voltage is applied to the system and is increased until the first bond fails. The conductance of the failed bond is set to zero, the applied voltage is kept content, and the voltage distribution in the network with its new configuration is recalculated. If another bond fails, its conductance is set to zero, and the procedure is repeated. If, at some stage, no more bond fails, the applied voltage is increased gradually until the next bond fails. This procedure is repeated until the entire sample fails. This model represents a slow breakdown process, since the characteristic time for a "hot" bond to suffer breakdown and become an insulator is assumed to be much larger than the time that it takes the system to relax and reach equilibrium.

Despite its apparent simplicity, the behavior of the system depends crucially on the value of $w$ and the linear size $L$ of the network, and exhibits interesting phenomena. Similar to the models with percolation disorder, the value of the external voltage to cause the network to fail decreases as $L$ increases, but at a rate that depends on $w$. If $w$ is small enough, then one of the first few bonds that fails triggers a path of failed bonds that propagates across the system. This is somewhat similar to brittle fracture of relatively homogeneous solids (in which

mechanical failure of the first few atomic bonds generates a path of broken bonds that eventually spans the materials), and hence we refer to this case as the "brittle" regime. In this case, the failure of the material is governed by the *weakest* (or at least one of the weakest) bonds in the initial system. For larger $w$, the disorder is stronger, and therefore the breakdown of the material is more gradual, as there is a large range over which individual bonds' failure is driven by an increases in the applied voltage. This situation somewhat resembles ductile fracture, and therefore we refer to it as such (without claiming that it actually represents the scalar analogue of ductile fracture). In a $d$-dimensional network of volume $L^d$, ductility is expected if the number of failed bonds exceeds $L^{d/2}$. Then, the behavior of the breaking voltage in the ductile regime parallels that of materials with percolation disorder. However, the average breaking voltage cannot be less than $v_-$, and therefore this leads to an eventual crossover to the brittle regime as the linear size $L$ of the network increases, except when $w = 2$.

Whether the network behaves as in the brittle or ductile regime depends on $w$ and $L$. Kahng *et al.* (1988) showed that there exists a critical value $w_c^{(1)}$ of $w$ such that, regardless of $L$, the material always fails in the brittle regime. The failure of the system in this case is trivial. For $w > w_c^{(1)}$, the network's failure is brittle for large $L$ and ductile for small $L$. The two regimes are separated by another critical value of $w$, $w_c^{(2)}(L)$, which is a function of $L$. For $L \to \infty$, one has $w_c^{(2)} \to 2$, and failure of the system is brittle.

More quantitatively (but approximately), we consider the sequence of the weakest bonds. The average failure voltage for the $N$th weakest bond to break can be shown to be (Kahng *et al.*, 1988)

$$V_1 = \langle v_{\text{weak}}(N) \rangle = v_- + \frac{wN}{L^2}, \qquad (42)$$

which predicts a linear dependence of $V_1$ on $N$, since the distribution of the thresholds is uniform and must be equal to $v_-$ for $N = 0$ and to $v_+$ for $N = L^2$. We now suppose that $N$ bonds have failed and formed $2N$ edge bonds, where there is an increase of the current due to enhancement effect (see above). It can then be shown that the average failure voltage for the $2N$ failed bonds is given by

$$V_2 = \langle v_{\text{edge}}(N) \rangle = v_+ + \frac{w}{2N + 1}. \qquad (43)$$

Observe that $V_2$ is a decreasing function of $N$, because as $N$ increases, the probability that a weak bond is included in the $2N$ edge bonds increases. An approximate criterion for brittleness of the system is then given by

$$\mathcal{E}V_1(N) > V_2(N), \qquad (44)$$

where $\mathcal{E}$ is the enhancement factor described earlier. Then, two possible situations may arise:

(1) If we plot $\mathcal{E}V_1(N)$ and $V_2(N)$ versus $N$, the two curves do not cross each other. In this case, the network becomes unstable (behaves as in the brittle regime) after the first bond fails, regardless of the network size $L$. For this to happen,

one must have $\mathcal{E} v_- > v_- + w = 1 + \frac{1}{2} w$, and therefore,

$$w_c^{(1)} = 2 \frac{\mathcal{E} - 1}{\mathcal{E} + 1}. \tag{45}$$

For example, as mentioned above, for the square network, $\mathcal{E} = 4/\pi$, and there-fore $w_c^{(1)} \simeq 0.24$. For $w < w_c^{(1)}$ the effect of the randomness is trivial, since the minimum voltage to break the first bond is just $v_- = 1 - \frac{1}{2} w$, which generates a voltage $\mathcal{E} v_-$ at its edge.

(2) In the second case, the curves $\mathcal{E} V_1(N)$ and $V_2(N)$ do cross each other. The crossing point defines the critical value $N_c$ of $N$ for failure of the system. During breakdown of the first $N_c$ bonds, the system is stable and behaves in the ductile regime, but it becomes unstable beyond $N_c$ and fails. However, if $L$ is small enough, then the system may stay in the ductile regime.

The mean failure voltage (per bond) $V_f$ was also determined by Kahng *et al.* (1988). In the brittle regime, one has

$$V_f = v_- + \frac{\alpha w}{L^2}, \tag{46}$$

where $\alpha$ is a constant. Since $V_f$ can never be less than $v_-$, Eq. (46) indicates clearly that by increasing $L$ the system will always eventually behave as brittle. For the ductile regime, we have

$$V_f \sim (\ln L)^{-y}, \tag{47}$$

where $y \simeq 0.8$ for 2D systems. Equation (47) was confirmed by the numerical simulations of Leath and Duxbury (1994).

Two points are worth mentioning here. One is that the qualitative features of the above results hold for a large class of voltage thresholds (see, for example, de Arcangelis and Herrmann, 1989). However, Stephens and Sahimi (1987) (see also Chan *et al.*, 1991) showed that, if the conductances of the bonds are distributed according to a probability density function, and if this function is of power-law type, then many of the above results do not hold, and the problem is more complex. The second point is that these qualitative features are also observed in discrete models of mechanical fracture, and in fact, prior to Kahng *et al.* (1988), had been predicted by Sahimi and Goddard (1986), who were the first to propose a class of discrete models for mechanical fracture.

## 5.2.7   Dynamical and Thermal Aspects of Electrical Breakdown

All the breakdown models discussed so far are quasi-static models, since they do not have an explicit time scale built in them. However, time-dependent effects in breakdown phenomena are very important. In particular, a highly important characteristic of a conducting material is its *failure time*, i.e., the time that it takes to suffer breakdown and become insulating. Similar to failure current and failure voltage, failure time is also not a self-averaged property of a material, as nominally identical samples exhibit completely different failure times. In practice, what is

usually done is to select *a priori* a distribution, such as log-normal or the Weibull distribution, and fit the experimental data for the failure time in order to estimate the distribution's parameters (see, for example, Ohring, 1998). However, failure time data measured in given test conditions are often sufficiently well fitted by several distributions, with the drawback that different distributions may predict widely different failure times when they are extrapolated to a specific application, hence resulting in serious error. In addition, due to cost limitations, the number of samples tested usually represents only a small fraction of the entire ensemble, and therefore there may be significant uncertainties in the estimated values of the distributions' parameters. Therefore, a dynamic model that can provide accurate predictions for failure times and other dynamical properties is of considerable interest. Another important dynamical aspect of the problem that has not been discussed so far is the behavior of the material in an AC field, whereas use of an AC field in experiments is very common. The question is, how does a material suffer electrical breakdown if one applies an AC voltage across it?

In addition, all the results presented so far have been derived based on purely geometrical considerations, whereas thermal (Joule) effects are in fact the main driving force for electrical breakdown of composite materials. The purpose of this section is to address these issues.

### 5.2.7.1   Discrete Dynamical Models

A few dynamical models have already been developed. We describe and discuss three of these models, one of which is deterministic, while the other two are stochastic. The deterministic model is due to Sornette and Vanneste (1992) and Vanneste and Sornette (1992) (see also Sornette and Vanneste, 1994), which is a generalization of the fuse model of de Arcangelis *et al.* (1985b) described above. In their model, the temperature $T$ of each conducting bond at time $t$ satisfies the following equation

$$C_p \frac{dT}{dt} = Ri^b - aT, \qquad (48)$$

where $C_p$ is the specific heat of the material at constant pressure, $R$ is its resistance, $i$ is the current in the bond, and $a$ and $b$ are two constants. The $Ri^b$ term accounts for a generalized Joule heating of the bond ($b = 2$ for real fuses), while $aT$ represents the heat lost to the substrate. To each conducting bond a critical temperature $T_c$ is assigned, such that the bond burns out and becomes an insulator once its temperature exceeds $T_c$. A current $I$ is injected into the system, and the current distribution throughout the network is calculated. Each bond's current is then used in Eq. (48) to calculate the time evolution of its temperature. The first bond burns out when its temperature reaches $T_c$. The current distribution in the new network is calculated and the next bond is allowed to burn out. Thus, one essential assumption of the model is that, the redistribution of the currents in the network is either instantaneous or happens much faster than the temperature evolution of the bonds. The limit $b \to \infty$ corresponds to the fuse model described earlier, because in this limit only the bond that carries the largest current is heated significantly

and reaches its critical temperature faster than any other bond. The opposite limit, $b \to 0$, corresponds to a percolation model, because in this limit the heating rate becomes independent of the current, and therefore the sequence in which the bonds burn out is essentially random. Note that there are two characteristic time scales in the system which are $t_1 = T_c / R i^b$ and $t_2 = 1/a$. If $I_c$ is the critical current for the emergence of the first sample-spanning cluster of the burnt-out bonds, then, three distinct regimes can be recognized.

(1)  If the current $I$ through the network is very close to $I_c$, then one has a number of growing clusters of burnt-out bonds, all nucleating from the same center, which is the first burnt-out bond in the network. The degree of branchiness of the clusters depends on the quenched disorder of the network (for example, the distribution of the resistances). The larger the disorder in the network, the more branched the clusters are.
(2)  If $I \gg I_c$, then there is only one relevant time scale, $t_1$, in the system. Initially, the bonds burn out more or less at random, a process that is dominated by the quenched disorder, and then at later times the growth of the burnt-out clusters becomes correlated as they become connected.
(3)  The third regime corresponds to a crossover between (1) and (2). In this case, the behavior of the system is extremely sensitive to the applied voltage or current. The model produces a hierarchy of evolving failure patterns at various length scales, as the applied current $I$ is varied. The breakdown patterns are also fractal with a fractal dimension $D_f$ which is a strong function of the parameter $b$. Experimental realization and confirmation of this model will be described and discussed shortly.

A stochastic model that takes into account the Joule effect was developed by Pennetta *et al.* (2000), which was intended for electrical breakdown of thin conducting films. An external current $I$, which is held constant, is injected into a 2D lattice. Each bond of the network is a resistor with a resistance $r(T) = r_0[1 + \alpha(T - T_0)]$, where $r_0$ is a constant resistance, $T$ is the resistor's present temperature, $T_0$ is a constant reference temperature, and $\alpha = (1/r)dr/dT$ is the temperature coefficient of resistance. A bond breaks down and becomes an insulator with a probability $p_b$ given by

$$p_b = \exp\left(-\frac{\mathcal{H}_0}{k_B T}\right), \tag{49}$$

where $\mathcal{H}_0$ is an activation energy characteristic, and $k_B$ is the Boltzmann's constant. The temperature in the $j$th resistor is updated according to the following equation

$$T_j = T_0 + a_1\left[r_j i_j^2 + \frac{a_2}{N}\sum_{k=1}^{N}(r_k i_k^2 - r_j i_j^2)\right], \tag{50}$$

where $i_j$ is the current in, and $N$ is the number of nearest neighbors of, the $j$th resistor. The parameter $a_1$ describes the heat coupling of each resistor with the

substrate to which the thin film is attached, and measures the importance of Joule heating effects. $a_2$ is a constant which was taken to be $3/4$.

Hence, starting from a resistor network in which all the bonds are conducting, the current and temperature distributions in the network are calculated. Conducting bonds are then converted to insulating ones with a probability given by Eq. (49). The current and temperature distributions are then recalculated, the next bonds to fail are identified, and so on. The simulations stop when a sample-spanning cluster of the failed bonds is formed. Computer simulations indicated that the effective resistance $R_e(t)$ of the sample at time $t$ follows the following power law,

$$R_e(t) \sim (t - t_f)^{-\mu_d}, \tag{51}$$

where $\mu_d \simeq 1/4$. Note that the failure time $t_f$ can be estimated from *two* measurements of $R_e(t)$ at two different times, namely,

$$t_f = \frac{ct_1 - t_2}{c - 1}, \tag{52}$$

where $c = [R_e(t_1)/R_e(t_2)]^{1/\mu_d}$ represents the ratio of the two measured resistances at two different times, raised to the power $1/\mu_d$. Therefore, once again, the concepts of scaling and universality seem to be quite useful to modeling of an important phenomenon, namely, electrical breakdown of thin solid films. Let us mention that another deterministic model that takes into account the Joule effect, but uses nonlinear, power-law, resistors (see Chapters 2 and 3) was developed by Martin and Heaney (2000).

The second stochastic dynamical model that we describe was developed by Hansen *et al.* (1990), and is a generalization of the dielectric breakdown model of Niemeyer *et al.* (1984) which will be studied shortly, but also has some similarities to the fuse model of de Arcangelis *et al.* (1985) described above. In their model, a conducting bond breaks down and becomes an insulator with a probability $p_b \sim i_{ij}^\eta$, where $\eta$ is a parameter of the model, and $i_{ij}$ is the current in the bond $ij$. Initially, all the bonds in the network are conducting. A macroscopic voltage drop is applied to the network, and the current distribution in the bonds is computed. The bond that breaks first is selected from among *all* the conducting bonds. The current distribution in the network with its new configuration, including the failed bond, is calculated, the next bond to be broken is selected, and so on.

This model provides some interesting predictions. Hansen *et al.* (1990) found that there is a critical value $\eta_c = 2$ of $\eta$, such that the breakdown patterns are qualitatively different for $\eta < \eta_c$ and $\eta > \eta_c$. For $\eta < \eta_c$ the breakdown pattern resembles a percolation cluster, in the sense that a finite fraction of the conducting bonds must breakdown before the system fails and becomes insulating. On the other hand, for $\eta > \eta_c$ the breakdown pattern is a fractal object with a fractal dimension that depends on $\eta$. The vector analogue of Hansen *et al.*'s model, i.e., one in which the bonds represent elastic elements that break with some probability (which might be applicable to mechanical fracture), was analyzed in detail by Curtin and Scher (1991,1992).

## 5.2.7.2    Breakdown in an AC Field: Thermal Effects

Suppose that the initial resistance of a sample material is $R_0$. If a current $I$ is injected into the material, its resistances will change by $\Delta R = R_0 \alpha \Delta T$, where $\alpha$ is the temperature coefficient of the resistance, and $\Delta T$ is the temperature rise in the sample as a result of injecting the current into the material. Since $\Delta T \sim R_0 I^2$, one obtains $\Delta R \sim (R_0 I)^2$. Corresponding to the current $I$ there exists a voltage $V$ across the material which is given by

$$V = R_0 I + c T_0^2 I^3,    \tag{53}$$

where $c$ is a constant. Then, if $I = I_0 \cos(\omega t) = i_0 \cos(2\pi f t)$, the voltage $V$ becomes

$$V = R_0 I_0 \cos(\omega t) + V_{3f} \cos(3\omega t),    \tag{54}$$

where $V_{3f} \sim \Delta R I_0$ is the *third harmonic voltage*. The *third harmonic coefficient* (THC) $B$ is then defined by

$$B = \frac{V_{3f}}{I_0^3}.    \tag{55}$$

As discussed by Dubson *et al.* (1989), the THC results from local Joule heating. Therefore, in effect $B$ measures the local temperature rise at the hot spots that are developed as a result of Joule heating.

If the material is a two-phase composite a fraction $p > p_c$ of which is conducting and the rest is insulating, then, as was pointed out by Yagil *et al.* (1992,1993), the failure current $I_f$ is related to the THC $B$. Yagil *et al.* (1992) suggested that breakdown occurs when a hot spot in the material reaches the melting temperature $T_m$ of the metallic (conducting) grains, at which a weak link in the system breaks down, an irreversible change occurs in the material, and its resistance is modified. To derive the relation between $I_f$ and $B$ (Yagil *et al.*, 1992), one notes that the temperature rise due to a weak link with resistance $r_0$ and current $i$ is $\Delta T = r_0 i^2 \mathcal{R}$, where $\mathcal{R}$ is the ratio of the temperature rise and the dissipated power at the hot spot. The resulting change in the local resistance is $\delta r = r_0 \alpha \Delta T$, where $\alpha$ is the temperature coefficient of resistance. If one applies an AC current, $I = I_0 \cos(\omega t)$, to the material, it results in the generation of a third harmonic voltage component $V_{3f}$, given by

$$V_{3f} = \frac{1}{4 I_0} \sum_j i_j^2 \delta r_j,    \tag{56}$$

where the sum is over all the hot spots. If we assume that the resistance $r$ and the ratio $\mathcal{R}$ are the same for all the links (in the percolation material), we obtain

$$B = \frac{\alpha r^2 \mathcal{R}}{4 I_0^4} \sum_j i_j^4    \tag{57}$$

which implies that $B$ is related to the fourth moment of the current distribution in the material, a subject that was discussed in Section 5.16 of Volume I. The current in

each resistor of an $L \times L$ resistor network is $I/L$ (where $L$ is measured in units of the bonds). For a resistor network near the percolation threshold, the current in the red bonds (i.e., those that, if cut, would break the backbone into two pieces) is much larger than the rest of the bonds. Since near $p_c$ the resistance follows the power law $(p - p_c)^{-\mu}$, and the fourth moment of the current as $(p - p_c)^{-2\kappa}$ (see Chapter 5 of Volume I), the third harmonic follows the power law $(p - p_c)^{-(2\mu+\kappa)}$, where, as discussed in Chapter 5 of Volume I, the exponent $\kappa$ is independent of all the percolation exponents.

As discussed in Chapter 5 of Volume I, one may obtain upper and lower bounds for the exponent $\kappa$. Similar ideas can be used for deriving a bound for $B$ (Yagil *et al.*, 1992). Consider, for example, deriving a lower bound to $B$ for a thin (2D) film. The bound is obtained by taking into account only the red bonds. For $L \gg \xi_p$, where $\xi_p$ is the correlation length of percolation, the number $M_r$ of the red bonds follows the power law, $M_r \sim (p - p_c)^{-1}$, and the current through each of such bounds is $I_r = (\xi_p/L)I_0$. Therefore,

$$B \geq \frac{1}{4I_0^4}\alpha r_0^2 \mathcal{R} M_r I_r^4. \tag{58}$$

On the other hand, the average AC component of the temperature increase in each of the red bonds is, $\Delta T_r = \frac{1}{2}r_0 \mathcal{R} I_0^2$, and thus

$$B \geq \frac{1}{2I_0^2}\alpha r_0 \Delta T_r (p - p_c)^{-1}. \tag{59}$$

Suppose now that $\Delta T_m$ is the temperature rise that the material needs to reach the melting temperature of its conducting portion. If one defines the failure or breakdown current $I_f$ as the current at which the melting temperature is reached, then

$$I_f \geq \sqrt{\frac{1}{2}\alpha r_0 \Delta T_m}(p - p_c)^{-1/2}B^{-1/2}. \tag{60}$$

The THC for the pure material (with no insulating region) is given by, $B_0 = \frac{1}{4}\alpha r_0^2 \mathcal{R}/L^2$, while its failure current is, $I_f^0 = L(2\Delta T_m/r_0\mathcal{R})^{1/2}$. Since, $(p - p_c)^{1/2} = (B/B_0)^{1/2(2\mu+\kappa)}$, one obtains the final result:

$$I_f \geq \frac{I_f^0}{B_0^{1/2(2\mu+\kappa)-1/2}}B^{1/2(2\mu+\kappa)-1/2}. \tag{61}$$

If we substitute the 2D lower bound, $\kappa = 2\nu + 1 - 2\mu$ (see Chapter 5 of Volume I), we obtain $I_f \sim (p - p_c)^{\nu}$, in agreement with Eqs. (24) and (28). Thus, taking the thermal effects into account, one obtains a refinement to Eqs. (24) and (28) which were derived earlier based on geometrical considerations alone. Since, typically, only a fraction of the red bonds contribute significantly to the sum $\sum i^4$, we expect to have

$$I_f \sim B^{-w}, \quad \text{with} \quad \frac{1}{2} - \frac{1}{2(2\mu + \kappa)} \leq w \leq \frac{1}{2}. \tag{62}$$

We are now ready to compare the above theoretical predictions to the relevant experimental data.

### 5.2.7.3    Comparison with the Experimental Data

An experimental realization of the dynamical model of Sornette and Vanneste (1992) was provided by Lamaignere *et al.* (1996). In their experiment, insulating epoxy resin was mixed with spherical carbon microbeads. The matrix was obtained by heating the solution for 2 hours, yielding a conducting composite with quenched disorder. The $I - V$ characteristic of the composite is linear when the applied voltage is small, $V < V_1$, signifying the fact that the connectivity properties of the composite are independent of the voltage $V$. For $V_1 < V < V_c$, where $V_c$ is the critical threshold, the $I - V$ curve bends over and the tangential conductivity decreases, indicating a significant change in the connectivity of the beads which is the result of local breakdown caused by Joule heating. If the volume fraction of the beads is above the percolation threshold, and if the temperature of the system is in the range $20 - 30°C$ above $120°C$, an additional factor decreases the conductivity of the composite. This factor is due to the thermal expansion of the polymer matrix that entails strain growth, leading to a redistribution of the stress field and modification of the connectivity, and thus the conductivity. Beyond $V_c$ and its corresponding current $I_c$ the tangential conductivity vanishes, and $I$ deceases as $V$ increases. For $I \geq I_c$ macroscopic breakdown occurs. These data are summarized in Figure 5.3. As $I \geq I_c$ increases the resistance of the composite also increases,



FIGURE 5.3. $I - V$ characteristics obtained under applied voltage (circles) or current (squares). Dashed line indicates the linear $I - V$ behavior, while the thick solid line indicates the critical value of the current (after Lamaignere *et al.*, 1996).

signaling the breakdown of more and more conducting fraction of the composite. Suppose that at time $t_f$ the composite fails and becomes insulating. Lamaignere *et al.* (1996) found that

$$t_f \sim I^{-2}, \tag{63}$$

and that the effective conductivity $g_e$ of the composite at times close to $t_f$ follows the power law,

$$g_e \sim (t_f - t)^{\mu_d}, \tag{64}$$

with $\mu_d$, which is a sort of dynamical analogue of the percolation conductivity exponent $\mu$, being about 2/3 for their 2D material.

Yagil *et al.* (1992,1993) measured failure current $I_f$ and the THC of thin, semi-continuous Ag and Au percolating films. The films were evaporated in vacuum at a rate of 0.1 nm/sec onto room temperature glass substrate. Several samples with different surface coverage (i.e., different fraction of the conducting material) were employed. The samples were then removed from the vacuum and measured at room temperature. The measured $I - V$ characteristic indicated Ohmic behavior at low currents, and nonlinear behavior at high currents, due to Joule heating. The failure current $I_f$ was defined as the current at which the first irreversible change in the resistance was measured. Figure 5.4 presents a sample of their results for the failure current $I_f$ versus the THC $B$, measured for the Ag samples.



FIGURE 5.4. Failure current $I_f$ versus the third harmonic coefficient $B$, indicating the slope $w$ ($I_f \sim B^{-w}$). The inset presents the data for the relation, $B \sim R^{2+w}$ with a slope of 3.2 (after Yagil *et al.*, 1992).

The 2D value of the exponent $w$ defined by (62) is bounded in $0.36 \leq w \leq 0.5$, if we use $\mu \simeq 1.3$ and $\kappa \simeq 1.12$. Furthermore, if $\Delta T_m = 10^3\,\mathrm{K}$, $r_0 = 1\,\Omega$, and $\alpha = 10^{-3}\,\mathrm{K}^{-1}$, one obtains, $\frac{1}{2}\alpha r_0 \Delta T_m \simeq 0.5$, and $I_f = I_f^0 (B/B_0)^{-w}$ with $I_f^0 B_0^w \simeq 0.7$, which is in good agreement with the measured value for both the Ag samples, $w = 0.48 \pm 0.05$ and $I_f^0 B_0^w \simeq 0.4$, and for the Au materials, $w = 0.41 \pm 0.01$ and $I_f^0 B_0^w \simeq 0.6$.

The experiments of Yagil *et al.* (1992,1993) shed light on the mechanism of electrical breakdown of composite materials. If the initial material has a low resistance, the breakdown usually results in an insulating composite, implying that *all* the links or red bonds that carry high currents burn out and become insulators. Applying a high voltage, on the other hand, protected by a very low current limit, causes the material to become reconnected again and produce a composite with a high resistance and a very low failure current, indicating that only a few of the red bonds were re-established. Thus, such a material is dominated by the red bonds. On the other hand, according to Yagil *et al.*'s experiments, the breakdown of a high resistance material may result in higher, lower, or infinite (insulating) resistance, implying that a few red bonds are either burnt out, established (dielectric breakdown), or improved (i.e., their width increases).

## 5.3   Electromigration Phenomena and the Minimum Gap

A dynamical model of electromigration was proposed by Bradley and Wu (1993) and Wu and Bradley (1994) which was intended for electromigration failure in polycrystalline metal films. In their model each bond of a lattice is either a conducting wire with probability $1 - p$ or an insulator with probability $p$. Suppose that a certain mass $m_w$ leaves a wire before it fails. The mass flux $j_m$ in the wire is given by

$$j_m = \frac{\rho D}{k_B T} Z^* e E \tag{65}$$

where $\rho$ is the atomic density, $D$ is the diffusivity, $Z^*e$ is the effective charge, $k_B$ is the Boltzmann's constant, and $E$ is the electric field. The total mass $m_w$ out of the wire is proportional to the magnitude of the current $I$, $m_w = a(T)I$, where $a(T)$ is a temperature-dependent constant. Since Joule heat in the wire is rapidly conducted away by the substrate, one can assume that the temperature of the wire and that of the substrate are equal. The lifetime $t_\ell$ of the wire is given by

$$\int_0^{t_\ell} |I(t)|dt = \frac{m_w}{a(T)} = q(T). \tag{66}$$

Therefore, once a charge $q(T)$ has flowed through the wire, it fails irreversibly and becomes an insulator. To a good degree of approximation, the charge $q(T)$ has an Arrhenius-type temperature dependence. Thus, the essentials of this model are as follows. A macroscopic voltage is applied to the network and the current

distribution in it is calculated. The wire that carries the most current fails first, after which the current distribution is calculated again, the next wire to fail is identified, and so on.

Several interesting results emerge from this model. For example, suppose that at time $t = 0$ a "crack" (i.e., insulating material) of length $2c$ is inserted in the metal film and its growth is monitored. Suppose also that $v_\infty(x, c)$ is the speed of the crack tip when the crack's length is $2x$. If a constant external current flows through the film, then for $x \gg c$ (Wu and Bradley, 1994)

$$v_\infty(x, c) \simeq \frac{i_0}{c} x^2,\tag{67}$$

where $i_0$ is the current density far from the crack. Thus, as the crack grows, its speed of propagation increases quadratically. The dependence on the time $t$ of the crack tip location for $x \gg c$ is obviously found from $t = \int_c^x dz/v_\infty(z, c)$. Near $p_c$ the mean failure time $t_f$ obeys the following power law

$$t_f(p) \sim (p_c - p)^\nu,\tag{68}$$

where $\nu$ is the exponent of percolation correlation length; clearly, $t_f = 0$ for $p \geq p_c$.

Electromigration motivates the introduction of a new percolation quantity, which is called the *minimum gap*. Consider a random resistor network in which a fraction $(1 - p)$ of the bonds are insulating. Suppose now that a random walker starts its walk from one side of the lattice, and jumps from one cell to an adjacent cell by crossing the bonds, regardless of whether these bonds are conducting or insulating. We also assume that the walk is self-avoiding, i.e., the walker never visits a cell more than once. After some steps, the walker finally arrives at the opposite face of the network; its path consists of all the bonds that were visited. Suppose then that the path consists of $N_c$ conducting and $N_i$ insulating bonds. The connection between this concept and electromigration becomes clear if we assume that, any bond that is crossed by the walker breaks down and becomes an insulator. Thus, in a 2D system, for example, when the walker has crossed the sample, the system breaks down and becomes an insulator. The shortest path is one that corresponds to the smallest number of resistors that burn out during the walk.

We now introduce the concept of minimum gap $g_m$ which, in an insulating material, is the minimum number of conducting bonds (per length of the system) that must be added to the system (or to the trail of the random walk) in order for the material to become conducting. Clearly, $g_m$ depends on $p$, the fraction of the conducting bonds already in the material. Chayes *et al.* (1986) and Stinchcombe *et al.* (1986) studied the properties of the minimum gap $g_m(p)$. Figure 5.5 present the dependence of $g_m(p)$ on $p$ in the square network. For $p \simeq 1$, the minimum gap decreases from 1, with the slope $dg_m/dp \simeq 3$ in the square network. Near $p_c$, the minimum gap vanishes according to the power law,

$$g_m \sim (p - p_c)^\nu.\tag{69}$$

Thus, Eqs. (68) and (69) suggest that the failure time is proportional to the mini-

FIGURE 5.5. Dependence of the minimum gap $g_m$, normalized by the linear size of the square lattice, on the fraction $p$ of the conducting bonds (after Manna and Chakrabarti, 1987).

mum gap $g_m(p)$ of the network. On an intuitive ground, the relation between the minimum gap and the time to failure in the electromigration problem is expected.

A problem related to electromigration phenomenon is one in which the line width of the metallic interconnects is comparable to, or smaller than, the grain size of the film. In this case, referred to as the *bamboo regime*, the grain boundaries no longer provide connected diffusion paths along the conductor line. Instead, electrical breakdown occurs due to intergranular voids which nucleate at the edges of the line, migrate in the current direction, and finally collapse into a slit which disconnects the conductor.

This problem was studied in detail by Schimschak and Krug (1998), and later by Mahadevan *et al.* (1999), whose analysis we briefly describe. The shape of the void changes due to the current $I$ along its inner surface. Two factors contribute to the current, the electromigration and capillary smoothing. Thus, one writes

$$I = \gamma \left[ \sigma \frac{\partial \mathcal{Y}(\mathcal{L})}{\partial \mathcal{L}} + q E(\mathcal{L}) \right], \tag{70}$$

where $\gamma$ and $\sigma$ are, respectively, the atom mobility and the surface tension, $\mathcal{L}$ is the arc length along the surface, $\mathcal{Y}$ is the surface curvature, $q$ is the charge, and $E$ is the tangential local electric field. Because of conservation of the void area (in

2D), the inner surface must move with a normal velocity $v_n$ which is given by

$$v_n + \frac{\partial I}{\partial \mathcal{L}} = 0. \tag{71}$$

Due to the growth of the void, this is a moving boundary-value problem, the numerical solution of which is typically difficult to obtain.

One must first determine the electric field $E$ by solving the Laplace's equation in the domain outside the void, subject to the boundary conditions that the normal electric field vanishes at the void surface, and a constant electric force $E_0$ is applied to the system far from the void. It is not difficult to see that the only relevant length scale in the problem is $\ell_s = \sqrt{\sigma/(q E_0)}$, and therefore the natural time scale is given by, $t_s = \ell_s^4/(\sigma \gamma)$, with which the governing equations can be made dimensionless. After determining the distribution of the electric field, Eq. (71) is iterated. A breakup procedure is triggered if two points that belong to different surface segments are closer than half the distance between neighboring points along the surface. In a similar way, merging of two voids can be treated.

Numerical simulations of this model indicated that, typically, the void disintegrates at long times by one of the two routes. If the void is initially elongated along the current direction, then, a protrusion develops at the leading end of the void, which subsequently forms a daughter void. Because the daughter void is smaller than the initial void, it moves more rapidly ahead of the mother void. If, on the other hand, the void is initially deformed perpendicular to the current, an invagination develops which eventually splits the void horizontally.

## 5.4   Dielectric Breakdown

We consider a heterogeneous material, consisting mostly of an insulating (dielectric) phase, in which a conducting material has been dispersed. The (volume) fraction of the conducting phase is $p < p_c$, so that, macroscopically, the material is insulating. The electric field $\mathbf{E}$ and its corresponding displacement field $\mathbf{D} = \epsilon(\mathbf{r})\mathbf{E}(\mathbf{r})$ satisfy the usual equations that we have used so far in this book:

$$\nabla \cdot \mathbf{D} = 0, \quad \nabla \times \mathbf{E} = \mathbf{0}, \tag{72}$$

where, as usual, $\epsilon(\mathbf{r})$ is the dielectric constant of the insulating phase.

### 5.4.1   Exact Duality Relation

The duality relations described in Chapter 2, and also in Chapters 4 and 5 of Volume I, can also be used here to relate the problem of dielectric breakdown in 2D to the electrical breakdown in 2D (see, for example, Bowman and Stroud, 1989). With $\mathbf{r} = (x, y)$, Eq. (72) implies that

$$\frac{\partial}{\partial x}\left[\epsilon(\mathbf{r})\frac{\partial \phi}{\partial x}\right] + \frac{\partial}{\partial y}\left[\epsilon(\mathbf{r})\frac{\partial \phi}{\partial y}\right] = 0, \tag{73}$$

where the potential $\phi$ is defined such that, $\mathbf{E} = -\nabla \phi$.

Consider now the dual of the 2D material which is obtained by replacing the conducting phase by the insulating material and vice versa. We also assume that the conductivity $g$ of the formerly-insulating parts is given by, $g = 1/\epsilon$. The dual material is conducting since the original material was assumed to be insulating or dielectric, and therefore the current $\mathbf{I}$ must satisfy the continuity equation, $\nabla \cdot \mathbf{I} = 0$, because of which one can write, $\mathbf{I} = \nabla \times \boldsymbol{\psi}$, where the potential vector $\boldsymbol{\psi}$ is selected such that only its $z$-component $\psi_z(x, y) \neq 0$. As $\mathbf{I} = g(\mathbf{r})\mathbf{E}$, we must have

$$\frac{\partial}{\partial x}\left[\frac{1}{g(\mathbf{r})}\frac{\partial \psi_z}{\partial x}\right] + \frac{\partial}{\partial y}\left[\frac{1}{g(\mathbf{r})}\frac{\partial \psi_z}{\partial y}\right] = 0, \tag{74}$$

or,

$$\frac{\partial}{\partial x}\left[\epsilon(\mathbf{r})\frac{\partial \psi_z}{\partial x}\right] + \frac{\partial}{\partial y}\left[\epsilon(\mathbf{r})\frac{\partial \psi_z}{\partial y}\right] = 0. \tag{75}$$

In view of Eq. (73), we see that the conductivity problem in the dual material is identical with the dielectric problem in the original composite, if $\partial \psi_z/\partial x = \partial \phi/\partial x$ and $\partial \psi_z/\partial y = \partial \phi/\partial y$. If so, one has, $I_x = \partial \psi_z/\partial y = E_y$ and $I_y = -\partial \psi_z/\partial x = -E_x$. Therefore, the magnitude of the current density $\mathbf{I}$ in the dual material is equal to that of $\mathbf{E}$ in the original composite, but its direction is rotated by 90° from the dielectric problem. Physically, while in the electrical breakdown problem the current is zero inside an insulating inclusion, in the dielectric breakdown problem the electric field is zero inside a conducting region. Moreover, the regions that experience an enhancement of the current (in the electrical breakdown problem) are perpendicular to those that feel the enhancement of the electric field (in the dielectric breakdown problem). The conclusion is that, in 2D, most of the results that were described above for the electrical breakdown problem can be immediately translated to corresponding predictions for the dielectric breakdown problem. We will discuss this important point shortly, but let us first describe discrete models of dielectric breakdown.

### 5.4.2   Stochastic Models

The main stochastic model of dielectric breakdown was proposed by Niemeyer *et al.* (1984). In their model, the central site of a square lattice was designated as one of the electrodes, while the other electrode was placed on a circle at a large distance from the center. The rules of the model were as follows.

(1) The electric potential distribution in the lattice is obtained by solving the Laplace equation for $V$, $\nabla^2 V = 0$, with the boundary conditions that $V = V_0 = 0$ for all the sites that belong to the dielectric pattern, and $V = V_\infty = 1$ outside the external circle.

(2) At each step one bond suffers dielectric breakdown and is added to the developing dielectric pattern. The failing bond is selected from amongst those that are at the interface between the dielectric pattern and the rest of the system,

with a breakdown probability $p_b$ given by

$$p_b \sim V_{ij}^{\eta}, \tag{76}$$

where $V_{ij} = V_i - V_j$ is the potential or voltage difference between sites $i$ and $j$ of the interface bond $ij$, with $i$ being on the interface and $j$ outside of, but next to, the interface. Since $V_i = 0$, $V_{ij}$ is simply the potential $V_j$ at $j$, and is proportional to the current in the bond $ij$. In this model, $\eta$ is an important parameter, so much so that this model is popularly known as the $\eta$-model.

(3) After a bond suffers breakdown, the potential distribution in the system with its new configuration is recalculated, a new bond is selected for breakdown, and so on.

Niemeyer *et al.* (1984) showed that their model leads to fractal breakdown patterns which, for $\eta = 1$, are similar to diffusion-limited aggregation (DLA) model of Witten and Sander (1981) (for a review of aggregation models see Meakin, 1998), who had already pointed out the similarity between their model and the breakdown patterns. To see the similarity between the two models, let us describe briefly the DLA model.

In the DLA model one starts with an occupied site (the "seed") of a lattice, located either at the center of the lattice or on its edges. Random walkers are released, one at a time, far from the seed particle and are allowed to move randomly on the lattice. If they visit an empty site adjacent to an occupied one, the aggregate of the occupied sites advances by one site and absorbs the last site visited by the walker (in effect one bond is added to the aggregate). The walker is removed, another one is released, and so on. After a large number of particles have joined the aggregate, it takes on a disordered structure with many branches, very similar to the dielectric pattern with $\eta = 1$. To see the analogy between the two models, note that the original seed particle represents the point at which dielectric breakdown starts. Since the particles perform their random walks on the empty sites, the probability $P(\mathbf{r})$ of finding them at a position $\mathbf{r}$ in this region satisfies the Laplace's equation, $\nabla^2 P = 0$, the same as the governing equation for the nodal potentials or voltages in the dielectric breakdown model. Because the walkers never move into the aggregate, the probability of finding them there is zero, $P = 0$, the same as the boundary condition, $V = V_0 = 0$ in the dielectric breakdown model. Finally, the probability with which the aggregate grows is proportional to the flux of particles between the empty region and the aggregate front, i.e., $\nabla P \simeq P_i - P_j$, the same as Eq. (76) in the limit $\eta = 1$.

In Niemeyer *et al.*'s model, the fractal dimension of the dielectric pattern depends on $\eta$. In 2D one has $D_f \simeq 2.0$, 1.9, and 1.7 for $\eta = 0.$, 0.5, and 1.0, respectively. The resulting 2D pattern for $\eta = 1.0$ is very similar to a Lichtenberg figure. Earlier, Sawada *et al.* (1982) had used a similar model, except that they had assigned *a priori* a larger probability for the growth of the tips with respect to side branching. This is, however, not realistic as the discharge pattern depends *non-locally* on the potential distribution throughout the system, which in turn is controlled by the distribution of the heterogeneities in the material.

However, Niemeyer *et al.*'s model does not have an explicit rule for breakdown. A bond with even a small probability $p_b$ can break down, which is not realistic. Moreover, the physical reason for Eq. (76) is not clear. Pietronero and Wiesmann (1988) did attempt to give a theoretical justification for Eq. (76) based on the time required for the establishment of a filamentary projection of the discharge as a sort of a "conducting fluid" in a given region of the local field. While their argument may justify use of Eq. (76), in the limit $\eta = 1$, for dielectric patterns in gases, its generality is not clear, and in addition, whereas the structure of the simulated discharge patterns is highly sensitive to $\eta$ (Barclay *et al.*, 1990; Sánchez *et al.*, 1992), the physical origin or significance of $\eta$ is not clear. Moreover, the breakdown patterns in solid materials are propagating damage structures, not the advancing front of an injected charge "fluid," as in Niemeyer *et al.*'s model. As such, their model is not, in general, suitable for dielectric breakdown in solids.

Wiesmann and Zeller (1986) (see also Noskov *et al.*, 1995) modified the $\eta$-model by incorporating two new features in it. One was that a critical field $V_c$ for the growth of the dielectric pattern was introduced, such that the breakdown probability $p_b$ is non-zero if $V_{ij} \geq V_c$, and $p_b = 0$ otherwise, an assumption that makes the model somewhat similar to the deterministic models discussed in the next section. The second feature was the introduction of an internal field $V_s$ in the structure, such that the potential in it is no longer $V_0$ but $V_0 + s V_s$, where $s$ is the length of the path (measured as the number of sites that it contains) along the structure which connects the point to the central electrode. The structure of the resulting dielectric pattern now depends on $V_c$ and $V_s$. Figure 5.6 shows two of the fractal patterns generated by this model which are somewhat similar to treeing in polymers. However, the accumulation of damage, which is known to be required for electrical tree formation in AC fields, is not allowed in the Wiesmann–Zeller model, and therefore their model is probably more appropriate for nano-



FIGURE 5.6. Dielectric trees with the ground plate and the needle voltage $V = 0$ and the top plate at $V = V_0$. The threshold field for growth is zero for the left pattern, and about the original field at the tip for the right pattern (after Wiesmann and Zaller, 1986).

second impulses. Even then the damage pattern situation is *not* fractal (Knaur and Budenstein, 1980), whereas the Wiesmann–Zeller model predicts it to be fractal. Dissado and Sweeney (1993) argued that fractal tree-like patterns should form only when the fields at the growth tips can fluctuate around their values obtained from the solution of the Laplace's equation. They showed that if one treats the local-field enhancement factor as a white noise generated by the breakdown mechanism itself, the amount of branching in the dielectric pattern would depend only on the range of the fluctuations allowed. Thus, the Wiesmann–Zeller model, though interesting, is not also completely suitable for modeling dielectric breakdown in solids.

### 5.4.3  Deterministic Models

Several, very similar, discrete deterministic models of dielectric breakdown have been proposed over the past decade. These models assume percolation-type disorder, and their essential features are as follows. Each bond of a lattice is either a conductor with probability $p$ or a capacitor (an insulator) with probability $1 - p$. Each capacitor can sustain a fixed voltage drop, say 1 volt, beyond which it breaks down and becomes a conductor. A macroscopic voltage drop is then applied to the lattice, and the voltage distribution throughout the lattice is computed. The capacitor that sustains the largest voltage drop greater than its threshold fails first. The voltage distribution is then recalculated, the next capacitor to fail is identified, and so on. If at any stage the applied voltage drop is not large enough to cause breakdown of any capacitor, it is increased gradually. The simulation stops when a sample-spanning conducting cluster is formed. The breakdown or failure field $E_b$ is defined as the *minimum* external voltage required to cause formation of a sample-spanning cluster of failed capacitors (conductors), divided by the length $L$ of the lattice. One important result of this model is that $E_b \to 0$ as $p \to p_c$. This is of course due to the tortuous nature of the percolation cluster near $p_c$. Another significant prediction of this model is that $E_b$ is smaller for larger lattice, so that very large samples break down easier than the smaller ones (see also below).

Various versions of this basic model (Beale and Duxbury, 1988) have been studied, the first of which was probably suggested by Takayasu (1985). In his model, the resistance of the lattice bonds are distributed randomly. Each bond breaks down if it suffers a voltage greater than a critical threshold voltage $v_c$. If a bond does break down, its resistance $r$ is reduced to $\delta r$, where $\delta$ is a small number. After a bond breaks down, it remains in that state forever. The breakdown pattern was found to be fractal with a fractal dimension $D_f \simeq 1.6$ in 2D. In the model of Family *et al.* (1986), which is essentially a deterministic version of the Niemeyer *et al.*'s, the bonds are insulating and carry a breakdown coefficient $\mathcal{B}$ which is randomly distributed in [0,1]. The voltage distribution throughout the lattice is then computed, with the boundary conditions that $V = 0$ on the conducting discharge and $V = 1$ far from the interface between the conducting and insulating parts. Two versions of the model were investigated. In one model, at each time step an interface bond $ij$ with the largest $\mathcal{B}V_{ij}^{\eta}$ breaks down, whereas in the second model an interface bond breaks down with a probability $\mathcal{B}V_{ij}^{\eta}/p_b^m$, where $p_b^m$

FIGURE 5.7. Initial breakdown field $E_b$ in the square lattice versus the fraction $p$ of the conducting bonds. Squares and circles show the results for $50 \times 51$ and $100 \times 101$ samples (after Bowman and Stroud, 1989).

is the largest value of $\mathcal{B}V_{ij}^{\eta}$ among all the interface bonds. The second model is clearly very similar to the model of Niemeyer *et al.* (1984). Breakdown patterns were found to be fractal again, with a fractal dimension that depended sensitively on $\eta$. In the model of Manna and Chakrabarti (1987), each bond or site of the lattice is either conducting with probability $p$ or insulating (dielectric) with probability $1 - p$. After determining the voltage distribution throughout the lattice, all the insulating bonds or sites break down if the voltage that they suffer is larger than a threshold voltage. Chakrabarti *et al.* (1987) and Barbosa and de Queiroz (1989) studied this model with small-cell position-space renormalization group approach. Bowman and Stroud (1989) studied the same model, except that in their work the insulating bond with the largest voltage difference between its end sites breaks down first. In a somewhat different model, Benguigui (1988) considered the case in which after a bond breaks down it becomes a superconductor. This was achieved by inserting light emitting diodes as the insulators in a host of conductors.

The most critical questions in dielectric breakdown phenomenon, that any reasonable model should be able to address, are as follows.

(1) How does the initial breakdown field $E_b$ (or the corresponding voltage $V_i$) depend on the volume fraction $p$ of the conducting material (bonds) in the initial dielectric material? A typical example is shown in Figure 5.7.

(2) How does the final voltage $V_f$ vary with $p$? For small $p$ one expects the final breakdown voltage $V_b = E_b L$ to be different from the initial breakdown voltage, but as $p$ increases the difference between the two decreases until very near $p_c$ where they are essentially identical. This has an important consequence in that, when these two voltages are equal, the breakdown proceeds by an avalanche (see the discussion above) in that, many bonds break down without any need for further increase in the applied macroscopic voltage drop.

(3) How do the two voltages $V_i$ and $V_b$ depend on the linear size $L$ of the sample? To understand the importance of the sample size, recall that breakdown starts

near the critical defect of the system, which is (roughly speaking) the largest pair of strongly interacting conducting clusters which are oriented parallel to the macroscopic electric field. The breakdown field is of the order of the inverse of the linear size of the defect, and since the largest defect in a large system is larger than the largest defect in a small sample, the breakdown field is smaller in the larger sample.

(4) How does the path length, i.e., the number of bonds in the breakdown path, vary with $p$? An example is shown in Figure 5.8.

(5) Do power laws govern the important properties of dielectric breakdown (such as the breakdown field $E_b$) near $p_c$, and if so, are such laws universal?

We now discuss the scaling laws that govern the dependence on $p$ of various properties of interest near the percolation threshold, and also on the sample size $L$.

### 5.4.3.1   Scaling Properties of Dielectric Breakdown

Before discussing scaling properties of dielectric breakdown, let us emphasize a very important point. *Unlike* percolation and similar types of critical phenomena, some of the scaling properties of electrical and dielectric breakdown phenomena are valid over a wide range of the parameter space, and therefore are very useful from a practical point of view. For example, as already described and discussed for electrical breakdown, one can consider, for a fixed $p$, the scaling properties of breakdown phenomena in terms of the linear dimension $L$ of the system. Not

only are such scaling properties important, but are in fact measured routinely in practical situations, and therefore a scaling theory of breakdown phenomena in terms of the sample size $L$ is a very useful tool for interpreting the experimental data.

The scaling properties of dielectric breakdown phenomena have been studied extensively. Let us first recall that, as discussed in Chapter 6 of Volume I, the static dielectric constant $\epsilon_0$ follows the following power law (Efros and Shklovskii, 1976) near the percolation threshold $p_c$,

$$\epsilon \sim (p_c - p)^{-s}, \tag{77}$$

where $s$ is the critical exponent that characterizes the effective conductivity of conductor-superconductor percolation composites near $p_c$, utilized extensively in Chapters 5, 6 and 9 of Volume I. The root mean square $E_{rms}$ of the electric field is given by $E_{rms} = \langle |\mathbf{E}|^2 \rangle^{1/2} \propto \epsilon^{1/2} |\mathbf{E}_0| \sim (p_c - p)^{-s/2}$, where $\mathbf{E}_0$ is the applied electric field on the external surface of the system. The maximum field $E_m$ in the system is certainly larger than $E_{rms}$. Suppose that $E_m \sim (p_c - p)^{-y}$. Because $E_m > E_{rms}$, we must have $y > s/2$ (Bowman and Stroud, 1989), and

$$E_b \sim (p_c - p)^y. \tag{78}$$

To estimate $y$, Beale and Duxbury (1988) used an argument based on the idea of the critical defect mentioned above. Suppose that the total length of the critical defect (the conducting path), made up of a pair of the largest interacting clusters of conducting material, separated by a small distance, is $\ell$. The electric field between these two clusters is enhanced by a factor of the order of $\ell$ times the applied macroscopic field. Far from $p_c$ the probability of finding a percolation cluster of linear size $\ell$ is given by Eq. (26). The largest cluster in a $d$-dimensional percolation system of volume $L^d$ is of the order of $\ell_m \sim \xi_p \ln L^d$ [see Eq. (27)]. Since $E_b \sim 1/\ell_m$, we obtain (Beale and Duxbury, 1988)

$$E_b \sim \frac{(p_c - p)^\nu}{\ln L}, \tag{79}$$

and therefore $y = \nu$, which is certainly greater than $s/2$. Equation (79) can also be derived based on the argument (Stinchcombe et al., 1986) that $E_b$ should be proportional to the minimum gap $g_m$ which is proportional to $\xi_p^{-1}$. The $\ln L$ term of Eq. (79) can also be derived from the fact that (Li and Duxbury, 1987) the maximum current $I_m$ in a percolation network of linear size $L$ that leads to its failure is given by $I_m \sim (\ln L)^\psi$, where $\psi$ is the same exponent that appears in (29) and (30) (for the problem of the largest currents in a random resistor network see also Machta and Guyer, 1987). Numerical simulations (Manna and Chakrabarti, 1987; Benguigui, 1988; Beale and Duxbury, 1988; Bowman and Stroud, 1989) seem to confirm Eq. (79). Lobb et al. (1987) and Chakrabarti et al. (1988) extended this analysis to the Swiss-cheese model of continuum percolation, in which spherical or circular grains of dielectric are distributed randomly in a conducting matrix,

and showed that

$$y = \nu + \frac{1}{2} \qquad (80)$$

in any dimension. One may also consider the inverted Swiss-Cheese model (see Chapter 2 of Volume I) in which the metallic grains that can freely interpenetrate are randomly distributed in a dielectric matrix. For this case Lobb *et al.* (1987) showed that

$$y = \nu + 1. \qquad (81)$$

Equations (80) and (81) are both different from $y = \nu$ for lattice models, Eq. (79), and indicate that, as far as dielectric breakdown is concerned, a continuum is *weaker* than a discrete system. This is understandable since in a lattice model the conductivity of the bonds is independent of $p$, whereas the state (geometrical configuration) of a continuum depends on $p$.

The next important issue is the size dependence of $E_b$. Beale and Duxbury (1988) proposed that

$$E_b \sim \frac{1}{A(p) + B(p)\ln L}, \qquad (82)$$

where $A(p)$ and $B(p)$ are simple functions. If we compare Eq. (82) to Eq. (79), we infer that $B(p) \sim (p_c - p)^{-\nu}$, and numerical simulations of Beale and Duxbury (1988) in 2D confirmed this expectation; see Figure 5.9.

### 5.4.3.2   Distribution of Breakdown Fields

Similar to electrical breakdown of solids, the breakdown field for dielectric breakdown is *not* a self-averaged property, because different materials with different types of heterogeneity, or even nominally identical materials, have different



FIGURE 5.9. Breakdown field $E_b$ versus the linear size $L$ of the square lattice. The results, from top to bottom, are for $p = 0.4, 0.35, 0.25$, and $0.1$ (after Beale and Duxbury, 1987).

breakdown fields $E_b$. Therefore, there should be a distribution of such fields for given values of $p$ and $L$. In a series of papers, Duxbury and co-workers (Duxbury *et al.*, 1986, 1987; Duxbury and Leath, 1987; Beale and Duxbury, 1988) derived this distribution for the dielectric (and electrical) breakdown. The resulting distribution is very similar to what we derived for the electrical breakdown problem. A summary of their derivation is as follows. Suppose that $P_L(\ell_m)$ is the probability that no defect (conducting region) larger than size $\ell_m$ exists in a $d$-dimensional cubic lattice of volume $L^d$. We divide the cubic network into smaller cubes of linear dimension $L_c$, and assume that the characteristic size of the largest defect is much smaller than $L_c$. Then

$$P_L(\ell_m) \sim [P_{L_c}(\ell_m)]^{(L/L_s)^d}.$$

Solving this equation and using the fact that for $p \ll p_c$ and $L \gg \xi_p$ the cluster size distribution of percolation systems is an exponentially decaying function of $\ell$, we obtain

$$P_L(\ell_m) = \exp\left[-cL^d \exp(-k\ell_m)\right]. \tag{83}$$

Near $p_c$, the cluster size distribution is of power-law type, in which case

$$P_L(\ell_m) = \exp(-cL^d \ell_m^{-m}), \tag{84}$$

which is of the same form as the classical Weibull distribution, and is appropriate for length scales $L \ll \xi_p$, where $m$ is a constant parameter. Since the breakdown field is of the order $1/\ell_m$, the distribution of the breakdown fields is given by

$$F_L(E_b) = 1 - \exp\left[-cL^d \exp\left(-\frac{k}{E_b}\right)\right], \tag{85}$$

a Gumble distribution which is appropriate for length scales $L \gg \xi_p$. In Eq. (85) the parameter $c$ depends only weakly on $p$, and $k \propto \xi_p^{-1}$. If we now define $E_{1/2}$ as that value of $E_b$ for which half of the system fails, we obtain

$$F_L(E_{1/2}) = \frac{1}{2} = 1 - \exp\left[-cL^d \exp\left(-\frac{k}{E_{1/2}}\right)\right], \tag{86}$$

which, when solved for $E_{1/2}$, yields an equation similar to (82) with $A(p) = [\ln c - \ln(\ln 2)]/k \sim \xi_p$ and $B(p) = d/k = d\xi_p$. The equivalent Weibull forms are

$$F_L(E_b) = 1 - \exp(-cL^d E_b^m), \tag{87}$$

$$E_{1/2} \sim L^{-d/m}. \tag{88}$$

To determine which one of the two distributions, Eq. (85) or (87), can fit the experimental data more accurately, we proceed as in the case of electrical breakdown, namely, we compute the quantities $A_W$ and $A_G$, analogous to Eqs. (40) and (41), and fit the data to them. We note that Sornette (1988) argued that in a continuum

system with percolation-type disorder, such as the Swiss-cheese model, Eq. (86) is no longer valid. Instead, one has a simple exponential, Weibull-like distribution.

### 5.4.4   Comparison with the Experimental Data

The above theoretical results have been tested against (at least) two sets of experimental data. Coppard *et al.* (1989) studied the dielectric breakdown of polyethylene plaques that contained a fixed volume fraction of aluminum particles. Each plaque was compression molded to a disc of thickness 0.7 mm with a depressed inner region of diameter 54 mm. The particles had a well-defined range (53–75 $\mu$m), and were distributed randomly within the polyethylene. The breakdown statistics were collected by stressing the metal-loaded plaques under a uniform AC field and ramping the field amplitude at a fixed rate until breakdown took place. Their data confirmed the validity of Eq. (79), and indicated that Eq. (85) is at least as accurate as Eq. (87).

Benguigui (1988) and Benguigui and Ron (1994) carried out experiments using a square lattice of random resistors and light-emitting diodes. The diodes had a very large resistance up to a voltage threshold $V_b$, but their resistance decreased very significantly above $V_b$, converting them to conductors. The transition between the two states was relatively sharp, but beyond $V_b$ the voltage across the diodes remained essentially constant (which is in contrast to the usual insulator-conductor transition in which the voltage after the transition would be almost zero). The advantages of using such diodes are that, (1) their breakdown is reversible, and (2) the breakdown becomes visible as the diodes, after becoming conductors, emit light.

Suppose that the lattice consists of resistors with fraction $p < p_c$ and the diodes with fraction $(1 - p)$. Figure 5.10 presents the dependence of the voltage $V_b$ on $(p_c - p)$. If we fit these data to Eq. (79), we obtain an exponent $y \simeq 1.1 \pm 0.05$, which reasonably close to the theoretical prediction $y = \nu = 4/3$. The difference is presumably due to the small size of the lattice ($L = 20$) used in the study.



FIGURE 5.10. Failure voltage $V_b$ as a function of the fraction of the resistors $p$ in a system with light-emitting diode (after Chakrabarti and Benguigui, 1997).

## Summary

Discrete models of electrical and dielectric breakdown of composite solids have provided very useful insights into the properties of these important phenomena, by demonstrating the significant role that defects of heterogeneities play in them. In particular, they have provided the important prediction that the statistics of these breakdown phenomena depend critically on the volume fraction of the defects or the broadness of the distribution of the heterogeneities. If the volume fraction of the defects is low, then, the probability distribution of the failure fields (voltage or current) is of Gumble type, rather than the classical Weibull distribution. Moreover, the discrete models have enabled us to obtain important predictions for the effect of sample size on breakdown properties of heterogeneous materials.

# 6
# Fracture: Basic Concepts and Experimental Techniques

## 6.0   Introduction

In Chapter 5 we studied electrical and dielectric breakdown of materials—phenomena that are well-known examples of nonlinear scalar transport processes with their nonlinearity manifested by the existence of a threshold in the linear (or possibly nonlinear) constitutive law that describes the relation between the flux and the potential gradient. Beginning with this chapter, we study a nonlinear vector transport process which is of immense significance to materials, and leads to their mechanical failure. This type of failure, which is a result of nucleation and propagation of fractures in materials, varies anywhere from brittle fracture, that represents a nonlinear vector transport process characterized by a threshold in the otherwise linear elasticity equations that govern the elastic behavior of the material, to ductile yielding and flow. Such failure phenomena are some of the most complex sets of phenomena in science and technology. The range of natural and industrial systems in which mechanical fracture occurs is very broad. Under a large stress or strain, a crack opens up in soils which grows with time, leading to complex phenomena such as soil liquefaction and eventually earthquake. Natural or man-made fractures in oil and geothermal reservoirs and aquifers are crucial to the flow of oil, heat and vapor, or groundwater, especially in those reservoirs that have a very small porosity, such as many oil fields in the Middle East. Other rock-like materials, such as concrete and asphaltenes, often develop large fractures, causing considerable damage to highways and buildings. Propagation of cracks in airplane wings and fuselages can cause an airliner to crash. An important, and undesirable, property of many high-temperature superconducting materials is their brittleness and mechanical instability. Polymers, glasses and ceramics often develop microcracks under a large enough stress or strain which can lead to their mechanical failure and eventual fragmentation. Composite materials can develop cracks due to thermal mismatch between their various constituents. Pressurized nuclear reactors can develop cracks in their structure which can create tremendous safety problems. Thus, a comprehensive understanding of fracture nucleation and propagation has tremendous practical implications.

## 6.1    Historical Background

Most of us have been familiar with the phenomenon of fracture of materials since our childhood, since most of us broke something like a glass or a doll when we were very young. Even if we did not break anything during our childhood, at least some of us might have heard a song like the following in a nursery:

<div align="center">

*Humpty Dumpty sat on a wall*

*Humpty Dumpty had a great fall*

*All the king's horses and all the king's men*

*couldn't put Humpty together again.*[1]

</div>

Any child who heard this song in a nursery was in fact introduced to the phenomenon of fracture, without, of course, knowing it. This simple song also points out two important aspects of fracture phenomenon, namely,

(1) a material develops fracture in response to a driving force which, in the case of Humpty, was the collapse of the church tower, and
(2) fracture is irreversible, since not even all the king's horses and men could put Humpty together after it had been broken into pieces!

The story about Humpty Dumpty also points out another important aspect of fracture of materials, namely, that because of its huge practical significance, the development of an understanding of how materials fracture and break has been of great interest for many *centuries*, and goes back at least 500 years to Leonardo da Vinci who studied fracture of iron wires and showed that a long wire breaks more easily than a short one. That is, long wires are, on average, weaker than short wires. Today, this behavior is known as the *size effect* and is a manifestation of the fact that often fracture is initiated by rare flaws in a material. Since a larger piece of a material is more likely to contain a rare defect, it is also more likely to break under an applied force than a smaller piece of the same material.

Marder and Fineberg (1996) presented a delightful discussion of the historical background of the development of solid mechanics that has led us to the present continuum fracture mechanics. According to them, this development goes back to at least Galileo Galilei who was almost 70 years old when he was working on this subject. His life had been nearly ruined by a trial for heresy before the Inquisition, when he retired in 1633 to his villa near Florence to construct the Dialogues Concerning lluo New Sciences. His first science was the study of the forces that hold objects together and the conditions that cause them to fall apart—the dialogue taking place in a shipyard, triggered by observations of craftsmen building the Venetian fleet. His second science concerned local motions—laws

---

[1]According to legends, Humpty Dumpty was a powerful cannon that was mounted on top of St. Mary's at the Wall Church in Colchester, defending the city against siege in the summer of 1648, during the English Civil War (1642–1649). The church tower was hit by the enemy, with its top blown off, hence sending Humpty to the ground.

governing the movement of projectiles. As we now know, these two subjects have fared differently over the centuries. The first subject, now known as the strength of material, is an integral part of the basic education that most engineering students receive, while the second one has become a core subject that physicists learn at the beginning of their education. Although now, as in Galileo's time, shipbuilders need good answers to questions about the strength of materials, the subject has never yielded easily to basic analysis. Galileo identified the main difficulty when he wrote: *One cannot reason from the small to the large, because many mechanical devices succeed on a small scale that cannot exist in great size*. Over 350 years after Galileo wrote these lines science reached the atomic scale and began to answer the questions that he had posed on the origins of strength and the relation between large and small. These wise words of Galileo also pointed out an important aspect of fracture of materials, namely, the fact that this is an inherently multiscale phenomenon, ranging from atomic to macroscopic length scales. While the vast majority of the theoretical and computer simulation studies of fracture have been concerned with only one of these length scales, the past few years have witnessed development of multiscale modeling approaches to fracture propagation in solid materials. We will describe such approaches in Chapter 10.

However, huge accidents in the 1800s and the first half of the twentieth century, that were caused by catastrophic fracture of materials, provided the motivation for intensive study of fracture phenomena. For example, the boiler of the *Soltana*, a steamboat that carried the Union soldiers during the American Civil War, exploded, resulting in the death of over 1,000 soldiers. In 1919, a molasses tank 50 feet high and 90 feet wide burst in Boston, killing 12 people and several horses. The court auditor concluded that, *the only rock to which he could safely cling was the obvious fact that at least one-half of the scientists must be wrong*.

One of the most important cases of material fracture in the twentieth century, that helped to establish the significance of fracture mechanics, occurred during World War II. Wartime demands for ocean freighters led to the production of the Liberty ship, the first to have an all-welded hull. Of the nearly 4,700 ships of the Liberty class launched during the war, over 200 suffered catastrophic failure, some splitting in two while lying at anchor in port, and over 1,200 suffered some sort of severe damage due to fractures. The discipline of fracture mechanics emerged from these catastrophes. The all-welded ships were redesigned, eliminating, for example, sharp corners on hatches, and systematic procedures were developed for testing the fracture resistance of materials. In the early 1950s, failure by fracture cursed the British airline industry's efforts to establish passenger service using jet aircraft. Ill-placed rivet holes destroyed two of Britain's Comet aircraft, and played an important role in transferring the center of gravity for building civilian jet aircrafts from Britain to the United States. Aircrafts are now subjected to a systematic program of inspection that acknowledges that every structure has flaws, but that flaws greater than a certain size are intolerable. Testing procedures have continued to evolve in response to accidents, most recently after an incident (in the 1980s) in which part of the top of the fuselage of an Aloha airliner separated during flight, killing two people.

## 6.2    Fracture of a Homogeneous Solid

The strength of a material is its ability to resist an applied load without breaking or changing its shape. Therefore, let us first ask the seemingly simple question, how does a perfect (defect-free) solid break? To answer this question, consider a block of material of height $h$ and cross-sectional area $S$, pulled by a force $F$. The block separates into halves when its atoms are pulled beyond the breaking point. To estimate the force $F_c$, or the corresponding stress $\sigma_c$, required to reach the breaking point, we recall that the Young's modulus $Y$ relates the stress $\sigma$ on a material to its extension $\delta h$ through the relation

$$\sigma = -\frac{F}{S} = \frac{\delta h}{h} Y. \tag{1}$$

The ideal or cohesive strength of a perfect solid, i.e., the critical stress to reach the breaking point, is typically

$$\sigma_c = \frac{1}{10} Y. \tag{2}$$

If the material is under shear, the same estimate of $\sigma_c$ should be used, except that the Young's modulus $Y$ should be replaced by the shear modulus $\mu$. Except for some rather exotic materials, such as micrometer-sized whiskers, however, most solids have strengths in the range $10^{-2}Y$ to $10^{-4}Y$. This lower strength is caused by various defects, such as vacancies, interstitials, impurity atoms (point defects), dislocations (line defects), grain boundaries, heterogeneous interfaces, microcracks (planar defects), chemically-heterogeneous precipitants, twins, and other strain-inducing phase transformations (volume defects). The defects promote plasticity and premature fracture (see below). The mechanisms of crack nucleation that are described below provide insight into the phenomena involved.

However, the lower strength of certain materials, such as silicate glasses, which represent three-dimensional (3D) covalent networks, cannot be explained by the above deformation processes, since their microstructure is homogeneous except perhaps at very small length scales, of the order of 10 nm. In this case, the smooth surface of the glass, when it comes into contact with another solid material, produces sub-microscopic cracks, as point contacts generate very large localized stresses that cannot be relieved by plastic (or viscoelastic) deformation. The severity of the contact also determines the length and distribution of the cracks.

This example demonstrates the fact that, in order to find the best material to build, for example, a house, it is not enough to simply pull out the Periodic Table and find the element with the highest bonding strength and melting point, as this "exercise" will point to diamond, too expensive a material to build a house with! If one were to use, for example, vitreous mixture of silicon and oxygen, raw materials that are abundant and safe and form strong bonds, the attempt will again be a failure as soon as the material is hit with, say, a piece of stone. The failure of the Periodic Table in telling us which material to use is due to the fact that the relation between bonding energies and strength of materials is far from direct.

## 6.3   Introduction of Heterogeneity

In most engineering materials (as well as natural materials, such as rock) the presence of flaws or defects with various sizes, shapes and orientations makes fracture a very complex phenomenon. In fact, disorder comes into play in many ways during a fracture process. The effect of even small initial disorder can be enormously amplified during fracture. This makes fracture a collective phenomenon in which disorder plays a fundamental role. Due to disorder, brittle materials generally exhibit large statistical fluctuations in their fracture strengths, when nominally identical samples are tested under identical loading, giving rise to a distribution of fracture strengths (similar to distribution of the breakdown fields in the electrical and dielectric breakdown phenomena described in Chapter 5). Because of these statistical fluctuations, it is insufficient, and indeed inappropriate, to represent the fracture behavior of a disordered material by only its *average* properties, an idea which, as the previous chapters should have made clear, is usually used in mean-field and effective-medium approximations: Fluctuations are important to fracture nucleation and propagation and cannot be neglected

The traditional approaches to fracture mechanics (see, for example, Ewalds and Wanhill, 1986; Freund, 1990; Lawn, 1993) have certainly provided the framework for analyzing a wide variety of phenomena *without* considering the effect of disorder. These approaches are based on continuum fracture mechanics, some of the most important contributions of which will be summarized, described and discussed in Chapter 7. The basis for most of these traditional approaches is the important criterion developed by Griffith (1920; see below and also Chapter 7). The analogue of Griffith's analysis for the dielectric breakdown problem was already described in Section 5.1.1. He proposed that a *single* crack becomes unstable to extension when the elastic energy released in the crack extension by a small length $dc$ becomes equal to the surface energy required to create a length $dc$ of crack surface. However, Griffith's criterion was derived under *quasi-static* conditions and, moreover, it is presumably valid for materials that are essentially homogeneous, so that strong disorder plays no important role. Once the crack begins to move, the prevailing dynamical conditions render this criterion useless. In addition, the extension of this criterion to heterogeneous materials, as simple as polycrystalline ceramics with various crystalline orientations and/or grain boundary energies, is not obvious.

Accompanying the traditional phenomenological theories has been direct numerical modeling using the finite-element method (FEM). With a combination of computers and adroit mesh constructions, the stress field of a configuration of grains, fibers or cracks may be calculated by the FEM. The mesh size of the FEM must be smaller than the scale on which the stress field is expected to vary, which is therefore much smaller than the relevant length scale of the disorder. Therefore, only a small portion of a disordered material can be analyzed using the FEM, and full calculations must be performed for each of the many local configurations which are required to understand the statistical nature of the problem. To extend such small-scale FEM studies to larger length scales is still a formidable, if not

impossible, computational problem. We will describe in Chapter 7 some typical FE simulations of fracture propagation in solid materials.

Such difficulties have inspired further development of continuum mechanics approach to fracture on one hand, and development of many discrete models of fracture of materials on the other hand. The discrete models are typically based on lattices of elastic elements, such as springs and beams. The advantage of such models is that, at least over certain length scales, they allow disorder to be explicitly included in the models. We will study such discrete models in Chapter 8. Another type of discrete model of fracture and failure of materials is based on molecular dynamics simulations that consider propagation of a fracture at the atomic scale. We will describe and study this approach in Chapters 9 and 10. Both the lattice models and the MD simulations have also necessitated use of large-scale computer simulations.

In the present chapter we lay the foundations for our discussions of fracture phenomena, and describe the basic concepts that will be employed heavily in the subsequent chapters. We also describe and discuss the experimental techniques for measuring the most important properties of interest in facture of materials, so that when in the subsequent chapters we compare the theoretical predictions with the relevant experimental data and mention the technique by which the data have been collected, the reader will have a clear understanding of, and familiarity with, the technique. Also described in this chapter are the basic features of several important classes of materials, as they relate to their fracture properties.

## 6.4   Brittle Versus Ductile Materials

The most important qualitative fact in the mechanical properties of solid materials is that some are brittle and shatter in response to an external force, while others are ductile and merely deform in response to the blow. If we take a piece of a solid material, make a saw cut in it, and pull it, then, if the material is brittle, the tip of the saw cut sharpens spontaneously down to atomic dimensions and, similar to a knife blade one atom wide, it slices its way forward. In a ductile material, on the other hand, the tip of the saw cut blunts, broadens and flows, so that great effort is required to make the cut progress. The question is, why? Posing the question in a new guise, we ask, what makes a crack grow and propagate?

There is no completely satisfactory answer to the question of why some materials are brittle and others are ductile, as the giant stars of the Milky Way Galaxy, the long-dead true manufacturers of atoms, forgot to specify this property when writing down their technical specifications. The most well-developed investigation of this problem considers stationary, atomically sharp cracks in otherwise perfect crystals, and asks what happens when slowly increasing stresses are imposed on them. Rice and Thomson (1974) were probably the first to show how to estimate whether the crack will move forward in response to such a stress, or whether, instead, a crystal dislocation (i.e., a line of defects) will pop out of

the crack tip, causing the tip to become blunt. Brittleness and ductility are not, in fact, inherent in the atoms that make up a solid. For most solid materials there is a definite temperature at which they make a transition from brittle to ductile behavior which for example, is about 500°C for silicon. In Chapter 7 we will briefly describe theories that attempt to predict this transition temperature.

## 6.5    Mechanisms of Fracture

To understand how a fracture propagates in a solid material, it is essential to understand how a fracture is nucleated. At the atomic level, a crack or fracture is the result of breaking the interatomic bonds of a material. However, the answer to the all important question, "when do the atomic bonds break," is mostly material-specific, and depends critically on the morphology of a material. Normally, fractures are generated as a result of a stress or strain imposed on a material which causes its deformation and breakage of its interatomic bonds. The stress or strain can be applied externally, or can be generated internally by differential changes within the material. The cracks in the latter case are usually referred to as the *pre-existing cracks*. The differential changes can be caused by a temperature gradient, a transport process such as diffusion, chemical changes and reactions, or by shrinkage. One must also distinguish between the *nucleation* of a crack and its *propagation*. In some cases, a crack propagates by growing alone, while in other cases the propagation process is the result of coalescence of a multitude of smaller cracks. What follows is a brief discussion of several mechanisms of deformation of a material which leads to nucleation of cracks.

### 6.5.1    Elastic Incompatibility

If a solid material consists of rigid phases or grains, then cracks nucleate at the interface between the grains (and also in the grains themselves). This is due to the elastic incompatibility of the neighboring grains, caused by the differences in their composition and orientations. These differences result in different elastic strains in the grains, when a stress is applied to the material, leading to formation of local high-stress areas in the material that can be relieved only by formation of a crack.

### 6.5.2    Plastic Deformation

First introduced as a mathematical concept in the 19th century, the idea of a dislocation as a crystal defect was hypothesized simultaneously by Orowan (1934), Polanyi (1934), and Taylor (1934), mainly to explain the less-than-ideal strength of crystalline materials. Only much later, in the 1950s, was the existence of dislocations experimentally confirmed (Hirsch *et al.*, 1956). Currently, such ubiquitous crystal defects are routinely observed by various means of electron microscopy.

Low-temperature shear deformation of crystalline materials (e.g., ceramics) occurs by gliding of individual dislocations or the coordinated movement of arrays of partial dislocations. The shear can be localized in a narrow band which, if it meets some sort of a microstructural barrier (e.g., a grain boundary or a particle from another phase of the material), leads to very high local stresses at the band's tip, resulting in the nucleation of a crack. The direction of the shear as well as the location of the crack are both influenced very strongly by the crystal structure and the strength of the interface between the shear band and the barrier. However, instead of nucleating a crack, the high stresses can also be relieved by some sort of generalized plastic deformation. Many materials are unable to relieve high stresses caused by plastic deformation, and therefore form cracks.

Over the last seven decades, experimental and theoretical developments have firmly established the principal role of dislocation mechanisms in defining material strength. It is now universally accepted that the macroscopic plasticity properties of crystalline materials are derivable, at least in principle, from the behavior of their constituent defects. However, this fundamental understanding has not translated into a quantitative theory of crystal plasticity based on dislocation mechanisms. One difficulty is the multiplicity and complexity of the mechanisms of dislocation motion and interactions, which leave little hope, if any, for a quantitative analytical approach. The situation is further exacerbated by the need to trace the evolution of a large number of interacting dislocations over long periods of time, which is required for any calculation of plastic response in a representative volume element of the material.

### 6.5.3   Coalescence of Plastic Cavities

An operating mechanism for crack nucleation, especially in ductile materials that contain rigid inclusions, is the coalescence of cavities. When a stress is applied to the material, the ductile matrix is deformed, with its mechanism of deformation being either slip (as in crystalline materials) or shear deformation (as in amorphous materials). The rigid inclusions do not deform, and therefore the interface between them and the matrix separates, followed by development of plastic cavities around the inclusions. Further deformation of the matrix forces the cavities to grow. Alternatively, if the temperature of the system is high enough, the cavities grow by a diffusion process. At some point the local cavities begin to interact with each other, and eventually merge and form a crack.

### 6.5.4   Cracks Initiated by Thin Brittle Films

If a strong material is covered by a thin brittle film, fracture of the film can lead to the fracture of the material itself in the bulk, even if the material is ductile. An example is a nitride layer on steel. In this case, the deformation of the film causes its fracture which then propagates at high speeds, penetrating the material itself. Degradation of the surface of materials can lead to the same effect.

### 6.5.5 *Crazing*

Crack nucleation by crazing occurs in amorphous polymeric materials. When such materials are deformed by an applied stress, the polymeric chains rotate and, if the strain is large enough, become aligned in the direction of the maximum extensional strain. Crazing then involves formation of planar arrays of fine voids that are normal to the tensile stress. The distance between the voids is filled by ligaments of aligned polymer chains. If the deformation is strong enough, the ligaments eventually break and help the voids to merge.

### 6.5.6 *Boundary Sliding*

If a material contains rigid blocks (as in polycrystalline materials), and if the temperature of the system is high enough, then, it is deformed by sliding of the rigid blocks. The sliding is stopped at the triple point grain corners, and cracks that are wedge-shaped are formed. In addition, rigid particles can help nucleate plastic cavities during sliding which then grow, coalesce and form cracks.

## 6.6   Conventional Fracture Modes

There are three symmetrical ways of loading a solid material with a crack. These are known as *modes*, and are illustrated in Figure 6.1. A generic loading situation produced by some combination of forces without any particular symmetry is usually referred to as *mixed mode* fracture. Although understanding mixed-mode fracture is obviously of practical importance, our focus will primarily be upon the physics of fracture propagation rather than upon engineering applications. Therefore, we will restrict our attention to the cases in which the loading has a high degree of symmetry, but will also briefly discuss the mixed mode case.

The fracture mode that we will mainly deal with in this book is Mode I (*opening mode*), where the fracture faces, under tension, are displaced in a direction normal to the fracture plane. In Mode II (*sliding mode*), the motion of the fracture faces is that of shear along the fracture plane. Mode III (*tearing mode*) fracture corresponds to an out of plane tearing motion where the direction of the stresses at the fracture faces is normal to the plane of the sample. One experimental difficulty of Modes II and III is that the fracture faces are not pulled away from one another, and thus contact along the fracture faces still occurs. The resulting friction between the fracture faces contributes to the forces acting on the crack, but its precise measurement is difficult.



FIGURE 6.1. The three basic fracture modes.

For these reasons, Mode I corresponds most closely to the conditions used in most experimental and theoretical work on brittle fracture of solids, since there is always a tendency for a brittle crack to seek an orientation that minimizes the shear loading. This is consistent with crack extension by progressive stretching and rupture of cohesive bonds across the crack plane. In 2D isotropic materials, Mode II fracture cannot easily be observed, because slowly propagating fractures spontaneously orient themselves so as to make the Mode II component of the loading vanish near the crack tip (Cotterell and Rice, 1980). Mode II fracture is, however, observed in strongly anisotropic materials. For example, friction and earthquakes along a pre-defined fault are examples of Mode II fracture where the binding across the fracture interface is considerably weaker than the strength of the bulk of the material. Pure Mode III fracture, although experimentally difficult to achieve, is sometimes used as a model system for theoretical studies, since in this case the equations of elasticity simplify considerably. Analytical solutions obtained in this mode (some of which will be described in Chapter 7) have provided considerable insight into the fracture process.

## 6.7  Stress Concentration and Griffith's Criterion

Inglis (1913) analyzed the stress distribution in a uniformly-stressed plate containing an elliptical cavity at its center. His work, which represents an important precursor to that of Griffith (1920), showed that the stress around a sharp notch or corner may be many times larger than the applied stress, hence providing the important clue that even sub-microscopic voids or flaws can weaken a material. Most importantly, his analysis established that the limiting case of an infinitesimally narrow ellipse can be considered as representing a crack. We summarize Inglis' analysis here.

Consider a plate that contains an elliptical cavity of semi-axes $c$ and $b$, which are small compared to the dimensions of the plate. We apply a uniform tension $\sigma^0$ along the $y$-axis. The system is shown in Figure 6.2. The cavity's boundary is stress-free, and Hooke's law of linear elasticity holds everywhere in the plate. The equation for the ellipse is given by

$$\frac{x^2}{c^2} + \frac{y^2}{b^2} = 1, \tag{3}$$

based on which it is easy to show that the radius of curvature of the ellipse's boundary given by, $\mathcal{Y} = b^2/c$, achieves its maximum at point $A$ shown in Figure 6.2. Point A is also where the stress is maximum and is given by

$$\sigma_m = \sigma^0 \left(1 + \frac{2c}{b}\right) = \sigma^0 \left(1 + 2\sqrt{\frac{c}{\mathcal{Y}}}\right), \tag{4}$$

which, in the limit $b \ll c$ that the cavity represents a crack, reduces to

$$\frac{\sigma_m}{\sigma^0} = \frac{2c}{b} = 2\sqrt{\frac{c}{\mathcal{Y}}}. \tag{5}$$

FIGURE 6.2. The elliptical cavity in a plate, subjected to a uniform applied stress. Point A represents the notch tip.





FIGURE 6.3. Stress concentration at the elliptical cavity for $c = 3b$.

The ratio $\sigma_m/\sigma^0$ is called the *stress-concentration factor*, which is the mechanical analogue of the field-multiplication factor defined in Section 5.1.1 for the problem of dielectric breakdown with an elliptical conductor. Since as $b \to 0$ the radius of curvature becomes very small, it is clear that $\sigma_m$ can become much larger than the applied stress $\sigma^0$.

Of particular interest is the local stresses along the $x$-axis. This is shown in Figure 6.3 for $c/b = 3$, where we present the stresses $\sigma_{xx}$ and $\sigma_{yy}$. The stress $\sigma_{yy}$ decreases from its maximum value of $7\sigma^0$ at point A to an asymptotic value of $\sigma^0$, while $\sigma_{xx}$ rises from a zero value at A, reaching its maximum value at a point very near the boundary of the cavity, beyond which it approaches 0 at large distances.

Note that the value of the stress depends on the shape of the cavity rather than its size. Therefore, although it appeared that Eq. (5) can be used for estimating the stress-concentration factors of such systems as the surface notch, a nagging

FIGURE 6.4. Incremental extension of a fracture of length $c$ through $dc$, under the applied stress.

question hindered further progress in understanding of fracture mechanics at Inglis time: If the analysis of Inglis is applicable to a crack system (predicting a size-independent stress), then why in practice large cracks appear to grow and propagate more easily than the small ones? In addition, since the result of Inglis was in terms of the radius of the curvature, the natural question to ask was, what is the physical significance of the radius of curvature at the tip of a real crack?

Inglis' work was followed up by Griffith (1920) who was interested in the strength of inorganic glasses. He showed that the low strength of these materials, compared to the theoretical estimates described earlier, was due to the presence of sub-microscopic cracks. To reach this conclusion, Griffith analyzed the system shown in Figure 6.4 which shows an elastic body that contains a plane-crack surface $S$ of length $c$, subjected to loads applied at its outer boundary. Griffith's main idea was to analyze this problem as a reversible thermodynamic system, seeking the configuration that minimizes the total free energy of the system. Under this condition, the crack would be in a state of equilibrium, and thus on the verge of propagation.

If the crack undergoes extension, the energy $\mathcal{H}$ of the system associated with this motion is the sum of the mechanical and surface energies. The mechanical energy $\mathcal{H}_M$ is itself the sum of two terms, the strain potential energy stored in the elastic material, and the potential energy of the outer applied loading system (which, in magnitude, is equal to the work associated with the displacements of the loading points). The surface contribution $\mathcal{H}_S$ is the free energy expended in generating the new crack surfaces. Thus,

$$\mathcal{H} = \mathcal{H}_M + \mathcal{H}_S. \tag{6}$$

Thermodynamic equilibrium is reached when the mechanical and surface energies for a virtual crack extension $dc$ (see Figure 6.4) are balanced. However, the mechanical energy favors the crack extension (i.e., $d\mathcal{H}_M/dc < 0$) while the surface energy opposes it ($d\mathcal{H}_S/dc > 0$). Thus, the Griffith energy-balance concept is expressed through the equilibrium requirement that

$$\frac{d\mathcal{H}}{dc} = 0. \tag{7}$$

Therefore, a crack would extend or contract reversibly for small displacements from the equilibrium length, according to whether $d\mathcal{H}/dc$ is negative or positive,

respectively. For over 80 years, Eq. (7) has remained a pillar of the classical continuum theory of brittle fracture.

To develop his theory further, Griffith took advantage of the Inglis' solution for an elliptical cavity described above. It can be shown that for a system under constant applied stress (during crack formation), $\mathcal{H}_M = -\mathcal{H}_E$, where $\mathcal{H}_E$ is the strain potential energy stored in the elastic material, mentioned above, and the negative sign is due to the fact that crack formation *reduces* the mechanical energy. Using the solution of Inglis, it is not difficult to compute the strain energy density, from which one obtains (by integrating the energy density over dimensions that are large compared with the length of the crack), $\mathcal{H}_E = -\mathcal{H}_M = \pi c^2 (\sigma^0)^2 / Y'$, where $Y'$ is equal to the Young's modulus $Y$ in plane stress (thin plates), and $Y' = Y/(1 - v_p^2)$ in plane strain (thick plates), with $v_p$ being the Poisson's ratio. Since, for a unit width of the crack front, one has $\mathcal{H}_S = 4c\Gamma$, where $\Gamma$ is the free surface energy per unit area, one obtains

$$\mathcal{H} = 4c\Gamma - \frac{\pi c^2 (\sigma^0)^2}{Y'}. \tag{8}$$

If we now apply Griffith's criterion, Eq. (7), and identify $\sigma^0 = \sigma_c^0$ as the critical stress, we obtain

$$\sigma_c^0 = \sqrt{\frac{2Y'\Gamma}{\pi c}}. \tag{9}$$

Equation (9) is the famous Griffith relation, and is the mechanical analogue of Eq. (5.8), the critical value of the far-field electric field for dielectric breakdown. Griffith also succeeded in qualitative verification of Eq. (9) by carrying out experiments on an inorganic glass.

Because $d^2\mathcal{H}/dc^2 < 0$, the energy of the system at equilibrium is maximum, and therefore its configuration is *unstable*. That is, for $\sigma^0 < \sigma_c^0$ the crack remains stationary at its initial size $c$, whereas for $\sigma^0 > \sigma_c^0$ it propagates spontaneously *without limit*. Note, however, that an unstable crack may ultimately be *arrested* at some point, which is often the case with cracks around contacts and inclusions. In this case, further increase in the applied loading may lead to a second, catastrophic instability configuration.

## 6.8   The Stress Intensity Factor and Fracture Toughness

An alternative, but equivalent, approach to determining the critical stress $\sigma_c^0$ was developed by Irwin (1958). He was the first to note that the stress field at a point $(r, \theta)$ near the fracture tip, measured in polar coordinates with the crack line corresponding to $\theta = 0$, can be determined analytically. This problem will be discussed in detail in Chapter 7, but for now it suffices to record the solution for the stress components for Mode I fracture:

$$\sigma_{xx} = \sigma^0 \sqrt{\frac{c}{2r}} \left[ 1 - \sin\left(\frac{1}{2}\theta\right) \sin\left(\frac{3}{2}\theta\right) \right] \cos\left(\frac{1}{2}\theta\right), \tag{10}$$

$$\sigma_{yy} = \sigma^0 \sqrt{\frac{c}{2r}} \left[ 1 + \sin\left(\frac{1}{2}\theta\right) \sin\left(\frac{3}{2}\theta\right) \right] \cos\left(\frac{1}{2}\theta\right), \tag{11}$$

$$\sigma_{xy} = \sigma^0 \sqrt{\frac{c}{2r}} \, \sin\left(\frac{1}{2}\theta\right) \cos\left(\frac{1}{2}\theta\right) \cos\left(\frac{3}{2}\theta\right), \tag{12}$$

and $\sigma_{zz} = \nu'(\sigma_{xx} + \sigma_{yy})$, where $\nu' = 0$ for plane stress and $\nu' = \nu_p$ for plane strain, with $\nu_p$ being the Poisson's ratio. The other components of the stress tensor are zero. These results can also be written in terms of $\sigma_{rr}, \sigma_{\theta\theta}, \sigma_{r\theta}$, etc. Qualitatively similar equations also hold for Mode II fracture. The results for Mode III fracture are particularly simple, as only $\sigma_{xz}$ and $\sigma_{yz}$ are non-zero. Therefore, Eqs. (10)–(12) can be written in a general form:

$$\sigma_{ij} = \sigma^0 \sqrt{\frac{c}{2r}} \, f_{ij}(\theta). \tag{13}$$

Irwin introduced the quantity $K = \sigma^0 (\pi c)^{1/2}$ as the *stress intensity factor*. Since, in general, the stress intensity factor and the function $f$ depend on the fracture mode (I, II, or III), and as $f$ also depends on the instantaneous crack velocity $v$, Eq. (13) is written in a very general form:

$$\sigma_{ij} = \frac{K_\beta}{\sqrt{2\pi r}} f_{ij}^\beta (v, \theta), \tag{14}$$

where $\beta$ indicates the fracture modes, $\beta =$I, II, and III. For each of the three symmetrical loading configurations, $f_{ij}^\beta(v, \theta)$ in Eq. (14) is a known universal function. The stress intensity factor $K_\beta$ contains all the detailed information about sample loading and history, and is determined by the elastic fields that develop throughout the material, but the stress that locally drives the fracture is one which is present at its tip. The stress intensity factors are related to the flow of energy into the crack tip. A fracture can be viewed as a sort of sink that dissipates built-up energy in a material. Therefore, the amount of energy flowing into a fracture tip influences its behavior. The theoretical aspects of this view will be discussed in Chapter 7. Thus, $K_\beta$ determines entirely the behavior of a fracture, and much of the study of fracture processes is focused on either calculating or measuring this quantity. The universal form of the stress intensity factor allows a complete description of the behavior of the tip of a fracture where one need only carry out the analysis of a given problem within the universal elastic region (see Chapter 7).

What happens if the material contains complicating factors, such as hetero-geneity and anisotropy? Such complications destroy the symmetry that exists in homogeneous and isotropic materials. For example, for a material in which the elastic properties on opposing sides of a plane-crack interface are asymmetric, the crack tip fields will also be asymmetric. Therefore, for example, a crack inter-face between two dissimilar materials, subjected to tensile loading, will exhibit not only Mode I behavior, but some Mode II and Mode III as well. Despite such complications, it is now generally accepted that the essential $r^{-1/2}$ singularity that Eqs. (10)–(14) exhibit is not changed by such complexities, and therefore the stress intensity factors can still be superposed. Therefore, for arbitrary loading configu-rations, the stress field around the crack tip is given by three stress intensity factors

$K_\beta$ which lead to a stress field that is a linear combination of the pure modes:

$$\sigma_{ij} = \sum_{\beta=1}^{3} \frac{K_\beta}{\sqrt{2\pi r}} f_{ij}^{\beta}(v, \theta). \tag{15}$$

The critical condition for crack propagation can now be expressed in terms of the critical value $K_c$ of the stress intensity factor, which is usually referred to as the *fracture toughness*. Thus, in terms of the critical energy $\mathcal{H}_c$, the fracture toughness is given by

$$K_c = \sqrt{\mathcal{H}_c Y'}. \tag{16}$$

We should emphasize, as already mentioned above, that the Griffith–Irwin prediction for the critical stress $\sigma_c^0$ (or the fracture toughness $K_c$) is valid for the *onset* of growth under static conditions, *and* for homogeneous materials. As soon as the crack begins to grow, the stress field around it changes dynamically. In particular, if the crack propagates at high speeds, the inertial effects substantially change the stress field. The Griffith–Irwin approach has nothing to offer for these changes. In other words, the Griffith–Irwin criterion can tell us *when* a brittle crack may extend, but has nothing to say about *how* it will extend. In addition, the $r^{-1/2}$ singularity at the tip of the crack cannot be reconciled with any real fracture process, as there is no solid that can resist an infinite stress anywhere in its structure. The root of this singularity is in the assumptions that the Hooke's law (linear elasticity) is operative everywhere in the material, and that a continuum approximation can describe the state of the system. These assumptions break down for the region in the vicinity of the crack tip, and necessitate a reclassification of the region around the tip; this is discussed in the next section.

## 6.9   Classification of the Regions Around the Crack Tip

Many complex phenomena are active in the vicinity of a crack tip that vary, depending on the material, from dislocation formation and emission in crystalline materials to the complex unraveling and fracture of intertangled polymer strands in amorphous polymers. Fracturing and the complex dissipative processes occurring in the vicinity of the crack tip occur due to large values of the stress field as one approaches the tip. As discussed above (and will also be considered in detail in Chapter 7 where we describe formation of fracture nucleation and propagation by continuum mechanics), if the material around the crack tip were to remain linearly elastic until fracture, the stress field at the crack tip would be singular. Since a real material cannot support such singular stresses, the assumption of linearly elastic behavior in the vicinity of the tip must break down and material-dependent dissipative processes must begin playing an important role. Thus, at first glance, a universal description of fracture, in terms of the stress intensity factor and the function $f_{ij}$ described above, may seem a hopeless task. However, a way for attacking this problem was proposed by Orowan (1955) and Irwin (1956) who suggested independently that the region around the crack tip should be divided into three

separate regions:

(1) *The cohesive zone,* which is the region immediately surrounding the crack tip, in which all the nonlinear dissipative processes that allow a crack to move (forward) are assumed to occur. In continuum fracture mechanics, detailed description of this zone is avoided, and is simply characterized by the energy $\Gamma$, per unit area of crack extension, that it will consume. The size of the cohesive zone is material-dependent, ranging from nanometers in glass to microns in brittle polymers. Its typical size is the radius at which an assumed linear elastic stress field surrounding the fracture tip would equal the yield stress of the material.

(2) *The universal elastic region,* which is the region outside of the cohesive zone for which the response of the material can be described by linear continuum elasticity. Outside of the cohesive zone, but in the vicinity of the fracture tip, the stress and strain fields take on universal singular forms which depend only on the symmetry of the externally applied loads. In 2D the singular fields surrounding the cohesive zone are completely described by the three stress intensity factors which incorporate all the information regarding the loading of the material. As discussed above, the stress intensity factors are related to the energy flux into the cohesive zone. The larger the overall size of the material containing the crack, the larger this region becomes. Roughly speaking, for given values of the stress intensity factors, the size of the universal elastic region scales as $\sqrt{L}$, where $L$ is the macroscopic length scale on which forces are applied to the material. Thus, as $L$ increases, the assumptions of continuum fracture mechanics become progressively more accurate.

(3) *Outer elastic region,* which is the region far from the crack tip in which stresses and strains are described by linear elasticity. Details of the solution of the equations, describing fracture propagation, in this region depend only on the locations and strengths of the loads, and the shape of the material. For some special cases, analytical solutions are available, but in general one must resort to numerical simulation. That deriving these solutions is possible is because, so far as linear elasticity is concerned, viewed on macroscopic scales, the cohesive zone shrinks to a point at the fracture tip, and the fracture itself becomes a branch cut. Thus, replacing the complex domain in which linear elasticity holds with an approximate one that needs no detailed knowledge of the cohesive zone is another approximation that becomes increasingly accurate as the dimensions of the sample, and hence the size of the universal elastic region, increase. The assumption that the cohesive zone in a material is encompassed within the universal elastic region is sometimes called the assumption of *small-scale yielding*.

The dissipative processes within the cohesive zone determine the fracture energy $\Gamma$. If no dissipative processes other than the direct breaking of the atomic bonds take place, then $\Gamma$ will be a constant that depends on the bond energy. In general though, $\Gamma$ is a complex function of both the crack velocity and history, and differs by orders of magnitude from the *surface energy*—the amount of energy required

to sever a unit area of atomic bonds. No general first principles description of the cohesive zone exists, although numerous models have been proposed (see, for example, Lawn, 1993). We will come back to this important issue in Chapter 7.

## 6.10    Dynamic Fracture

Our discussion so far has been limited to static fractures. However, in practice dynamical effects are important and must be considered. To understand how a dynamic situation may come about, suppose that an unbalanced force acts on any volume elements within a material that contains cracks. Then, that element will be accelerated, thereby acquiring kinetic energy. The system will then be in a dynamic state so that, as pointed out and emphasized above, the Griffith–Irwin static equilibrium condition will no longer apply. Under certain conditions, the growth of the crack may be slow (for example, when, compared to the mechanical energy, the contribution of the kinetic energy is insignificant), in which case the material may be considered as being in a quasi-steady-state condition.

There are two scenarios by which the state of a cracked material may become dynamical. One is when a crack reaches an unstable state in its length: The material receives kinetic energy contributed by the inertia of the material that surrounds the rapidly-separating walls of the crack. One then has a *running* crack which is characterized by a rapid acceleration toward a *terminal velocity* $v_c$, and is governed by the speed of elastic waves in the solid. As will be discussed in Chapter 7, the prediction of linear continuum fracture mechanics for the value of the terminal velocity $v_c$ did not agree with experimental observations, and because of this the subject was controversial for a long time and was resolved only recently. In the second scenario a dynamical state arises when the applied loading changes rapidly with the time, as in, for example, impact loading.

A very common dynamical effect is *fatigue*. It has been seen in many experiments that, often a material that has resisted the same external load many times without developing cracks, suddenly does so after the external load has been applied a certain number of times. If the external load is applied periodically in time, then the phenomenon is called *cyclic fatigue*, and the number of times that the external load must be applied for the crack to develop is called *failure life*. It has been found empirically that the number $N_c$ of the cycles that the external load must be applied scales with the amplitude $\mathcal{A}$ of the load as

$$N_c \propto (\mathcal{A} - \mathcal{A}_e)^{-\alpha}, \tag{17}$$

where $\mathcal{A}_e$ is called the *endurance limit*. Clearly, if the amplitude of the applied load is less than $\mathcal{A}_e$, then the material will not break at all. The value of the exponent $\alpha$ has been found to be around $8 - 10$. Equation (17) is usually called the *Besquin law*.

Another important dynamical effect is *stress corrosion cracking*. Baker and Preston (1946) first reported that the toughness of glass reduces considerably if it is in a humid environment, since water penetrates the glass at the crack tip where

the crystalline structure of the glass is relatively open. Once inside the glass, water forms a base with existing sodium ions which corrodes the region in the vicinity of the crack tip, hence lowering its toughness and increasing the likelihood of brittle fracture. Aluminum and titanium, two heavily-used metals in aircraft, suffer most from stress corrosion cracking.

The mechanism that leads to stress corrosion cracking is either anodic or cathodic. That is, the phenomenon can be suppressed in an electrolytic environment by placing either the anode or the cathode on the material and the corroding agent as electrolytic medium. For example, hydrogen embrittlement of metals is the most common cathodic process. The anodic process also occurs in metals that are coated by a layer of oxide to be protected from the environment. If the coating is opened at the crack tip, the metal will be exposed to the anodic agent at the tip. Under this condition, the velocity of fracture propagation would be controlled by the rate of the chemical reactions. Since these reactions are typically slow, the fracture propagates slowly, which is why, for example, it takes aircraft a long time to develop stress corrosion cracks in their fuselage.

# 6.11   Experimental Methods in Dynamic Fracture

We now describe and discuss some of the main experimental methods that are used in studies of dynamic fracture. These methods vary greatly and their use depends on both the specific phenomenon that is under study and on the experimental resources at hand. In a typical experiment stress is applied externally at the boundary of the system and its response and the resulting behavior of the fracture are observed and measured. During the time that the crack propagates one can measure its position and velocity, the time-dependent stress field at the crack's tip, the acoustic emissions resulting from the crack motion, as well as the resulting fracture surface. In what follows we describe the typical ways by which the various quantities of interest are measured. Our discussion in this section follows closely that of Fineberg and Marder (1999).

## 6.11.1   Application of External Stress

The externally-applied stress distribution determines the stress field in the close vicinity of the crack tip or, equivalently, the stress intensity factor, and hence, is the driving force for advancement of a fracture. Two basic types of loading are typically used in fracture experiments, static and dynamic, and what follows is the description of each type.

### 6.11.1.1   Static Stress

In such experiments either the boundary conditions or the applied stresses are constant, thus imprinting an initial static stress distribution onto the sample material. Depending on the applied loading and boundary conditions, the stress intensity

FIGURE 6.5. Three typical experimental configurations.

factor (or stored energy density) along the prospective path of a crack can increase, resulting in a continuously accelerating crack, or decrease, leading to a decelerating and possibly arrested crack. A few examples of the common loading conditions that are used are shown in Figure 6.5 where *single-edge notched* (SEN), *double-cantilever beam* (DCB), and *infinite strip* (IS) loading conditions are presented. The SEN condition is sometimes used to approximate fracture propagation in a semi-infinite system. When the external loading is a constant stress applied at the vertical boundaries of the sample, then for a large enough sample the stress intensity factor $K_I$ is proportional to $\sigma \sqrt{l}$, and therefore the energy release rate $\mathcal{H}$ is given by, $\mathcal{H} \propto \sigma^2 l$ (see Chapter 7 for additional theoretical details), where $\sigma$ is the applied stress and $l$ is the length of the crack. This configuration is used, for example, to study the behavior of an accelerating crack.

In the IS configuration, the sample is loaded by displacing its vertical boundaries by a constant amount. Under this condition, the energy release rate is constant for a crack that is sufficiently far from the horizontal boundaries of the sample, and thus this loading configuration is amenable to the study of a crack moving in steady-state.

In the DCB configuration, a constant separation of the crack faces is imposed at $l = 0$, $\mathcal{H} \propto l^{-4}$ is a decreasing function of $l$, and hence can be used to cause crack arrest. How is the DCB configuration used to study dynamic fracture? An initially imposed seed crack of length $l = l_0$ would propagate as soon as $\mathcal{H}$ exceeds the limit imposed by the Griffith condition, i.e., when $d\mathcal{H}/dl = 0$. Under ideal conditions, the crack propagates for an *infinitesimal* distance and then stops, because in DCB configuration $\mathcal{H}$ is a decreasing function of $l$. Although the Griffith criterion assumes that the initial crack is as sharp as possible, what is prepared in the laboratory by cutting rarely yields a tip that meets this condition. We may view the initial seed crack as having a finite radius at its tip, thereby blunting the stress singularity and allowing a substantially higher energy density to be imposed in the system prior to fracture than what is allowed by a sharp crack. Thus, there is excess elastic energy that drives the crack beyond the constraints imposed by an initially sharp crack which, in the case of the DCB configuration with constant

separation imposed, can cause a crack to propagate well into the sample before crack arrest occurs. Nonlinear material deformation around the tip, plastic flow induced by the large stress build-up, and from crack-tip shielding that results from the formation of either micro-cracks or small bridges across the crack faces in the near vicinity of the tip, can also cause blunting of the singularity around a crack tip. The DCB configuration can also generate an accelerating crack by imposing a constant stress (instead of constant separation) at the crack faces. Under this condition, the (quasi-static) energy release rate increases with the crack length $l$ as

$$\mathcal{H} = \frac{12\sigma^2 l^2}{Y w^2 d^3},  \tag{18}$$

where $\sigma$ is the stress applied at opposite points on the crack faces at the edge of the sample, and $w$ and $d$ are, respectively, the thickness and half-width of the sample.

### 6.11.1.2  Initiation of Fractures

The stress singularity at the tip of a crack, as its radius of curvature approaches zero (see above), implies that initiation of a fracture under static loading configurations is strongly dependent on the initial radius of the crack tip and hence on the preparation of the initial crack. However, the stress build up that precedes fracture initiation can be taken advantage of for loading a material with an initial energy density before the onset of fracture. This is, however, extremely difficult as experimental reproducibility of the stress at fracture initiation is non-trivial. In some materials one can achieve a reproducible stress at fracture initiation by first loading the system to the desired stress and then either waiting for some time for the material to fracture as a result of noise-induced perturbations, or by sharpening the initial crack, once the desired initial conditions have been reached. These tricks do not, however, work very well in such brittle materials as ceramics.

### 6.11.1.3  Dynamic Stress

In some applications, such as the study of crack initiation before the material surrounding the crack tip has had time to react to the applied stress, very high loading rates are desirable. A common way to achieve this is by loading an initially seeded sample by collision with a guided projectile. In this way loading rates as high as $\dot{K}_I \sim 10^9$ MPa$\sqrt{\text{m}}$/s (Prakash and Clifton, 1992) have been achieved. An alternative way for producing high loading rate is by sending a very large current through a folded conducting strip, inserted between the two faces of an initial crack, which induces magnetic repulsion between adjacent parts of the strip, enabling direct loading of the crack faces (Ravi-Chandar and Knauss, 1982). A high loading rate can also be produced by discharging a capacitor-inductor bank through the strips. This technique has been utilized for producing a pressure pulse with a step function profile on the crack faces having loading rates of the order of $\dot{K}_I \sim 10^5$ MPa$\sqrt{\text{m}}$/s in experiments designed to investigate the response of a moving crack to rapidly changing stresses.

## 6.11.2  Direct Measurement of the Stress Intensity Factor

The stress intensity factor can be directly measured by optical methods, which can measure the energy release rate. Two common methods are the method of caustics, and photoelasticity.

### 6.11.2.1  The Method of Caustics

The technique was originally proposed by Manogg (1966) with significant contributions by Theocaris and Gdoutos (1972) and Kalthoff (1987) (who also provided a review of this method) for transparent materials, and by Rosakis *et al.* (1984) for opaque materials. The method, applicable to thin quasi-2D plates, uses the deflection of an incident collimated beam of light as it either passes through transparent material, or is reflected by an opaque material, that surrounds the crack's tip. Due to so-called Poisson contraction generated by the high tensile stresses near the tip, the initially flat faces of a plate will deform inwardly which creates a lensing effect and diverts light away from the crack's tip. The diverted rays form a 3D surface in space in which no light propagates. When this light is imaged on a screen, a shadow (hence the name *shadow-spot* that is sometimes used) is observed which is bounded by a caustic surface or a region of high luminescence formed by the locus of the diverted rays. From the shape of the caustic surface, which is recorded by a high speed camera, the instantaneous value of the stress intensity factor is estimated. This method works well with the caveat that estimating the stress intensity factor is based on a certain assumption that, as discussed in Chapter 7 [see the discussion before and after Eqs. (7.50)–(7.52)], must, in the immediate vicinity of the crack tip, break down as the material's yield stress is approached. Therefore, care must be taken that the curve on the material that maps onto the caustic is well away from the cohesive zone surrounding the crack tip; see, for example, Rosakis and Freund (1981).

### 6.11.2.2  Photoelasticity

This method, coupled with high-speed photography, is also used for measuring the stress distribution, and hence the stress intensity factor, induced by a moving crack (Kobayashi, 1987). It is based on the birefringence induced in most materials under an imposed stress, which causes the rotation of the plane of polarization light moving through the material. The induced polarization depends on the properties of the stress tensor which are rotationally invariant, and therefore can depend only on the two principal stresses $\sigma_1$ and $\sigma_2$. Moreover, there should be no rotation of polarization when the material is stretched uniformly in all directions, in which case the two principal stresses are equal, and therefore the angular rotation of the plane of polarization must be of the form, $c(\sigma_1 - \sigma_2)$, where $c$ is a constant that is determined experimentally. If stresses of a 2D problem are calculated analytically, the results can be substituted into this expression and compared with experimental fringe patterns obtained by viewing a reflected or transmitted beam of incident polarized light through a polarizer. The observed intensity depends on the phase

difference picked up while traversing the material, and hence provides a quantitative measure of the local value of the stress field. The application of this method to transparent materials is straightforward. These methods have also been extended to opaque materials by the use of birefringent coatings which, when sufficiently thin, mirror the stress field at the surface of the underlying material. Dally (1987) reviewed the applications of these methods. Similar to the method of caustics, quantitative interpretation of these measurements is limited to the region outside of the plastic zone.

### 6.11.3   Direct Measurement of Energy

Direct measurement of the energy release rate, as a function of the velocity of a moving crack, can be obtained by constraining a crack to propagate along a long and narrow strip; see Figure 6.5. The advantage of this method is that it relies only on symmetry properties of the system, and hence does not require additional assumptions regarding, for example, the size or properties of the cohesive zone. A series of experiments, using a long strip geometry and varying the value of $\delta$ (as shown in the Figure 6.5), results in a direct measure of $\mathcal{H}(v)$, where $v$ is the velocity of the crack. In the experiments of Sharon *et al.* (1995) using polymethylmethacrylate (PMMA) (see Chapter 7), steady-state mean velocities were attained when the crack length exceeded roughly half the strip height. Their measurements of $\mathcal{H}(v)$ agreed well with results previously obtained with PMMA by means of the methods of caustics (see above) reported by Pratt and Green (1974).

### 6.11.4   Measurement of Fracture Velocity

Under dynamic conditions, the velocity of the tip of a crack generally accelerates to values of the order of the sound speed in the material. Since the duration of a typical experiment is of the order of 100 $\mu$s, one needs relatively high-speed measurement techniques. Three common methods, based on either high-speed photography, resistance measurements, or the interaction of a moving crack with ultrasonic waves, have been used in the past which are now briefly described.

#### 6.11.4.1   High-Speed Photography

This method is the most straightforward technique for measuring the velocity of a moving crack. It can be used in conjunction with instantaneous measurements of the stress intensity factor by means of the method of caustics or photoelasticity discussed above. It also has some major shortcomings. For example, although the frame rates of high speed cameras are typically between 200 kHz and 10 MHz, the cameras are capable of photographing only a limited number, say 30, of frames. Thus, this method can either provide measurements of the *mean velocity* (with the average taken over the interval between the frames) at a few points, or can, at the highest photographic rates, provide a detailed measurement of the crack

velocity over a short, say about 3 $\mu$s, interval. Moreover, precision of this method is obviously limited by the accuracy at which the location of the crack tip can be determined from a photograph.

These problems can, to some extent, be overcome by using a streak camera (Bergkvist, 1974). In this method, a film is pulled past the camera's aperture at high speed. The material is illuminated from behind so that, at a given instant, only the light passing through the crack is photographed. Since one can force a crack to propagate along an essentially straight line, the exposed film provides a continuous record of its length as a function of time. The basic resolution of the measurements depends on the film's velocity and that of the high-speed film used, and on the post-processing performed on the film in order to extract the velocity measurement and the stability of the film's travel velocity. The same type of experiments have also been carried out (Döll, 1975) by high-speed measurements of the total beam intensity that penetrates the material. If the crack does not change its shape, the beam intensity depends linearly on the crack's length.

### 6.11.4.2  Measurement of Resistivity

Another method of measuring the velocity of a rapidly moving crack is by adhering a grid of thin, electrically-conductive strips to a sample material prior to fracture. Crack propagation causes the crack faces, and therefore the conducting strips, to separate. Therefore, if, for example, the strips are connected in parallel to a current source, measurement of the grid's electric resistance with time will provide a jump at each instant that the crack tip traverses the end of a strip, yielding the precise location of the crack tip at a number of discrete times. To ensure that the crack tip is not significantly ahead of the fracture of the strip, the strip's thickness must be at least an order of magnitude less than the crack face separation. The disadvantage of this method is that the discrete measurements can only provide a measure of the mean velocity *between* the strips. By extending the method to a continuous coating (instead of discrete strips), one has the advantage that the crack tip's location is obtained as quickly as the voltage drop across the coating can be digitized. The precision of the measurements is limited only by the background noise and the uniformity of the coating. It can be improved with an evaporated coating which provides precise velocity data near the sample faces. Thinness of the sample does not present a limitation and can, in fact, be taken advantage of if one wishes to correlate the instantaneous velocity with localized features formed on the fracture surface. This method has been used widely (see, for example, Brickstad and Nilsson, 1980; Fineberg *et al.*, 1991, 1992) with considerable success, using a variety of materials.

### 6.11.4.3  Ultrasonic Measurements

In this method (Kerkhof, 1973), which has been used both with glass and brittle polymers, a moving fracture is perturbed by an ultrasonic wave generated by a sample boundary in a direction orthogonal to that of crack propagation. The interaction of the sound with the crack tip causes the sound to be deflected periodically

as it traverses the sample, the trace of which is imprinted onto the resulting fracture surface. Since the temporal frequency of the modulation is that of the ultrasonic driving, measuring the distance between neighboring surface modulations provides a nearly continuous data set for the instantaneous velocity of the crack tip. The method's precision is limited only by the ultrasonic frequency used, which is typically in the MHz range, and also by the precision of the surface measurement. The disadvantage of this method, relative to the other techniques, is that it is a perturbative method, since the crack deflection is accomplished by altering the stress field at its tip, and hence externally-induced oscillations can potentially mask intrinsic, time-dependent effects.

### 6.11.5   Measurement of the Thermal Effects

A propagating fracture transforms the elastic energy stored in the material to either kinetic energy, the energy needed for breaking the atomic bonds, or to dissipated heat. The dissipated heat can be measured by two types of measurements. In one method one places small temperature sensors at a given distance from the path of a crack and measures the temperature rise in the material as a function of time after fracture has occurred. Since the time scale of fracture is orders of magnitude shorter than the typical times for thermal diffusion within the material, one can approximate the problem by assuming that an instantaneous planar heat source is created along the fracture plane, and that the radiative losses are negligible over the period of measurements. Then, the measured time-dependence of the temperature at a single point can be fitted to the solution of the heat conduction equation. Measurements of this sort were carried out in PMMA (Döll, 1973), in glass (Weichert and Schonert, 1974), and in steel (Zimmerman *et al.*, 1984). Moreover, it is possible to estimate the temperature rise in the vicinity of the crack tip by use of IR detectors (Fuller *et al.*, 1983; Zehnder and Rosakis, 1991; Kallivayalil and Zahnder, 1994), assuming that the emission spectrum of a crack corresponds to a black body spectrum, although this assumption may be suspect, at least in the immediate vicinity of the tip.

### 6.11.6   Measurement of Acoustic Emissions of Fractures

Measurements of acoustic emissions have long been used (see, for example, Scott, 1991) as a means of detecting either the onset of, or the precursors to, fracture, where the existence, the frequency of events and their locations can be measured. Although these techniques, due to their relatively limited precision, have not been used extensively in dynamic fracture experiments, they provide a sensitive method for determining whether changes in the stress field are taking place during fracture, because any rapid changes invariably release stress waves, and therefore can be used for detection of fracture and its onset. Such methods utilize arrays of resonant acoustic transducers since the advantage of their high sensitivity more than offsets the loss of information about the signal's spectral content. In fact, since the spectral content of the acoustic signal broadcast by a moving crack carries

important information (see Chapter 7 for theoretical discussion of this point), broadband transducers should be used together with relatively high amplification to offset the transducers' lack of sensitivity. The emissions are then correlated (Gross *et al.*, 1993; Boudet *et al.*, 1995, 1996) with velocity and fracture surface measurements. That this is a sensible method even when deflections of the 2D sample normal to their surface that are measured are due to the fact that the probe is sensitive to both longitudinal and shear waves due to mode conversion (Kolsky, 1953).

## 6.12   Oscillatory Fracture Patterns

One fundamental prediction of linear continuum fracture mechanics is that, as a crack propagates, its speed should increase until it reaches its asymptotic value, the Rayleigh sound speed $c_R$—the speed of sound on a free surface. However, experimental observations of fracture propagation in many heterogeneous materials indicate that, in the vast majority of cases, the ultimate velocity of a propagating crack is not more than about $0.5c_R$ (unless the material is strongly anisotropic; see Chapter 7). For example, oscillatory fracture patterns that have been observed in many materials strongly violate this fundamental prediction, and our goal in this section is to briefly describe and discuss these patterns and how they have been created in the laboratory.

These patterns were observed in the beautiful experiments of Yuse and Sano (1993). They imposed a temperature gradient along a thin glass plate, from a hot region to a cold one. A microcrack was introduced in the glass, and the glass was pushed. As the plate started to move the crack jumped ahead of the thermal gradient and stayed there. It was observed that if the plate moves slowly, the growing crack remains straight and stable. However, increasing the velocity to a critical value $v_c$ gives rise to a transition whereby the fracture path begins to oscillate and an instability appears. At still higher velocities crack branching appears; see Figure 6.6. Ronsin *et al.* (1995) also provided experimental data for brittle fracture propagation in thin glass strips, using a thermally-induced stress field. In their experiments the temperature field was controlled by the width $w$ of the plate, and induced thermal expansion in the sample. It was observed that for widths below a critical value $w_c$ no fracture was formed. For $w_c < w < w_o$, where $w_o$ is a second critical width for the onset of oscillatory cracks, straight fractures were formed and propagated with a constant speed. For $w > w_o$ oscillatory fractures were generated which became more irregular as $w$ was increased beyond $w_o$.

These predictions are in agreement with the results of several sets of spectacular experiments by Fineberg *et al.* (1991,1992) and Gross *et al.* (1993). Many earlier experiments had already reported several interesting features of dynamic crack propagation in materials (see, for example, Mecholsky, 1985). For example, it had been reported (see, for example, Döll, 1975; Kusy and Turner, 1977) that in some brittle materials, such as PMMA, the fracture pattern exhibits characteristic wavelength, that surface roughness increases with crack speed (see, for example,

FIGURE 6.6. Fracture pattern formation in the experiments of Yuse and Sano (1993).

Langford *et al.*, 1989, and references therein), and that periodic stress waves are emitted from the tip of the rapidly moving cracks in a wide variety of materials (see, for example, Rosakis and Zehnder, 1985; Dally *et al.*, 1985, and references therein). Fineberg *et al.* (1991,1992) carried out beautiful and precise experiments to study fracture propagation in brittle plastic PMMA and showed, that there is a critical velocity $v_c$ beyond which the velocity of crack tip begins to oscillate, the dynamics of the crack changes abruptly, and a periodic fracture pattern is formed. For $v > v_c$ the amplitude of the oscillations depends linearly on the mean velocity of the propagating crack. Thus, the dynamics of cracks is governed by a dynamical instability, and explains why the crack tip velocity does not attain the limiting Rayleigh velocity predicted by the linear elastic theory. Although Yoffe (1951) had already predicted the existence of a sort of dynamical instability in fracture, showing that a fracture that moves along a straight line will branch off if its speed becomes larger than a critical value, her predicted critical velocity was too large, and therefore the type of instability that was considered by her could not provide a complete explanation for Fineberg *et al.*'s experiments. The theoretical studies of such fracture patterns will be discussed in Chapter 7.

In another set of beautiful experiments, Gross *et al.* (1993) used two materials, the PMMA and soda-lime glass, to show that all features of dynamics of crack propagation in the two materials, such as acoustic emission, crack velocity, and surface structure, exhibit quantitative similarity with each other. Thus, there exists universal characteristics of fracture energy in most materials that are the result of energy dissipation in a dynamical instability. Perhaps the most spectacular experiments were carried out by Sharon *et al.* (1995) and Sharon and Fineberg (1996) using the brittle plastic PMMA. They identified the origin of the dynamical instability

during fracture propagation as being the nucleation and growth of the daughter cracks which limit the speed of the propagating crack tip. The daughter fracture carries away a fraction of the energy concentrated at the tip of the moving crack, thus lowering the velocity of the tip. After some time, the daughter crack stops growing, and thus the crack tip velocity increases, until a new daughter fracture starts to grow, and so on. They also observed that the branching angle for a longer daughter fracture was smaller than that of the shorter daughter fractures. Theoretical modeling and computer simulations of dynamic fracture that can reproduce these features will be described in detail in Chapters 7 and 8.

## 6.13   Mirror, Mist, and Hackle Pattern on a Fracture Surface

Studies of fracture surfaces of amorphous brittle materials indicate that they have a characteristic structure that is popularly referred to as *mirror, mist,* and *hackle*. This pattern has provided an important tool for studying a number of important fracture phenomena, and at the same time has raised a number of fundamental questions. Figure 6.7 presents the original pattern reported by Johnson and Holloway (1966), which is the fracture surface of an inorganic glass, soda-lime-silica glass rod with



FIGURE 6.7. Light microscope photograph of mirror, mist, and hackle regions on fracture surface of a 5 mm diameter soda-lime-silica glass rod, tested in uniaxial tension. The mirror region is roughly circular, surrounded by the narrow band of mist that gradually develops into the hackle (after Johnson and Holloway, 1966).

a diameter of 5 mm, tested in uniaxial tension. A crack nucleated at a small surface flaw that was generated by contact during handling, and then propagated normal to the tensile axis, i.e., under Mode I fracture. In the initial stages of the experiment, the crack growth led to a very smooth fracture surface, which is called mirror. The crossing of the rupture front with elastic waves can leave behind ripples in the mirror zone which are called *Wallner lines*.

This region is surrounded by a slightly rougher and less reflective region, which is referred to as the mist. It consists of fine striations that look like microscopic blades that are oblique to the crack plane. This zone appears when the velocity of the crack is about half of the velocity of transverse elastic waves. Finally, the mist region merges into a very rough fracture surface with irregularly oriented facets, which comprise the hackle region. The facets are separated by large steps that are aligned parallel to the main direction of crack propagation. As Figure 6.7 indicates, the transitions between the neighboring regions are not sharp; rather they represent progressive changes in the surface roughness.

Since the transitions from mirror to mist to hackle regions are not sharp but gradual and diffused, the answer to the question of where one region ends and another one starts cannot be precise. Johnson and Holloway (1966), who analyzed these regions for the first time, stated that, "The position assigned to the boundary between mirror and mist zones depends upon illumination and the magnification at which the fracture is examined, even within the range of the optical microscope. With an electron microscope mist can readily be resolved in the region seen as mirror under optical conditions." However, a better way of distinguishing between the three zones is by measuring the changes that occur in the surface roughness of the fractured material. While the height of the roughness remains essentially constant in the mirror region, it increases sharply and monotonically as the transition to the mist zone is made. Measurement of roughness of fracture surfaces and its significance will be discussed in the next section.

In the light of our discussions earlier in this chapter, it is not difficult to understand the development of the mirror, mist and hackle pattern. Suppose that the length of the initial flaw is $c$. In uniaxial tension, the stress concentration is large at the tip of the flaw. If the stress is large enough, the Griffith criterion is satisfied and the fracture begins to grow. If the loading condition is held constant, the increase in the fracture length implies fracture instability and the existence of excess energy that drives the fracture. Thus, the crack accelerates very rapidly, with which the rate of energy release also increases rapidly, resulting in higher stress intensities at the tip. The large stress intensity and rate of energy release also imply a corresponding increase in the micro-mechanical activity at the tip of the fracture, and hence a corresponding increase in the roughness of the fracture surface. Note that, depending on the test conditions, a fourth region of the fracture surface may also develop. This region would be the result of having the main fracture bifurcate into two or more branches. Normally, bifurcation occurs in high-stress failures.

The boundaries between the mirror, mist, and hackle regions are roughly circular, implying that the crack accelerates outward in all directions with essentially the same rate. Experiments have indicated that if $R$ is the radius of a boundary between

two zones, then the fracture strength $\sigma_f$ of the material, i.e., the stress at which the crack starts to move (see also Chapters 7 and 8) is related to $R$ through

$$\sigma_f \sqrt{R} = a, \quad R = R_{\text{mirror}}, R_{\text{mist}}, R_{\text{hackle}}, \tag{19}$$

where $a$ is a constant. Observe that Eq. (19) has the same form as the Griffith condition, Eq. (9) [if we rewrite Eq. (9) as, $\sigma_c \sqrt{c} \propto \sqrt{Y'\Gamma}$ ], and therefore the constant $a$ is related to the quantity $\sqrt{Y'\Gamma}$ that appears in the Griffith condition. Moreover, in view of Eq. (16), the constant $a$ can also be related to the fracture toughness $K_c$. Experiments have also indicated that if $R_0$ is the radius of the initial flaw at which the crack nucleates, then the radius $R_{\text{mirror}}$ of the mirror zone is related to $R_0$ through, $R_{\text{mirror}}/R_0 \simeq 10$. Clearly, the circular boundaries between the three zones will not develop if the crack cannot accelerate in all directions with the same rate. The deviation from circularity depends partly on the boundary conditions used in the test. For example, a material in a bending experiment develops a stress distribution that is quite different from one that it experiences in a uniaxial tension experiment. Moreover, the mechanism of crack growth in amorphous materials is different from that of crystalline materials, so that the shape of the boundaries between the mirror, mist and hackle zones also depends on the material.

Before closing this section, let us point out that in the fracture literature one often finds references to *twist hackle* and *stress* or *velocity hackle*. The former refers to a rough surface that is generated by a Mode I/III fracture experiment, whereas the latter is the result of a crack propagating at very high speeds or under a large stress. The phrase mirror has also been used occasionally for describing the initial stage of the development of a fracture surface, whereas careful examination of the surface would reveal that it is too rough to be classified as the mirror zone. To make the distinction between a mirror zone and a rougher region, one may define the mirror region as the zone in which the average height of the roughness is less that the wavelength of light.

## 6.14    Roughness of Fracture Surfaces

The development of mirror, mist and hackle zones makes it clear that, as a crack propagates, the fracture surface develops roughness, the intensity of which increases with the extent of the crack propagation, which in turn depends on the loading condition, and the shape, morphology and composition of the material. Therefore, measurement of the roughness of a fracture surface may provide insight into dynamics of fracture propagation in a material. However, because of the dearth of comprehensive experimental data, i.e., data sets that contain *simultaneous* measurements of the roughness, the (dynamic) stress intensity factor $K_d$ and the speed $v$ of fracture propagation, the relation between the three quantities is not clear at present, and is the subject of ongoing research by many groups around the world. Arakawa and Takahashi (1991, where references to their earlier work in the Japanese literature can also be found) carried out one such set of measure-

FIGURE 6.8. Dynamic stress intensity factor $K_d$, the fracture velocity $v$, and the roughness of the fracture surface versus the fracture length for a brittle epoxy resin. $K_{Ic}$ is the critical value of the stress intensity factor $K_I$. The data are from Arakawa and Takahashi (1991) (after Hull, 1999).

ments which is summarized in Figure 6.8. In their experiments, they used 6 mm thick plates of various transparent plastics, including a thermosetting epoxy and a thermoplastic PMMA, and measured the velocity of the propagating crack, the dynamic stress intensity factor $K_d$, and the roughness $w$ of the surface. There seem to be general correlations between the crack speed and the stress intensity factor on one hand, and the roughness of the fracture surface on the other hand. At the same time, another feature of this figure indicates that there may not be a unique relation between the crack speed and the intensity factor, since the two quantities have not reached their maximum at the same point, whereas the maxima of the intensity factor and surface roughness seem to happen at the same crack length, and therefore these two quantities are probably better correlated than $K_d$ and $v$.

However, this is *not* a completely universal rule. Under certain circumstances, the surface may become *smoother* as fracture propagation proceeds. An example is provided by elastomers, where in some range of the crack speed their fracture surface is rough at low speeds, while it is smooth and mirror-like at high speeds. Thus, the increase or decrease in $K_d$ and $v$ is not directly linked to the roughness of the surface. Moreover, it must be mentioned that many materials do not develop mirror smooth surfaces at all. For example, if sharp pre-existing cracks are not

present on the surface, or are blunted by deformation, the mirror surface will not develop. In addition, the presence of grain boundaries, multiple phases of the material, and reinforcing particles force the crack paths into irregular shapes.

Systematic investigation of roughness of fracture surfaces and their scaling properties were first undertaken by Mandelbrot $et\ al.$ (1984), although Passoja and Amborski (1978) and Chermant and Coster (1979) had already suggested that fracture surface of metals may have fractal and scale-invariant properties. As discussed in Chapter 1, if the width $w$ of a rough surface follows the scaling law (1.34), then the surface is a self-affine fractal with a fractal dimension $D_f$ which, in $d$ dimensions, is given by

$$D_f = d - \alpha, \tag{20}$$

where $\alpha$ is the roughness exponent, which is usually the same as the Hurst exponent $H$ introduced and discussed in Chapter 1, although, theoretically, the two exponents can be different. Mandelbrot $et\ al.$ (1984) studied fracture surface of steel and concluded that the surface possessed fractal morphology. They estimated the fractal dimension of the fracture surface of their material to be $D_f \simeq 1.28$, implying a roughness exponent $\alpha \simeq 0.72$. If we assume that the roughness exponent $\alpha$ is equivalent to the Hurst exponent $H$ for the fractional Brownian motion described in Section 1.4, a roughness exponent of 0.72 implies long-range positive correlations on the fracture surface. Indeed, the profiles of such fracture surfaces are very similar to fBm with a Hurst exponent $H > 0.5$ (see Figure 1.2). Since the original work of Mandelbrot $et\ al.$ (1984), many other measurements of fractal and self-affine properties of fracture surface of a wide variety materials have been reported. In particular, several experimental techniques have been used for measuring and characterizing the roughness of fracture surfaces and estimating its roughness exponent, which we now describe and discuss.

### 6.14.1  Measurement of Roughness of Fracture Surface

Underwood and Banerji (1986) measured fractal dimension of fracture surface of AISI 4340 steels over the temperature range of 200 to 7000°C, and found that the lowest value of $D_f$ is at 500°C, generally believed to correspond to temper brittleness. Pande $et\ al.$ (1987) disputed the accuracy of Mandelbrot $et\ al.$'s result, and measured the apparent fractal dimension of fracture surfaces of titanium alloys. Fractal dimensions of about 1.2 were obtained, implying a roughness exponent $\alpha \simeq 0.8$. This value is, however, in agreement with many other measurements on a wide variety of materials discussed below, and with the Molecular Dynamics simulation results described in Chapter 9, and thus it does not cast doubt on the measurements of Mandelbrot $et\ al.$ (1984). Wang $et\ al.$ (1988) investigated the relationship between the fractal dimension of a fracture surface and its fatigue threshold using dual-phase steel, and found roughly a linear relation between the two. Mu and Lung (1988) measured the fractal dimension $D_f$ of fracture surface of $24SiMnCrNi_2Mo$ and $30CrMnSiNI_2A$ steels under plane strain. A linear relationship was found between the fractal dimension of fracture surface of these

metals and their fracture toughness, such that $D_f$ decreased smoothly as the fracture toughness increased. These issues and the progress up to 1988 were reviewed by Williford (1988).

Mecholsky *et al.* (1988,1989) and Passoja (1988) studied fracture surfaces of many solid materials, including several different aluminum and five glass ceramics, all of which had distinct microstructures. They found that as the toughness of the materials increases, so does also the roughness of the fracture surface. The fractal dimension $D_f$ was found to be in the range $1.15 - 1.30$, with an average of about 1.22, implying an average roughness exponent $\alpha \simeq 0.78$. They also investigated the relation between fracture energy and the geometry of fracture surface in many different brittle materials and proposed the following equation

$$\Gamma = \frac{1}{2} Y \xi (D_f - 1), \tag{21}$$

where $\Gamma$ is the fracture energy, $Y$ is an elastic modulus, and $\xi$ is a characteristic length scale of the material.

Dauskardt *et al.* (1990) undertook a systematic study of five samples of brittle and ductile transgranular cleavage, intergranular fracture, microvoid coalescence, quasi-cleavage, and intergranular microvoid coalescence in various steels. These materials were fractured both at room temperature and also a very low temperature. They analyzed the measured length $L$ of the surface versus the measuring step length $L_s$ which are related through, $L \sim L_s^{1-D_f}$. In many cases, a fractal dimension $D_f \simeq 1.2$ was obtained, in agreement with the previous estimates discussed earlier. However, in several other cases the relation between $L$ and $L_s$ was more complex. Bouchaud *et al.* (1990) studied fracture of an aluminum alloy in 4 different heat treatment regimes. The fracture surface was elecro-coated with nickel, then polished and digitized. The correlation function $C(r)$, Eqs. (1.5)–(1.8), was then constructed for the aluminum-nickel boundary for a large number of samples. Even though quite different mechanisms of fracture were dominant in these materials, in all cases the roughness exponent was $\alpha \simeq 0.8$.

Zhenyi *et al.* (1990) and Dickinson (1991) studied fracture surface of polymers and ceramics, measuring both surface roughness and light emission signals. Fractal dimensions of 1.2–1.3 were measured for the rough surfaces, resulting in roughness exponents of about 0.7–0.8. The photon emission signals also had fractal characteristics, and measurement of their fractal dimensions yielded values between 1.24 to 1.42, implying roughness exponents in the range 0.6–0.75. Note that, there appears to be a close relationship between the fractal dimensions of the fracture surface and that of the emission signals. If the exact nature of this relationship can be identified, then photon emission signals may provide an accurate probe of fracture surfaces and their morphology.

Fractures on carbon surfaces were analyzed by Miller and Reifenberger (1992), who reported that $\alpha \simeq 0.75$. Poon *et al.* (1992) studied fracture surface of natural rock, such as sandstone, limestone, and carbonates. For each sample roughness profiles of several thousand points were constructed, and for all cases studied a roughness exponent of about 0.8 was obtained. Måløy *et al.* (1992) investi-

gated fracture surfaces of six different brittle materials, ranging from Al-Si alloy AA4253 to porcelain. The materials were notched and then fractured at the temperature at which nitrogen becomes liquid. Many profiles of the rough fracture surfaces were then obtained and analyzed. Two methods of analysis, including the power-spectrum method described in Section 1.4.1, were used. The roughness exponent was estimated to be $\alpha \simeq 0.87 \pm 0.07$ for all the six samples. Baran *et al.* (1992) analyzed fracture surface of several brittle materials, including glass and dental porcelain, and reported large roughness exponents, ranging from 0.65 to 0.93. Poirier *et al.* (1992) studied deformation of regular packings of equal parallel cylinders. The local stress-strain characteristics, at the contact between the cylinders, exhibited a softening part which localized the deformation. The deformation band was rough with a roughness exponent $\alpha \simeq 0.73 \pm 0.07$.

An interesting method for studying fracture surface was developed by Imre *et al.* (1992) who determined the fractal dimension of the surface electrochemically by measuring the diffusion current, also called Cottrell current, at a gold replica of the fractured metal electrode. (It is interesting to find research groups that are rich enough to afford gold in their investigations, while others starve for research funds!) The replicas were prepared by pressing gold wafers into the fractured steel surfaces in a hydraulic press at high pressure. The gold surfaces were then cleaned, and the gold electrodes were immersed in an aqueous electrolyte with a calomel reference electrode. The potential was switched from 0 V to 650 mV for a short period of time, and then was switched back to 0 V. According to Nyikos and Pajkossy (1985) the current $I(t)$ should scale with the time $t$ as

$$I(t) \sim t^{(\alpha-2)/2}, \tag{22}$$

so that simple measurements of $I(t)$ versus $t$ should yield $\alpha$ (and hence $D_f$). Roughness exponents of about 0.8 were measured by this method.

Another interesting method for measuring roughness properties of a fracture surface was developed by Friel and Pande (1993). In their method pairs of electron micrograph images of fracture surface of titanium 6211 at two different inclination angles ($30°$ and $36°$) were constructed using a scanning electron microscope (SEM). The surfaces were fractured under tension. The SEM images were obtained under various magnifications, ranging from 50 to 10,000. The surface fractal dimension was then estimated by measuring the surface area as a function of the length scale (or measurement resolution), and was found to be about 2.22, implying a roughness exponent $\alpha = 3 - 2.22 = 0.78$. Schmittbuhl *et al.* (1993) measured roughness exponent of several granitic faults and found $\alpha \simeq 0.85$, close to the values obtained by others for various materials. E. Bouchaud *et al.* (1993b) analyzed the statistics of fracture surfaces of polycrystalline intermediate compound Ni$_3$Al. Such fracture surfaces also contain secondary branches, as opposed to most of the fracture surfaces discussed above which had no side branches. Despite this, E. Bouchaud *et al.* (1993) could define a roughness exponent for fracture surface of these materials, and their measurements indicated that $\alpha \simeq 0.8$. Lemaire *et al.* (1993) put a viscoelastic paste made of sand and resin between two plates which were driven away from each other at a given velocity until the paste broke.

Five different velocities were used, and after fracture the hardened paste was sliced parallel to the tensile direction. The fractal dimension of the profiles was then determined by two methods, the standard box-counting method, and by the power-spectrum methods, both of which were described in Chapter 1. A roughness exponent $\alpha \simeq 0.88 \pm 0.05$ was measured which was independent of the velocity.

Daguier *et al.* (1995) studied the morphology of fractures in two different metallic alloys. The fractures had been stopped during their propagation by pinning microstructural obstacles to the surface. One of the alloys was the 8090-Al-Li which is very anisotropic, for which the roughness exponent was found to be $\alpha \simeq 0.6 \pm 0.04$. The other alloy was Super $\alpha_2$ Ti$_3$Al with a 3D fatigue fracture for which $\alpha \simeq 0.54 \pm 0.03$. Daguier *et al.* (1996) used atomic force microscopy and SEM methods to study fracture surface of Ti$_3$Al-based alloys. They found that at large length scales, and over *several decades* in length scales, the roughness exponent was $\alpha \simeq 0.8$, whereas at much shorter length scales the roughness exponent was close to 0.5. Daguier *et al.* (1997) also studied fracture surface of a silicate glass as a function of the fracture velocity. At large length scales the roughness exponent was $\alpha \simeq 0.78$, whereas at smaller length scales $\alpha \simeq 0.5$. The crossover length scale $\xi_{co}$ that separated the two scaling regimes was shown to be proportional to the inverse of the fracture velocity. If $h_{max}$ is the difference between the maximum and minimum heights $h$ within a given window on the surface, then the two scaling regimes could be combined into a single scaling law

$$h_{max} \sim r^{0.5}\Psi(r/\xi_{co}), \tag{23}$$

where $\Psi$ is a scaling function with the properties that $\Psi(x) \sim 1$ as $x \to 0$, and $\Psi(x) \sim x^{0.28}$ for $x \gg 1$.

Thus, summarizing all the experimental data discussed so far, it appears that at large enough length scales a roughness exponent $\alpha \simeq 0.8$ represents a universal value, regardless of the material or even the mechanism of fracture. The possibility of universality of $\alpha$ was first pointed out by Bouchaud *et al.* (1990). We should, however, point out that if a fracture surface is analyzed on relatively short length scales, then the *effective value* of $\alpha$ may be smaller than 0.8. For example, Mitchell and Bonnell (1990) analyzed fracture surface of fatigued polycrystalline copper and reported that $\alpha \simeq 0.65$, while for a single crystal silicon $\alpha \simeq 0.7$ was obtained. Metallic materials, the roughness exponents of which have been determined through scanning tunneling microscopy, usually operate in the nanometer range and have $\alpha < 0.8$. For example, Milman *et al.* (1993, 1994) reported a roughness exponent of about 0.6 for fractured tungstene, and close to 0.5 for graphite. Low cycle fatigue experiments on steel samples on micrometer scales yielded a roughness exponent close to 0.6 (McAnulty *et al.*, 1992). Low values of the roughness exponents are interesting because they might be explained based on models of minimum energy surfaces in disordered environments. Such concepts were first discussed by Chudnovsky and Kunin (1987), Kardar (1990), Roux and Francois (1991), and Ertas and Kardar (1992,1993,1994,1996). For example, Roux and Francois (1991) argued that the path that is selected by a propagating fracture should be such that the overall fracture energy is minimized. Their simulations un-

der such a condition led to a roughness exponent in the range 0.4–0.5. The apparent length-scale dependence of the roughness exponent $\alpha$ may also be explained in another way based on the velocity of fracture propagation, and whether one is in the regime of quasi-static or rapid fracture (Bouchaud and Navéos, 1995). This distinction, its theoretical treatment, and the corresponding roughness exponents will be described in Chapter 7, where we will also discuss the implication of the self-affine structure of fracture surface for crack propagation.

Now that there is little doubt that fracture surface of a wide variety of materials is rough with well-defined characteristics, let us briefly describe how such surfaces are studied experimentally. This subject has been discussed in detail by Hull (1999), and what follows is a summary of his discussion. Roughness is typically characterized by measuring the height $h$ of the roughness profile. A "primitive" method would be based on using a raster scan of parallel traverses across the surface using a stylus which traverses parallel to the $x$-axis—the axis that is parallel to the mean position of the roughness profile—and measures the height. The stylus is typically a fine, diamond-tipped needle which is in contact with the surface by a small external load. The height of the needle is measured using a transducer. The disadvantage of this method is that the stylus may damage the surface, and hence create traces that do not belong to the original fracture surface.

Atomic force microscope can also be used which has a highly fine silicon nitride stylus with a tip radius of about 20–30 nm. The probe is held at a fixed position from the base of the rough surface, and the surface itself is moved parallel to this base. The height of the probe is measured from the reflection of light from mirror on the stylus beam. A powerful feature of this method is that it can determine roughness parameters on specific sections of the roughness profile.

In a modern version of the stylus technique, the mechanical stylus is replaced by a fine laser beam that is held at a constant distant from a references surface. The size of the spot is typically 1 $\mu$m in diameter, and the rough surface traverses under the beam light. The surface shape is then determined from the change in the length of the light's path that is reflected from the surface.

## 6.14.2   Mechanisms of Surface Roughness Generation

There are at least three main mechanisms that give rise to a rough fracture surface. What follows is a brief description of each mechanism.

### 6.14.2.1   Growth of Microcracks

In thermoplastic polymers (as well as other materials) the high stresses around the main crack cause micro-cracking in the material ahead of the main fracture. These smaller cracks grow and eventually become connected to each other and to the main crack. As the stress intensity increases, there are corresponding increases in the size of the damage zone and the out-of-the plane crack nucleation. The net result is a rough fracture surface. Natural materials, particularly rock, exhibit intense micro-cracking and surface roughness (see Sahimi, 1993b, 1995b, for

detailed discussions), with the scale of their roughness being equal to at least the scale of the microstructure. We will come back to this mechanism in Chapter 7, where we discuss the relation between micro-cracking and dynamics of fracture propagation.

### 6.14.2.2    Plastic Deformation

If plastic deformation occurs ahead of the tip of the main growing crack, crack growth takes place in a zone of deformed material. If the deformation zone is not homogeneous, the crack path is deflected out of the plane in which it is propagating, leading to surface roughness. The interaction between deformation processes and the growing crack depends on the dynamics of growth of the deformation zones and cracks, which in turn depends on the stress field in the material, and the stress level at which these phenomena are activated.

### 6.14.2.3    Macroscopic Branching and Bifurcation

Roughness of fracture surface in isotropic, homogeneous, amorphous and brittle materials, such as inorganic glasses, might be the result of local changes in the path of the growing crack. These changes are the result of local instabilities at the tip of the growing crack. The nature of these instabilities will be discussed in detail in Chapters 7 and 8. For now it suffices to say that micro-cracks are formed ahead of, and interact with, the main crack, the nucleation of which can be explained based on the Griffith criterion. Due to the high stresses that are distributed around the main growing crack, the micro-cracks are deflected out of the plane of the main crack by micro-branching or micro-bifurcation, hence giving rise to roughness in the fracture surface.

However, the growth of micro-cracks ahead of the main crack in brittle glasses has been disputed by some researchers, who argue that in such materials the stress to activate very small flaws and grow them into micro-cracks approaches the theoretical strength of the material, in which case only the main crack grows by breaking the interatomic bonds. It has been suggested instead that local tilting of the crack out of its main plane is the cause of micro-branching. These tilted cracks grow a short distance, but their size increases with the dynamic stress intensity factor $K_d$ and the crack velocity $v$. When the dimensions of the tilted cracks become comparable to the dimensions of the test sample, macroscopic bifurcation takes place. Experimental evidence for this mechanism was reported by Johnson and Holloway (1968) and Kulawansa *et al.* (1993).

## 6.15    Cleavage of Crystalline Materials

The discussions so far are mostly relevant to brittle fracture of amorphous materials. Another important subject is cleavage of crystalline materials. Single crystals are homogeneous, but they also contain a degree of anisotropy which assists their cleavage. To understand this phenomenon, not only does one need information on

the effective properties of the material, such as their elastic moduli and fracture toughness, but also an understanding of such micro-deformation processes as slip that usually precedes and accompanies fracture in crystalline materials. The degree of symmetry that the crystalline material exhibits also plays an important role, because the strength of the anisotropy of micro-deformation processes depends on such symmetries. The most important effect of anisotropy is that cleavage may occur parallel to planes in a crystal that are not normal to the maximum tensile stress. This is particularly true in crystalline materials that exhibit a low degree of symmetry, such as mica in which cleavage is only in a single set of planes. In addition, temperature and strain rates also play important roles by influencing the mobility of dislocations.

The low surface energy of crystallographic planes, which in turn depends on the strength of the interatomic bonds, is the main cause of cleavage in crystals. If cleavage occurs along a single plane, it would produce a featureless surface. However, often one observes well-defined and crystallographically oriented features on the fracture surface of a crystalline material. These features are usually caused by the generation and presence of dislocations that interact with the propagating fracture. In metals with body-centered cubic symmetry, such as chromium, tungsten, and iron, the main cleavage occurs on {001} planes, of which there are three, (001), (010), and (100). If a cyrstal is tested in an arbitrary direction, the {001} plane with the largest tensile stress normal to the plane is the most likely place for cleavage. If a crystal is tested in tension parallel to [011], the (001) and (010) planes have the same resolved normal tensile stress. In this orientation the stress on the (011) plane is much greater than on the {100} planes. Thus, fracture may occur either on an (011) plane, or along the two equally stressed {001} planes.

On the other hand, crystals with the zinc-blende structure, such as gallium arsenide, can be described as a cubic unit cell that consists of two interpenetrating FCC lattices of the two elements (Ga and As). The center of one lattice is at the position (1/4,1/4,1/4) of the other. These materials are of great industrial importance because of their use in producing semi-conductors. They cleave on {001} planes, of which there are three equivalent pairs of orthogonal planes. Slip is restricted to {111} planes. Such materials usually exhibit strong brittleness.

If polished (001) faces of GaAs crystals are coated with an epitaxial layer of GaAs that contains a small amount of carbon, tensile stresses are generated in the surface layers. These stresses then lead to the formation of very fine, atomically sharp surface cracks (see, for example, Murray *et al.*, 1996). The cracks form on two orthogonal {011} planes that intersect the (001) surface at right angles, remain sharp, and grow at very low stresses. The fracture surface is mirror smooth and flat. However, if GaAs crystals are tested in complex loading conditions, the fracture surface becomes very rough.

Layered materials usually have very strong bonding within the layers and weak bonding between the layers. An example is muscovite mica that consists of an ordered stack of double layers, about 2 nm thick, of strongly bonded planar arrays of silica tetrahedra held together by Coulomb attraction caused by the potassium ions between the layers. In such materials cleavage occurs between the weakly-bonded

layers, and may also occur through the center of the double layer. Deformation is restricted to the sliding of layers over each other. The reader should consult Hull (1999) for extensive discussions of other crystalline materials.

So far we have discussed the cleavage of single crystals. In practice, cleavage of polycrystalline materials, such as ceramics and rock, is also very important. Let us briefly discuss these phenomena. We assume that the bonding between the crystals is very strong, and that the grain boundary interface does not experience failure.

In polycrystalline materials, each grain is surrounded by many other grains of different orientations. Therefore, such materials fracture by successive nucleation and propagation of several cleavage cracks across the boundaries between neighboring crystals. There is a change in the orientation at the grain boundary. If the angle between the neighboring grains is small, the cleavage crack in a crystal can propagate across the boundary between the neighboring crystals, in which case the cleavage plane is tilted and twisted. However, if the orientations of the crystal grains are very different, the propagation of cleavage from one crystal to another depends on the relative orientation of the cleavage planes in the crystals.

Consider, for example, two adjacent grains with a common boundary between them, and suppose that a crack in one of the grains reaches the boundary. Then, it may stop there with no further crack propagation. Alternatively, the crack may stop at the boundary, but the high stress at its tip may help nucleate another crack in the adjacent grain with a different orientation. The two cracks have a common point at the boundary. The third possibility is having a cleavage plane in the second grain that is tilted relative to the cleavage plane in the first grain, in which case the crack propagates continuously across the boundary. Therefore, fracture propagation in polycrystalline materials depends critically on the distribution of their grains or single crystals. Even if an array of grains is distributed randomly, the local direction of crack propagation depends on the relative orientations of the grains at the crack tip. On the scale of the single crystal size, the main crack path is not straight. It is also possible that local regions of the crack "tunnel" ahead of the main crack front because of the existence of a path of favorably oriented single crystals in the region. If a polycrystalline material contains preferred orientations, then crack growth in it is easier in some directions than others.

## 6.16   Fracture Properties of Materials

Let us now describe and discuss important fracture properties of several classes of materials. In general, one may divide most materials into three distinct classes which are polymeric materials, metals, and rock-like materials which include concrete, rock, glass, and ceramics. We already described fracture properties of glass when we discussed the mirror, mist and hackle patterns. We do not consider concretes here, and fracture properties of natural rock have been described in detail elsewhere (Sahimi, 1993b, 1995b). What follows is a brief summary of the properties of the remaining important materials. Our discussion is not, and cannot be,

exhaustive, as the mechanical properties of each of these materials are subjects of separate books.

### 6.16.1   Polymeric Materials

We described in Chapter 9 of Volume I many important properties of polymeric materials, and therefore the discussion in this section must be considered as complementary to what was presented there. Since we already discussed the difference between brittle fracture of amorphous materials and cleavage of crystalline materials, it is important to understand to what extent a polymeric material can be crystalline. Although homopolymers are crystalline, due to the length of the chains in their structure, polymeric materials do not usually have a completely crystalline structure. Instead, they usually consist of a mixture of crystalline and amorphous regions. On the other hand, many industrial polymers, such as PMMA, are completely amorphous, as already mentioned above. Moreover, generally speaking, random copolymers and cross-linked polymers are also amorphous.

To discuss mechanical and fracture properties of polymers, we consider amorphous polymers below the glass transition temperature $T_g$ and crystalline polymers below the melting temperature $T_m$. Figure 6.9 shows typical stress-strain curves (in tension) for polymeric materials. The top curve represents brittle behavior. The tensile strain is typically about 1–5%. The middle curve exhibits a yield point and represents ductile fracture. The lowest curve indicates that the yield point is followed by a strain softening region in which the stress reaches a minimum, beyond which one has stress hardening which then leads to brittle fracture. The yield point $\sigma_y$ defines the onset of irreversible plastic deformation, and is proportional to the maximum of the true stress in a compression test. Its value depends, of course, on the composition of the material and the stress configuration. It increases



FIGURE 6.9. A typical stress-strain diagram for polymers. The top curve corresponds to brittle behavior, while the middle curve leads to ductile behavior. In the lowest curve, strain hardening leads to brittle behavior.

logarithmically with the strain rate, and slowly decreases with increasing temperature, eventually vanishing at $T_g$.

Between the elastic limit and the yield point, many polymers that are under tension exhibit a series of crazes that are normal to the tensile stress. Both amorphous and crystalline polymers generate crazes with the same features. In particular, it is easy to see crazes in amorphous polymers as they strongly scatter visible light. The inside of a polymer craze is typically filled with polymer fibrils, as a result of which the effective moduli of the material after crazing is only slightly smaller than before, implying that the onset of crazing cannot be detected on the stress-strain diagram.

Under tension or compression, polymeric materials can also develop shear bands, i.e., zones of highly localized shear. The bands are diffuse at high temperatures or low strain rates, but are localized at lower temperatures or higher strain rates. If the diffuse bands are further deformed, it will lead to ductile fracture, whereas deformation of localized shear bands leads to brittle fracture. If two shear bands intersect, it usually leads to a craze. The stress at the craze tip can also lead to the formation of shear bands.

## 6.16.2   Ceramics

The British Ceramic Society defines ceramic materials as, "All solid manufactured materials or products that are chemically inorganic, except for metals and their alloys, and which are usually rendered serviceable through high temperature processing." Ceramic materials include borides, carbides, halides, nitrides, oxides, and cermets, which are ceramic metals. They usually have a crystalline structure, but can also be found in amorphous form. The interatomic bonds in ceramics may be ionic, covalent, metallic, and van der Waals. It is clear how the first two types of bonds may form in ceramics. Metal transition carbides have bonds which have a metallic characteristic in that, valence electrons are freely shared by all the atoms in the structure.

Relative to metals, ceramics have large elastic moduli, ranging from 70 to 400 GPa. The moduli decrease very slowly with increasing temperature. They also have a large cohesive strength which is due to the fact that their interatomic bonds require high energies to be broken. However, as discussed earlier in this chapter, the presence of defects, which results in stress concentration, reduces the actual strength of these materials. In fact, the fracture strength $\sigma_f$ of ceramics is very sensitive to the presence of defects, the porosity, the shape and size of the grains, as wells as the pore-crack combination. Most importantly, $\sigma_f$ depends on the size of the defects, for which there is a critical size that, at a given stress, leads to fracture. For ceramics this size can be as large as a single crystal. The Weibull distribution [see Eq. (5.37); see also Chapter 8] usually describes well the statistical distribution of the fracture strength of ceramics. Moreover, if the defects are uniformly distributed in the material, the probability of having the critical condition in the material for fracture is relatively large. Experiments have indicated that in many ceramics, especially those that have a secondary phase, the crack velocity $v$ is related to the stress intensity factor $K_I$ by, $v \sim K_I^n$, where $n$ is a constant.

Experiments by Buresch *et al.* (1983) and others have also shown that the fracture strength of certain ceramics depends on the critical value $\sigma_n$ of the notch fracture stress, and also on the size of the cohesive zone (see above). The cohesive zone in ceramics is somewhat similar to the plastic zone in metals in that, the microcracked zone in the immediate vicinity of a crack tip causes the nonlinear behavior of ceramics. In this zone, there is a constant stress $\sigma_n$ for breaking either the grain boundaries or the crystal themselves, which depends on the cohesive stress $\sigma_c$ (see above). If the average stress in the cohesive zone reaches $\sigma_n$, instability occurs in the material and the main crack propagates.

The behavior of the fracture strength $\sigma_f$ of ceramics with variations in the temperature can be divided into two groups. In one group, $\sigma_f$ decreases monotonically with increasing temperature. Nitrides typically exhibit this behavior. In the second group, the fracture strength either stays constant with increasing temperature, or first experiences a small increase and then decreases. Ceramics that do not have a secondary phase at their grain boundaries exhibit this behavior.

### 6.16.3  Metals

Most metals have simple crystalline structures in the form of BCC or FCC lattices or a hexagonal close-packed (HCP). At the atomic scale the interatomic bonds break either along crystallographic plane in Mode I fracture (i.e., in a direction normal to the plane), or in Modes II and III fracture (i.e., in a direction parallel to that plane), which is also the mechanism for cleavage fracture already described above. Alternatively, metals fracture at high temperatures by coalescence of cavities. Single crystals of a HCP metal (for example, zinc) can slip on a single plane until the two parts completely separate. Usually, however, multiple slip occurs in single crystals which generates a neck in the material which is under tension. These necks usually initiate at inclusions which do not deform in the same way as the metallic matrix. In polycrystalline metals, necking occurs in a more diffuse fashion, but can lead to the complete separation of the two halves of the material when the neck's cross section vanishes. This mechanism is, however, rare. In most cases, the material breaks much sooner by developing a crack in the middle of the neck which is perpendicular to the tensile axis, which at the end tilts to a 45° orientation. This crack is the result of coalescence of vacancies which grow due to plastic deformation and elongate in the direction of the maximum principal strain.

This mechanism of fracture in metals is ductile because it involves large local slip deformation. It also often corresponds to a large macroscopic plastic deformation. However, coalescence of the vacancies does not always need large macroscopic deformation, such as when the volume fraction of the inclusions in a metal-matrix composite is large. Hydrostatic pressure prevents the growth of the cavities, whereas a tensile positive hydrostatic stress increases it, and thus reduces greatly the fracture strain.

Another mechanism of fracture of metals is intergranular cracking, which happens when the grain boundaries are weaker than their interior. The weakness is caused by impurities that have accumulated at the grain boundaries. In such a situation, the cracks preferentially follow the grain boundaries, leading to intergranular

fracture, sometimes referred to as *dimple fracture*. If the temperature of the system is high enough, then the vacancies migrate by diffusion, and then coalesce to create cavities and ultimately cracks. This mechanism is called *creep cracking*.

Cyclic straining of metals also results in fatigue fracture of metals which usually starts on the surface, and is generated by irreversible localized shear deformations (see above). The surface gradually develops roughness which, if strong enough, develops into a crack which then penetrates into the material along the shear direction.

### 6.16.4  Fiber-Reinforced Composites

These materials exhibit a wide variety of fracture modes, including rupture of individual fibers, interfacial debonding, matrix cracking and delamination. Various experiments involving X-ray radiography and optical and scanning microscopy indicate that if a unidirectional fiber-reinforced composite is loaded in the longitudinal direction (parallel to the fibers), the fracture process consists of four main stages.

(1) At less than 50% of the final load, individual fibers break at random.
(2) As the broken fibers accumulate, they join and form macroscopic cracks throughout the material.
(3) Delamination begins parallel to the fibers, starting at the large cracks.
(4) Delamination propagates parallel to the direction of the fibers.

If the composite material is subjected to a static tensile load in the longitudinal direction, the breaking of a fiber generates tensile stress concentration in the first unbroken fiber, which may lead to their breaking. In addition, shear stress concentration is generated at the interface between the broken fiber and the matrix which contributes to shear debonding along the fiber surface. Thus, breaking of fibers induces two types of fracture modes that proceed simultaneously. The volume fraction of the fibers and their orientations control which of the two modes is the dominating one. If the fibers are distributed closely, the fracture propagates from fiber to fiber, whereas when they are relatively far apart, the failure process proceeds along individual fiber surfaces in the shear fracture mode. Fiber misalignment, or *fiber waviness*, also influences the tensile strength of the composite. In fact, the broader the distribution of fiber misalignment, the smaller is the tensile strength of the composite materials.

### 6.16.5  Metal-Matrix Composites

The fracture strength of metallic materials can be improved by inserting into them short fibers or particles. A typical failure process in such materials involves,

(1) failure of the interface between the fibers and the matrix at the tip of the fiber;
(2) growth of a cavity within the matrix, beginning at the fiber tip;

(3) coalescence of the cavities due to plastic deformation and formation of a crack, and

(4) propagation of the crack.

Inside the metallic matrix the failure is ductile, but it appears brittle at the macroscopic length scales.

An important factor is the aspect ratio of the fibers, i.e., the ratio of their lengths and diameters. For example, fibers reinforce a material better than spherical particles. The larger the aspect ratio, the higher is the fracture strength of the composite. However, if the fibers become too long, they will no longer influence the strength of the material. The properties of the interface between the fibers and the matrix also have a very strong effect on the strength of the composite. If the interface is stronger, the composite material will have a lower ductility and a higher fracture strength. During production of the composite, internal stresses may be produced by mismatch between the thermal expansion coefficients of the matrix and the fibers. Thus, when the temperature of the system is reduced, residual stresses are produced in the composite which, however, disappear at high enough temperatures. In addition, one may have chemical reaction at the interface. If, for example, the fibers are oxidized, the fracture strength of the composite will reduce.

## Summary

The aim of this chapter was to define the basic concepts of fracture mechanics, and describe and discuss the basic phenomena that occur during fracture of materials. These concepts will be utilized in the next few chapters where we describe and discuss modeling of brittle fracture of heterogeneous materials and its transition to ductile behavior. We also described the experimental techniques that are used for measuring important characteristics of fracture of materials, such as the speed of crack propagation, and measurement and analysis of roughness of a fracture surface.

# 7
# Brittle Fracture: The Continuum Approach

## 7.0   Introduction

As discussed in the last chapter, fracture of brittle amorphous materials is a difficult problem, because the way a large piece of a material breaks is closely related to details of cohesion at microscopic length scales. For this reason alone, description of brittle fracture of materials has been plagued by conceptual puzzles. What made matters worse for a long time was the fact that many past experiments seemed to contradict the most firmly-established theoretical results. However, considerable progress has been made over the past decade, and one main aim of this chapter is to demonstrate that the theory and experiments fit within a consistent picture. This has become possible by the realization that dynamic instabilities of the tip of a fracture play a critical role in determining the fracture behavior of amorphous materials. To accomplish this goal, we follow our by-now-familiar path, namely, we first describe and summarize the central results of continuum theories of linear elastic dynamic fracture mechanics which provides an elegant and powerful description of fracture propagation. However, the continuum theory is unable to make quantitative predictions without additional information that must be provided by experiments, or be supplied by other types of theories. We already discussed in the last chapter some of the most important experimental observations and data, and the techniques that were used for obtaining them. These experiments teach us that when the flux of energy to a fracture tip exceeds a critical value, the fracture becomes unstable and hence propagates in an increasingly complex manner. As a result, the moving crack cannot travel as quickly as the linear continuum theory predicts or assumes, the fracture surface becomes rough and begins to branch out and radiate sound, and the energy cost for the motion of the crack increases significantly. These observations are completely consistent with the continuum theory, but *cannot* be described by it. Therefore, to complete the emerging theoretical picture and the fundamental understanding of this phenomenon, we continue this chapter with an account of theoretical and numerical work of the past decade or so that attempts to explain the dynamic instabilities in fracture propagation. As discussed in the last chapter, our current experimental understanding of instabilities in fracture tip in brittle amorphous materials is fairly detailed. We also have a rather detailed theoretical understanding of these instabilities in crystals which reproduces many qualitative features of the experiments. Recent numerical work

is attempting to establish the missing connections between the experiments and the theory.

Up until a decade or so ago, most engineers and materials scientists believed that the development of continuum fracture mechanics is largely complete. Why? Because this field is in fact one of the most heavily developed branches of engineering science. We only need to consider how many books and review articles have been written on this subject to appreciate this fact. The development of continuum fracture mechanics actually emerged from mathematical exercises in the early part of the 20th century into a coherent collection of theoretical concepts and experimental techniques that are now widely used to ensure the safety of critical structures, ranging from aircraft to microelectronic devices. Despite considerable progress, two important and puzzling features of the problem kept researchers attracted to fracture of brittle materials. The first feature is that it is often stated that propagating fractures do not reach the limiting velocity predicted by linear continuum mechanics of fracture propagation, and that they have a seemingly unexplained instability at a critical velocity of propagation which is between the prediction of the linear theory and the experimental data. In fact, only about a decade ago, Freund (1990) specifically mentioned in his book (pp. 37–38) in a short list of phenomena (associated with dynamic fracture) entitled "not yet completely understood" the *apparent terminal fracture speed well below the Rayleigh wave speed in glass and some other very brittle materials.* The Rayleigh wave speed $c_R$ is the speed at which sound travels over a free surface. The root cause of this apparent inconsistency is in the energy dissipation at the fracture tip and, as we discuss in this chapter, recent work indicates that when energy flux into a crack tip exceeds a certain critical value, efficient and steady motion of the tip becomes unstable to the formation of microfractures that propagate away from the main fracture. In fact, the tip undergoes a hierarchy of instabilities which increases enormously its ability to absorb energy. The second feature is the need for understanding how materials break at the atomic length scale. To understand this aspect of the problem one must resort to molecular dynamics (MD) simulations which enable one to model generation of fracture and their motion one atomic bond at a time. However, MD simulations require extensive and very time consuming computations. To make the simulations efficient and cost effective, a sound strategy is perhaps to study the existing analytical results so as to understand the qualitative effect of atomic discreteness on crack motion. Once this understanding is acquired, many experimental results become understandable, the relation between simulations and experiments becomes clearer, and therefore MD simulations will be much more efficient. We will describe MD simulations of fracture propagation in Chapters 9 and 10.

These puzzling, and theoretically challenging, features of dynamics of brittle fracture of materials have motivated a considerable amount of work in this research area, especially by physicists and their allied scientists. Their work has helped the emergence of a much clearer picture of fracture dynamics which indicates that the two puzzling features of fracture dynamics, at both atomic and macroscopic length scales, are in fact manifestations of the same underlying phenomenon. One

goal of this chapter is to explain how these puzzles have arisen, and how to recast them in new terms and explain them. We do not intend to provide a complete review of fracture mechanics as it will require a book by itself. Instead, we focus on brittle materials. Ductility and dynamic elasto-plastic fracture, which is a well-developed field, have been described well by others (see, for example, Freund, 1990; Chan, 1997). Therefore, we will discuss only the brittle-to-ductile transition. The emphasize in this chapter is first to describe and summarize the most important predictions of the conventional continuum fracture mechanics, and then answer some fundamental and interesting questions that this type of models do not ask or, if they do, cannot answer. To write a significant portion of this chapter, we relied heavily on the excellent review of this subject by Fineberg and Marder (1999). Some of the developments that we discuss had been described in an earlier article by the author (Sahimi, 1998), and thus have also been utilized in this chapter.

## 7.1   Scaling Analysis

Before embarking on a detailed analysis of fracture of materials, we carry out some preliminary scaling analysis of this problem. Although our analysis is too simple-minded, it does point to some fundamental properties of materials, and does exhibit some basic problems that a detailed analysis of fracture propagation must address. We consider both the static and dynamical cases.

### 7.1.1   Scaling Analysis of Materials Strength

Despite what most of us believe (and apparently feel), the world is farther from equilibrium than we realize. To see this, consider a piece of rock of area $S$ and height $h$. Equilibrium principles teach us that the rock should not be able to sustain its own weight under the force of gravity, if it becomes too tall. To estimate the critical height, recall that the gravitational potential energy of the rock is $\frac{1}{2}\rho S h^2 g$ where $\rho$ is the rock's density. If we cut the rock into two equal blocks of height $\frac{1}{2}h$ and set them side by side, this energy is reduced to $\frac{1}{4}\rho S h^2 g$, resulting in an energy gain of $\frac{1}{4}\rho S h^2 g$. The cost in energy of the cut is the same as the cost of creating new rock surface, the characteristic value of which per unit area is, $\mathcal{H} = 1 \text{ J/m}^2$. If we assume a typical value, $\rho = 2000 \text{ kg/m}^3$, the critical height $h_c$ at which it pays to divide the rock in two is, $h = \sqrt{4\mathcal{H}/\rho g} \sim 1.4$ cm, so that every block of rock more than about 2 cm tall is unstable under its own weight. Similar scaling analyses are applicable to steel or concrete. Thus, although things fall apart when they reach equilibrium, the time to reach this state is fortunately long.

Since the fact that most objects do not fall apart easily is an indication that they are out of mechanical equilibrium, one must estimate the size of the energy barriers that hold them in place. A rough estimate is obtained by imagining what happens to the atoms of a solid material as one pulls it uniformly at two ends. Initially, the forces between the atoms increase, but they eventually reach a maximum value, at which the material breaks into pieces.

As is well known, interatomic forces vary greatly between different elements and molecules (see Chapter 9), but they typically attain their maximum value when the distance between atoms increases by about 20% of its original value. The force needed to stretch a solid material slightly is, $F = YS\delta/L$, where $Y$ is the Young's modulus, $L$ is its initial length, and $\delta$ is the amount (in length) that the material has been stretched. Therefore, the force per unit area needed to reach the breaking point is about, $\sigma_c = F/S = Y\delta/L \simeq Y/5$, where we have used $\delta/L = 0.2$. We list in Table 7.1 values of $Y$ for several materials, the theoretical strength $\sigma_c$, and its comparison with the experimental data. As this table indicates, the theoretical estimate of $\sigma_c$ is in error by *orders of magnitude*. The scaling estimate of $h_c$ greatly underestimates the practical resistance of solid materials to fracturing, while the estimate of $\sigma_c$ too large. What is the problem? The only way to discuss the correct orders of magnitude is to account for the actual dynamical mode by which brittle materials fail mechanically, which is by propagation of a fracture.

As described in Chapter 6, and shown later in this chapter, the presence of a fracture in an otherwise perfect material results in a stress singularity at the fracture tip. If the fracture tip is atomically sharp, a single fracture which is a few microns long suffices for explaining the large discrepancies between the theoretical and experimental material strengths that are shown in Table 7.1. The stress singularity that develops at the tip of a fracture focuses the energy that is stored in the surrounding material and uses it efficiently for breaking one atomic bond after another. Thus, continuous fracture propagation provides an efficient way of overcoming the energy barrier between two equilibrium states of the system that have different amounts of mechanical energy. We now turn to a scaling analysis of dynamic fracture (Fineberg and Marder, 1999).

### 7.1.2  Scaling Analysis of Dynamic Fracture

An analysis of rapid fracture was first carried out by Mott (1948) whose analysis was slightly improved by Dulaney and Brace (1960). Mott's work is a dimensional

TABLE 7.1. The experimental strength $\sigma_e$ of a number polycrystalline or amorphous materials, and their comparison with the corresponding theoretical strength $\sigma_c$. $Y$ is the materials' Young's modulus (adopted from Fineberg and Marder, 1999).

| Material | $Y$ (GPa) | $\sigma_c$ (GPa) | $\sigma_e$ (GPa) | $\sigma_e/\sigma_c$ |
|---|---|---|---|---|
| Iron | 195-205 | 43-56 | 0.3 | 0.006 |
| Copper | 110-130 | 24-55 | 0.2 | 0.005 |
| Titanium | 110 | 31 | 0.3 | 0.009 |
| Silicon | 110-160 | 45 | 0.7 | 0.01 |
| Glass | 70 | 37 | 0.4 | 0.01 |
| Plexiglas | 3.6 | 3 | 0.05 | 0.01 |

FIGURE 7.1. Propagation of a fracture of length $l$ at velocity $v$ in an infinite plate disturbs the material up to a distance $l$ (after Fineberg and Marder, 1999).

analysis which, despite being wrong in many of its details, clarifies the basic physical processes. It consists of writing down an energy balance for the motion of a fracture. Consider a fracture of length $l(t)$ growing at time $t$ at rate $v(t)$ in a very large plate to which a stress $\sigma_\infty$ is applied at its far boundaries; see Figure 7.1. As the fracture extends, its faces separate, causing the plate to relax within a circular region centered at the middle of the crack with a diameter which is of the order of $l$. The kinetic energy $\mathcal{H}_k$ involved in moving a piece of material of this size is $\frac{1}{2}mv^2$, where $m$ is the total mass, and $v$ is a characteristic velocity. Since the mass of the moving material is proportional to $l^2$, the kinetic energy should be given by

$$\mathcal{H}_k(l, v) = c_k l^2 v^2, \tag{1}$$

where $c_k$ is a constant. The moving portion of the material is also where elastic potential energy is being released as the crack propagates. This stress release results in a gain in the potential energy which is given by

$$\mathcal{H}_p(l) = -c_p l^2, \tag{2}$$

where $c_p$ is another constant. Equations (1) and (2) are correct if the crack moves slowly, but they fail even qualitatively if the fracture velocity approaches the speed of sound, in which case $\mathcal{H}_k$ and $\mathcal{H}_p$ both diverge. Their divergence will be demonstrated below, but let us assume for now that it is true. The final piece of the energy balance is the contribution of creation of new fracture surfaces. This contribution is $\Gamma l$, where $\Gamma$ is the *fracture energy* that, as described in Chapter 6 (see Section 6.7), accounts for the minimum energy needed to break the atomic bonds and any other dissipative processes that the material may need in order for the fracture to propagate, and is often orders of magnitude greater than the thermodynamic surface energy. Therefore, the total energy $\mathcal{H}$ of the system containing a fracture

is given by

$$\mathcal{H}(l, v) = c_k l^2 v^2 + \mathcal{H}_{qs}(l), \tag{3}$$

where $\mathcal{H}_{qs}$ is the quasi-static part of the total energy given by

$$\mathcal{H}_{qs}(l) = -c_p l^2 + \Gamma l. \tag{4}$$

If a crack moves forward slowly, its kinetic energy will be negligible, and therefore only $\mathcal{H}_{qs}$ will be important. For small fractures, $\Gamma l$, the linear cost of fracture energy, is always greater than the quadratic gain of the potential energy, $\mathcal{H}_p = c_p l^2$. In fact, such fractures would heal (move backward) if such irreversible processes as oxidation of the crack surface did not prevent them from healing. The fact that the fracture grows is due to additional irreversible processes, such as chemical attack on the crack tip (see Chapter 6), or vibration and other irregular mechanical stresses. Eventually, at a critical length $l_c$, the energy gained by relieving elastic stresses in the material exceeds the cost of creating new fracture surfaces, in which case the crack is able to extend spontaneously. Clearly, at $l_c$, the energy functional $\mathcal{H}_{qs}(l)$ has a quadratic maximum. The Griffith criterion (Griffith, 1920; see Chapter 6 and also below) for the onset of fracture is that fracture occurs when the potential energy released per unit crack extension equals the fracture energy $\Gamma$. Thus, fracture in this system occurs at a critical crack length $l_c$ such that, $d\mathcal{H}_{qs}/dl = 0$ at $l = l_c$. Using Eq. (4) we find that,

$$l_c = \frac{\Gamma}{2c_p}, \tag{5}$$

so that

$$\mathcal{H}_{qs}(l) = \mathcal{H}_{qs}(l_c) - c_p(l - l_c)^2. \tag{6}$$

The most important issue in engineering fracture mechanics is calculating $l_c$, given such information as the external stresses which, in the present case, is represented by the constant $c_p$. Dynamic fracture begins in the next instant, and since it is very rapid, the energy $\mathcal{H}$ of the system is conserved, remaining at $\mathcal{H}_{qs}(l_c)$. Thus, from Eqs. (3) and (6) we obtain

$$v(t) = \frac{c_p}{c_k} \left(1 - \frac{l_c}{l}\right) = v_m \left(1 - \frac{l_c}{l}\right), \tag{7}$$

which predicts that fracture propagation will accelerate until it approaches the maximum speed $v_m$. Equation (7), and more generally the above scaling analysis, cannot by themselves predict $v_m$, but Stroh (1957) argued correctly that $v_m$ should be the Rayleigh wave speed $c_R$, although his suggestion was implicitly contained in the earlier calculations of Yoffe (1951) (see below).

In this system, one needs only to know the length $l_c$ at which a fracture begins to propagate in order to predict all the ensuing dynamics. As we discuss later in this chapter, Eq. (7) is actually very close to anticipating the results of a far more sophisticated analysis, which is surprising since the Eqs. (1), (3) and (4) for the kinetic and potential energy are in fact incorrect because they actually

diverge as the speed of fracture propagation approaches the Rayleigh wave speed $c_R$. However, the success of Eq. (7) is due to the fact that it involves the ratio $\mathcal{H}_p/\mathcal{H}_k$. Since the divergence of the kinetic and potential energy are according to exactly the same forms, the errors involved in their estimation cancel each other out. We now attempt to review and discuss the background, basic formalism and underlying assumptions that form the basis for continuum fracture mechanics.

## 7.2    Continuum Formulation of Fracture Mechanics

The general strategy in continuum fracture mechanics is to solve for the displacement fields in the material subject to both the boundary conditions and the externally applied stresses. The elastic energy transmitted by the displacement fields is then matched to the amount of energy dissipated throughout the material, which results in an equation of motion. The only energy sink in a single moving fracture is at the tip of the fracture itself. Thus, an equation of motion for a moving fracture is obtained if detailed knowledge of the dissiption mechanisms in the vicinity of the fracture tip is available.

### 7.2.1    Dissipation and the Cohesive Zone

As discussed in Chapter 6, the processes that give rise to energy dissipation in the vicinity of the crack tip are complex and, depending on the material, vary from dislocation formation and emission in crystalline materials to the complex unraveling and fracture of intertangled polymer strands in amorphous polymers. Fracturing and the complex dissipative processes occurring in the vicinity of the crack tip occur due to very large values of the stress field as one approaches the tip. As discussed below, if the material around the crack tip were to remain linearly elastic until fracture, the stress field at the crack tip would actually diverge. Since a real material cannot support such singular stresses, the assumption of linearly elastic behavior in the vicinity of the tip must break down and *material-dependent* dissipative processes must begin playing an important role. Given the enormous variety of materials, the emergence of material-dependent dissipative processes might indicate that a *universal* description of fracture is impossible. However, as described and discussed in Chapter 6, Orowan (1955) and Irwin (1956) developed a way around this difficulty by suggesting independently that the region around the fracture tip should be divided into three separate regions which, as described in Section 6.9, are as follows.

(1) The cohesive zone (also called the *process zone*), which is the region immediately surrounding the fracture tip in which all the nonlinear dissipative processes that allow a crack to move (forward) are assumed to occur. In continuum fracture mechanics detailed description of this zone is avoided. Instead, this zone is simply characterized by the energy $\Gamma$, per unit area of crack extension, that it consumes during fracture propagation. As discussed in Chapter 6,

the size of the cohesive zone depends on the material, ranging from nanometers in glass to microns in brittle polymers.

(2) The universal elastic region, which is the region outside the cohesive zone for which the response of the material can be described by linear continuum mechanics. Outside the cohesive zone, but in the vicinity of the fracture tip, the stress and strain fields take on *universal* singular forms which depend only on the symmetry of the externally applied loads. In two dimensions (2D) the singular fields surrounding the cohesive zone are completely described by three constants which are the stress intensity factors introduced and discussed in Section 6.8 (see also below). They incorporate all the information regarding the loading of the material.

(3) The outer elastic region far from the crack tip in which stresses and strains are described by linear elasticity. Details of the solution to the stress field in this region of materials depend only on the locations and strengths of the loads, and the shape of the material. For some special cases analytical solutions have been derived. Deriving such solutions is made possible by the fact that, so far as linear elasticity is concerned, viewed on macroscopic scales, the cohesive zone can be represented by just a point at the fracture tip, while the fracture itself is equivalent to a branch cut. In general, however, one must resort to numerical simulations and solutions.

The dissipative processes within the cohesive zone determine the fracture energy $\Gamma$. If no dissipative processes other than the direct breaking of the atomic bonds take place, then $\Gamma$ is a constant which depends on the bond energy. In general though, $\Gamma$ is a complex function of both the fracture velocity and history, and differs by orders of magnitude from the *surface energy*—the amount of energy required to sever a unit area of atomic bonds. No general first principle description of the cohesive zone exists, although numerous models have been proposed (see, for example, Lawn, 1993).

## 7.2.2   *Universal Singularities near the Fracture Tip*

As one approaches the tip of a fracture in a linearly elastic material, the stress field surrounding the tip develops a square root singularity (in the distance $r$). As mentioned in Section 6.8, Irwin (1958) noted that the stress field at a point $(r, \theta)$ near the fracture tip, measured in polar coordinates with the fracture line corresponding to $\theta = 0$, can be represented by

$$\sigma_{ij} = \frac{K_\beta}{\sqrt{2\pi r}} f_{ij}^\beta(v, \theta), \tag{8}$$

where $v$ is the instantaneous crack velocity, and $\beta$ is an index that represents Modes I, II and III of fracture described in Section 6.6. For each of these three symmetrical loading configurations, $f_{ij}^\beta(v, \theta)$ is a known *universal* function. The coefficients $K_\beta$ is the stress intensity factor, introduced in Chapter 6, that contains all the detailed information about sample loading and history, and is determined by

the elastic fields throughout the material. However, the stress that locally drives the fracture is one which is present at its tip. Thus, $K_\beta$ determines entirely the behavior of a fracture, and much of the study of fracture processes is aimed toward either calculating or measuring this quantity. The universal form of the stress intensity factor allows a complete description of the behavior of the tip of a fracture where one needs only carry out the analysis of a given problem within the universal elastic region (see below). For arbitrary loading configurations, the stress field around the fracture tip is given by three stress intensity factors $K_\beta$ which lead to a stress field that is a linear combination of the pure Modes:

$$\sigma_{ij} = \sum_{\beta=1}^{3} \frac{K_\beta}{\sqrt{2\pi r}} f_{ij}^\beta (v, \theta). \tag{9}$$

As mentioned above, the stress intensity factors are related to the flow of energy into the fracture tip. Since a fracture may be viewed as a means of dissipating built-up energy in a material, the amount of energy flowing into its tip must influence its behavior. Irwin (1956) showed that the stress intensity factor is related to the energy release rate $\mathcal{H}$, defined as the amount of energy flowing into the crack tip per unit fracture surface formed. The relation between the two quantities is given by

$$\mathcal{H} = \sum_{\beta=1}^{3} \frac{1 - v_p^2}{Y} A_\beta(v) K_\beta^2, \tag{10}$$

where $v_p$ is the Poisson's ratio of the material, and the three functions $A_\beta(v)$ depend only on the fracture velocity $v$. Equation (10) is accurate when the stress field near the tip of a fracture can be accurately described by Eq. (8), which is the case as the dimensions of the sample increase.

## 7.3   Linear Continuum Theory of Elasticity

Since most of the theoretical work that we describe in this chapter is carried out in 2D (or quasi-2D) systems, we follow the analysis presented by Fineberg and Marder (1999) who performed a reduction of the full 3D elastic description of a fracture to 2D in three important cases: For Mode III fracture, and Mode I fracture in very thin and very thick plates. As noted in Section 6.6, Mode III fracture is an important model system for which much analytical work has been carried out, resulting in deeper gains in understanding qualitative features of fracture. The second case, Mode I fracture of a thick plate, describes stress and strain conditions of importance in describing the phenomenon in the immediate vicinity of the fracture tip. The third case, Mode I fracture in thin plates, corresponds to much of the experimental work that was described in Chapter 6, some of which will also be considered in the present chapter.

As already described and discussed in detail in Chapter 7 of Volume I, the starting point is the Navier equation of motion for an isotropic elastic material:

$$\rho \frac{\partial^2 \mathbf{u}}{\partial t^2} = (\lambda + \mu)\nabla(\nabla \cdot \mathbf{u}) + \mu\nabla^2\mathbf{u}, \tag{11}$$

where $\mathbf{u}$ is the displacement field for each mass point relative to its original location in an unstrained material, and $\rho$ is the density. The constants $\mu$ and $\lambda$ are the usual Lamé constants (with dimensions of energy per volume and typical values of order of $10^{10}$ erg/cm$^3$). We also define the linear elastic strain tensor with components

$$\epsilon_{ij} = \frac{1}{2}\left(\frac{\partial u_i}{\partial x_j} + \frac{\partial u_j}{\partial x_i}\right). \tag{12}$$

If a linear stress-strain relation exists in a homogeneous and isotropic material, the components $\sigma_{ij}$ of the stress tensor are defined by

$$\sigma_{ij} = \lambda\delta_{ij}\sum_k \varepsilon_{kk} + 2\mu\epsilon_{ij}. \tag{13}$$

The simplest analytical results are obtained for pure Mode III. The only non-zero displacement is $u_z = u_z(x, y)$ alone. Thus, the only non-vanishing stresses are, $\sigma_{xz} = \mu\partial u_z/\partial x$, and, $\sigma_{yz} = \mu\partial u_z/\partial y$. The governing equation for $u_z$ is the ordinary wave equation,

$$\frac{1}{c^2}\frac{\partial^2 u_z}{\partial t^2} = \nabla^2 u_z, \tag{14}$$

where $c = \sqrt{\mu/\rho}$.

Consider now Mode I fracture in a sample material that is extremely thick along the $z$-direction. All the applied forces are uniform in this direction. Because all the derivatives with respect to $z$ vanish, all the fields are functions of $x$ and $y$ alone, so that one deals with a plane strain problem. The reduction of the problem to 2D is simple, but this geometry is not convenient for experiments. A third case in which the equations of elasticity reduce to 2D is the plane stress problem in which one pulls on a thin plate in Mode I. If the length scale over which the stresses vary in $x$ and $y$ is large compared with the thickness of the plate along the $z$-direction, then we might expect the displacements in that direction to quickly reach equilibrium with the local stresses. If the Poisson's ratio is positive, then when the material is stretched, the plate will contract in the $z$-direction, and if it is compressed, the plate will thicken. (Counter-examples, when the material expands under stretching, were described in Section 9.8 of Volume I.) Under this condition, $u_x$ and $u_y$ are independent of $z$, and therefore it is reasonable to assume that,

$$u_z = zf(u_x, u_y). \tag{15}$$

The function $f$ can be found by realizing that the stress $\sigma_{zz}$ must vanish on the face of the plate, implying that at the surface of the plate we must have

$$\lambda\left(\frac{\partial u_x}{\partial x} + \frac{\partial u_y}{\partial y}\right) + (\lambda + 2\mu)\frac{\partial u_z}{\partial z} = 0, \tag{16}$$

which means that

$$f(u_x, u_y) = \frac{\partial u_z}{\partial z} = -\frac{\lambda}{\lambda + 2\mu} \left( \frac{\partial u_x}{\partial x} + \frac{\partial u_y}{\partial y} \right), \tag{17}$$

so that

$$\frac{\partial u_x}{\partial x} + \frac{\partial u_y}{\partial y} + \frac{\partial u_z}{\partial z} = \frac{2\mu}{\lambda + 2\mu} \left( \frac{\partial u_x}{\partial x} + \frac{\partial u_y}{\partial y} \right). \tag{18}$$

Therefore,

$$\sigma_{\alpha\beta} = \tilde{\lambda} \delta_{\alpha\beta} \frac{\partial u_\gamma}{\partial x_\gamma} + \mu \left( \frac{\partial u_\alpha}{\partial x_\beta} + \frac{\partial u_\beta}{\partial x_\alpha} \right), \tag{19}$$

with

$$\tilde{\lambda} = \frac{2\mu\lambda}{\lambda + 2\mu}, \tag{20}$$

and $\alpha$ and $\beta$ now run only over $x$ and $y$. Therefore, a thin plate satisfies the equations of 2D elasticity, with an *effective constant* $\tilde{\lambda}$, so long as $u_z$ is dependent upon $u_x$ and $u_y$ according to Eqs. (15)–(17). In the following discussion, the tilde over $\lambda$ is dropped with the understanding that the relation to 3D materials properties is given by Eq. (20). The equation of motion is still the Navier equation, but is restricted to 2D.

Note that, as described in Chapter 7 of Volume I, materials are frequently described by the Young's modulus $Y$ and the Poisson's ratio $\nu_p$, in terms of which we have

$$\lambda = \frac{Y\nu_p}{(1 + \nu_p)(1 - 2\nu_p)}, \quad \tilde{\lambda} = \frac{Y}{2(1 - \nu_p^2)}, \quad \mu = \frac{Y}{2(1 + \nu_p)}. \tag{21}$$

Moreover, note that

$$\nabla \cdot \mathbf{u} = (\lambda + 2\mu) \sum_\alpha \sigma_{\alpha\alpha}, \tag{22}$$

and that from Eq. (11) one finds that

$$\frac{\rho}{\lambda + 2\mu} \frac{\partial^2 \sigma_{\alpha\alpha}}{\partial t^2} = \nabla^2 \sigma_{\alpha\alpha}. \tag{23}$$

Therefore, $\nabla \cdot \mathbf{u}$ satisfies the wave equation with the longitudinal wave speed

$$c_l = \frac{1}{\rho}(\lambda + 2\mu), \tag{24}$$

whereas, while $\nabla \times \mathbf{u}$ also satisfies the wave equation, it does so with the shear (compressional) wave speed

$$c_t = \sqrt{\mu/\rho}. \tag{25}$$

One must also consider the transition from 2D to 3D. Near the tip of a fracture in a plate stresses become large enough that the approximations leading to 2D plane

stress elasticity fail (Nakamura and Parks, 1988). If the thickness of the plate along the $z$-direction is denoted by $d$, then at distances from the fracture tip that are much larger than $d$ all fields are described by equations of plane stress. At distances from the fracture tip that are much less than $d$, and away from the $x - y$ surfaces of the plate, the fields solve the equations of plane strain.

### 7.3.1   Static Fractures in Mode III

If one inserts an elliptical crack in a plate and pulls it, then, as discussed in Section 6.7, Inglis (1913) was the first to derive the expression for the stresses at the crack's narrow ends, and found that they are much larger than those exerted off at infinity. Therefore, a crack acts as an amplifier of the stresses and causes the elastic energy to be preferentially focused into its tip, implying that the existence of a crack leads to a large decrease in the effective strength of a material. The ratio of the maximum to the applied stress is

$$\frac{\text{Maximum stress}}{\text{Applied stress}} = 2\frac{l}{\mathcal{Y}}, \tag{26}$$

where $l$ is the crack's length and $\mathcal{Y}$ the radius of curvature at its tip. Thus, if one assumes that typical solids have fracture tips of size 1 Å and length of $10^4$ Å, then one can account for the discrepancies shown in Table 7.1. To derive Eq. (26) we assume that the stresses applied to the plate coincides with the conditions of anti-plane shear stress, so that the only non-zero displacement is $u_z$. From Eq. (14) one sees that the static equation of linear elasticity is now simply the Laplace's equation, $\nabla^2 u_z = 0$. For our boundary value problem conformal mapping is the appropriate technique. Since $u_z$ is a solution of the Laplace's equation, it can be represented by

$$u_z = \frac{1}{2}[\phi(\zeta) + \overline{\phi(\zeta)}], \tag{27}$$

where $\phi$ is analytic, $\zeta = x + iy$, and $\bar{\phi}$ is the complex conjugate of $\phi$.

Far from the crack, the displacement $u_z(x, y)$ increases linearly with $y$, and therefore we must have the asymptotic property that

$$\phi = -ic\zeta. \tag{28}$$

Although the constant $c$ of Eq. (28) is dimensionless, in essence it measures the stress in units of the Lamé constant $\mu$. Because the crack's edges are free, the stress normal to the edge must vanish. It can then be shown that

$$\phi(\zeta) = \overline{\phi(\zeta)}, \tag{29}$$

when $\zeta$ lies on the boundary. To illustrate the use of Eq. (29), let us define $\omega$ such that

$$\zeta = \frac{l}{2(\omega + m/\omega)}. \tag{30}$$

When $\omega$ lies on the unit circle (i.e., $\omega = e^{i\theta}$, with $\theta$ real), $\zeta$ traces out an elliptical

boundary. When $m = 0$, the boundary is a circle of radius $l/2$, whereas when $m = 1$, the boundary is a cut, i.e., a straight fracture along the real axis extending from $-l$ to $+l$. The function $\phi(\omega)$ has the properties that, $\phi(\omega) = \overline{\phi(\omega)} = \bar{\phi}(1/\omega)$. The last property is due to the fact that, on the unit circles, $\bar{\omega} = 1/\omega$. These properties can be analytically continued outside the unit circle, where $\phi$ must be completely regular except that, for large $\zeta$, it should diverge as $-ic\zeta$. From Eq. (30) we see that for large $\zeta$ we must have $\omega \simeq \zeta$ and that $\phi \sim -ic\omega$ as $\omega \to -\infty$, implying (using the above properties of $\phi$) that, as $\omega \to 0$, we must have $\phi(\omega) \sim -ic/\omega$, and therefore $\phi(\omega) = -ic\omega + ic/\omega$. It is then straightforward to show (Fineberg and Marder, 1999) that the displacement $u_z$ is finite as one approaches the fracture tip, but the stress

$$\sigma_{yz} = \mu \frac{\partial u_z}{\partial y} \sim (z-1)^{-1/2}, \quad z \to 1, \tag{31}$$

diverges as one approaches the crack tip.

Although Eq. (31) was derived for a particular case, its main feature, namely, the existence of a square root stress singularity at the fracture tip, is of general validity and confirms Eq. (8), a feature that was already mentioned in Chapter 6. Thus, if a fracture is given a finite radius of curvature, the singularity is effectively removed. An amazing, and counterintuitive, application of this idea, that was pointed out by Fineberg and Marder (1999), is to arresting the advance of a fracture in a damaged material by drilling a small hole at the fracture tip, since the hole increases the tip's radius of curvature and hence blunts the singularity in the stress field. As a result, the strength of the material increases sharply!

The conformal mapping method outlined above for Mode III cracks was extended to Mode I by Muskhelishvili (1953). The problem in this case is more complex as one must solve the biharmonic equation rather than the Laplace's equation, and solve for two complex functions not one. Since Muskhelishvili's work hundreds of papers have been devoted to solutions of fracture problems using these methods, a review of which will occupy a book by itself.

## 7.3.2   *Dynamic Fractures in Mode I*

According to Eq. (31), in an elastic material to which a uniform stress is applied at its boundaries, the stress field at the tip of a static fracture is singular. Let us now consider the case of a propagating fracture and examine the structure of the stress field at its tip in Mode I. The dynamical equation for the displacement field **u** of a steady state in a frame moving with a constant velocity $v$ in the $x$-direction is given by

$$(\lambda + \mu)\nabla(\nabla \cdot \mathbf{u}) + \mu\nabla^2\mathbf{u} = \rho v^2 \frac{\partial^2 \mathbf{u}}{\partial x^2}. \tag{32}$$

If we decompose **u** into longitudinal and transverse parts, $\mathbf{u} = \mathbf{u}_l + \mathbf{u}_t$, with

$$\mathbf{u}_l = \nabla v_l, \quad \mathbf{u}_t = \left(\frac{\partial v_t}{\partial y} - \frac{\partial v_t}{\partial x}\right), \tag{33}$$

it follows immediately that $\mathbf{u}_l$ satisfies the following equation

$$\left[(\lambda + 2\mu)\nabla^2 - \rho v^2 \frac{\partial^2}{\partial x^2}\right]\mathbf{u}_l = -\left(\mu\nabla^2 - \rho v^2 \frac{\partial^2}{\partial x^2}\right)\mathbf{u}_t = \mathbf{f}(x, y). \tag{34}$$

It can be shown that, $\mathbf{f} = \mathbf{0}$. If

$$\alpha^2 = 1 - \frac{\rho v^2}{\lambda + 2\mu} = 1 - \frac{v^2}{c_l^2}, \tag{35}$$

$$\beta^2 = 1 - \frac{\rho v^2}{\mu} = 1 - \frac{v^2}{c_t^2}, \tag{36}$$

then, the general forms of $v_l$ and $v_t$ are (Fineberg and Marder, 1999)

$$v_l = v_l^0(z) + \overline{v_l^0(z)} + v_l^1(x + i\alpha y) + \overline{v_l^1(x + i\alpha y)}, \tag{37}$$

$$v_t = v_t^0(z) + \overline{v_t^0(z)} + v_t^1(x + i\beta y) + \overline{v_t^1(x + i\beta y)}. \tag{38}$$

However, it can be shown that, $v_l^0 = v_t^0 = 0$. Therefore, if we define $\phi(z) = \partial v_l^1/\partial z$ and $\psi(z) = \partial v_t^1/\partial z$, the components of $\mathbf{u} = (u_x, u_y)$ are given by,

$$u_x = \phi(z_\alpha) + \overline{\phi(z_\alpha)} + i\beta[\psi(z_\beta) - \overline{\psi(z_\beta)}], \tag{39}$$

$$u_y = i\alpha[\phi(z_\alpha) - \overline{\phi(z_\alpha)}] - [\psi(z_\beta) + \overline{\psi(z_\beta)}], \tag{40}$$

where, $z_\alpha = x + i\alpha y$, and $z_\beta = x + i\beta y$.

Equations (37) and (38) provide general expressions for steady-state elastic problems in which a fracture propagates with a velocity $v$. If we define $\Phi = \partial\phi(z)/\partial z$ and $\Psi = \partial\psi(z)/\partial z$, then the stresses are given by

$$\sigma_{xx} = \mu(1 + 2\alpha^2 - \beta^2)[\Phi(z_\alpha) + \overline{\Phi(z_\alpha)}] + 2i\beta\mu[\Psi(z_\beta) - \overline{\Psi(z_\beta)}], \tag{41}$$

$$\sigma_{yy} = -\mu(1 + \beta^2)[\Phi(z_\alpha) + \overline{\Phi(z_\alpha)}] - 2i\beta\mu[\Psi(z_\beta) - \overline{\Psi(z_\beta)}], \tag{42}$$

$$2\sigma_{xy} = 2\mu\left\{2i\alpha[\Phi(z_\alpha) - \overline{\Phi(z_\alpha)}] - (\beta^2 + 1)[\Psi(z_\beta) + \overline{\Psi(z_\beta)}]\right\}. \tag{43}$$

Equations (41)–(43) represent the general solutions in which the functions $\phi$ and $\psi$ must match the boundary conditions that are specified. Since one wishes to find the potentials from given stresses at the boundaries, $\Phi$ must diverge as $1/v$, and the right-hand sides of Eqs. (41)–(43) turn into the derivative of $\Phi$ with respect to $\alpha$, implying that the static theory has a different structure than the dynamical theory which is in fact more straightforward.

Let us now derive, as an application of Eqs. (37)–(43), the expressions for the stresses around the tip of a fracture moving under Mode I loading. We assume that the fracture lies along the negative $x$-axis (terminating at $x = 0$) and propagates forward. The only assumption is that the problem is symmetric under reflection about the $x$-axis. As discussed above (and also in Chapter 6), in the static case, the stress fields have a square root singularity at the crack tip. We assume the same to be true in the dynamic case (which can be verified in all cases for which the

expressions have been derived). Therefore, we assume that near the fracture tip (Fineberg and Marder, 1999)

$$\phi(z) \sim (b_r + ib_i)z^{-1/2}, \tag{44}$$

$$\psi(z) \sim (d_r + id_i)z^{-1/2}. \tag{45}$$

Since we are considering Mode I fracture, then by symmetry the displacements satisfy

$$u_x(-y) = u_x(y), \quad u_y(-y) = -u(y). \tag{46}$$

If we substitute Eqs. (44) and (45) into (39) and (40) and use Eq. (46), we find that $b_i = d_r = 0$, and therefore

$$\Phi(z) \sim b_r z^{-1/2}, \tag{47}$$

$$\Psi(z) \sim id_i z^{-1/2}. \tag{48}$$

Observe that the square roots in Eqs. (44) and (45) must be interpreted as having their cuts along the negative $x$-axis, where the fracture is located. Since on the crack surface the stresses are relaxed, $\sigma_{xy}$ and $\sigma_{yy}$ vanish there. If we substitute Eqs. (47) and (48) into Eqs. (41)–(43), we find that the condition for $\sigma_{yy}$ is satisfied identically for $x < 0$, $y = 0$, and that at $y = 0$

$$\sigma_{xy} = i\mu \left[ 2\alpha b_r - (\beta^2 + 1)d_i \right] (1/\sqrt{x} - \overline{1/\sqrt{x}}), \tag{49}$$

and therefore

$$\frac{d_i}{b_r} = \frac{2\alpha}{\beta^2 + 1}, \tag{50}$$

which, when used in Eqs. (41)–(43), (47) and (48), yields

$$\sigma_{xx} = \frac{K_I}{\sqrt{2\pi D}} \left[ (\beta^2 + 1)(1 + 2\alpha^2 - \beta^2) \left( \frac{1}{\sqrt{z_\alpha}} + \frac{1}{\sqrt{\bar{z}_\alpha}} \right) - 4\alpha\beta \left( \frac{1}{\sqrt{z_\beta}} + \frac{1}{\sqrt{\bar{z}_\beta}} \right) \right], \tag{51}$$

$$\sigma_{yy} = \frac{K_I}{2\sqrt{2\pi D}} \left[ 4\alpha\beta \left( \frac{1}{\sqrt{z_\beta}} + \frac{1}{\sqrt{\bar{z}_\beta}} \right) - (1 + \beta^2)^2 \left( \frac{1}{\sqrt{z_\alpha}} + \frac{1}{\sqrt{\bar{z}_\alpha}} \right) \right], \tag{52}$$

$$\sigma_{xy} = \frac{2i\alpha K_I}{2\sqrt{2\pi D}} (\beta^2 + 1) \left( \frac{1}{\sqrt{z_\alpha}} - \frac{1}{\sqrt{\bar{z}_\alpha}} - \frac{1}{\sqrt{z_\beta}} + \frac{1}{\sqrt{\bar{z}_\beta}} \right), \tag{53}$$

with

$$D = 4\alpha\beta - (1 + \beta^2)^2. \tag{54}$$

Note that the Rayleigh wave speed is in fact the root of $D = 0$, when Eqs. (35) and (36) are used in (54). The most important physical feature of Eqs. (51)–(53) is the overall scale of the stress singularity, which is characterized by the Mode I stress intensity factor which, at $y = 0$, is given by

$$K_I = \lim_{x \to 0^+} \sqrt{2\pi x} \sigma_{yy}, \tag{55}$$

FIGURE 7.2. Behind its tip, a fracture is pulled apart by two stresses (after Fineberg and Marder, 1999).

which, as will be shown below, is directly related to energy flux into a fracture tip. Moreover, Eqs. (51)–(53) contain information about the angular structure of the stress fields which can be used in both theoretical and experimental analyses. Theoretically, one can use these equations for predicting the direction of fracture motion, and the conditions under which a fracture branches out. Experimentally, one can utilize these equations for assessing the accuracy of the predictions of continuum fracture mechanics, and for obtaining measurements of the stress fields surrounding rapidly-propagating fractures; we will discuss these matters later in this chapter. It is important to recognize, as pointed out by Freund (1990), that although Eqs. (51)–(53) were derived for fractures moving at a constant speed, the same equations are also true for those that, during propagation, accelerate and/or decelerate, so long as the derivative $dv/dt$ is small during the time needed for sound to travel across the region of the universal elastic singularity.

We now suppose that a fracture is loaded by two stresses, located a distance $l_0$ behind its tip, moving with it in steady state at velocity $v$, and of strength $-\sigma_c$ (see Figure 7.2) such that

$$\lim_{y \to 0^+} \sigma_{yy}(x, y) = -\sigma_c \delta(x + l_0), \quad x < 0. \tag{56}$$

If the fracture tip is at the origin, the stress and displacement fields are continuous and differentiable everywhere, except along a branch cut starting at the origin and running backwards along the negative $x$-axis. If we define $\Phi_\pm(x)$ and $\Psi_\pm(x)$ by

$$\Phi_\pm(x) \equiv \lim_{y \to 0^+} \Phi(x \pm iy), \quad \Psi_\pm(x) = \Psi_\pm(x \pm iy), \tag{57}$$

then because of the branch cut, for $x < 0$, $\Phi_+(x) = -\Phi_-(x)$. As shown above, for Mode I loading, $\sigma_{xy} = 0$ for $y \to 0^+$ and $\forall x$. Therefore, from Eq. (43) we obtain

$$2i\alpha(\Phi_+ - \bar{\Phi}_-) = (\beta^2 + 1)(\Psi_+ + \bar{\Psi}_-), \tag{58}$$

using the fact that $\overline{\Psi(x + i\varepsilon)} = \bar{\Psi}(x - i\varepsilon)$. The function

$$f_+(x) = 2i\alpha\Phi_+(x) - (1 + \beta^2)\Psi_+(x) \tag{59}$$

is defined for all $x$, and can be analytically continued above the $x$-axis, where it is related to stresses and must be free of singularities, whereas $f_-$, defined in a manner similar to $f_+$, must contain no singularities below the real axis. If two complex functions are equal, one without singularities for $x > 0$ and the other without singularities for $x < 0$, the two functions must individually equal a constant which, in fact, is zero since all the stresses are zero far from the fracture. Therefore, $f_+ = f_- = 0$, and one has

$$2i\Phi_+(x) = (1 + \beta^2)\Psi_+(x), \quad 2i\bar{\Phi}_-(x) = (1 + \beta^2)\bar{\Psi}_-(x). \tag{60}$$

The boundary condition for $\sigma_{yy}$ for $x < 0$ is [see Eq. (42)]

$$\sigma_{yy} = -\mu(1 + \beta^2)(\Phi_+ + \bar{\Phi}_-) - 2i\beta\mu(\Psi_+ - \bar{\Psi}_-) = -\sigma_c\delta(x + l_0). \tag{61}$$

Using Eqs. (51)–(54), Eq. (61) becomes

$$\sigma_{yy} = -\sigma_c\delta(x + l_0) = \mu D(\Psi_+ - \bar{\Psi}_-)/(2i\alpha). \tag{62}$$

Since the delta function can be represented as

$$\delta(x + l_0) = \frac{i}{\pi}\frac{1}{x + l_0 + i\varepsilon} \tag{63}$$

one can argue that the only complex function that decays properly at infinity, has a singularity no worse than a square root at the origin, and satisfies Eq. (61), is

$$\Psi_+(x) = \frac{i\alpha}{\pi\mu D}\frac{\sigma_c}{x + l_0 + i\varepsilon}\sqrt{\frac{l_0}{x}}. \tag{64}$$

The function $\Psi(z)$ can now be obtained by analytical continuation of $\Psi_+(x)$. In particular, for $x > 0$ the stress $\sigma_{yy}$ is easily found from Eq. (62) to be

$$\sigma_{yy} = \frac{1}{\pi}\sqrt{\frac{l_0}{x}}\frac{\sigma_c}{x + l_0}, \tag{65}$$

which means that the stress intensity factor associated with $\sigma_{yy}$ is given by

$$K_I = \sigma_c\sqrt{2/\pi l_0}. \tag{66}$$

## 7.4  The Onset of Fracture Propagation: Griffith's Criterion

What are the conditions under which a fracture propagates? Calculations such as those outlined above yield the value of the stress fields at the tip of a propagating fracture, but have nothing to say about the conditions under which a fracture actually propagates. As already discussed in Chapter 6, Griffith (1920) proposed that fracture occurs when the energy per unit area released by a small extension of a

crack is equal to $\Gamma$, the energy required for creating new fracture surface. Griffith's idea, which is the final assumption of continuum fracture mechanics, states that the dynamics of a fracture tip depends only on the total energy flux $\mathcal{H}$ per unit area into the cohesive zone, and that all the details about the spatial structure of the stress fields are irrelevant. The energy $\mathcal{H}$ creates new fracture surfaces, and is also dissipated near the fracture tip. In general, the fracture velocity $v$ is a function of $\mathcal{H}$. It is common to use $\Gamma(v)$ for representing the energy consumed by a fracture in the cohesive zone, in which case the equation of motion for a fracture is

$$\mathcal{H} = \Gamma(v). \tag{67}$$

The central question of interest to continuum fracture mechanics is the conditions under which a static fracture begins to move. For this to happen, a critical fracture energy $\mathcal{H}_c$, the minimum energy per unit area needed for a fracture to propagate forward, irrespective of its velocity, is needed. The standard assumption is that the velocity consuming the minimum energy is very small, although this assumption is not necessarily correct. Equivalently [see Eq. (10) and Chapter 6], one may define a critical stress intensity factor $K_{Ic}$ at which the fracture first begins to propagate. We now derive this equivalence, following Fineberg and Marder (1999).

In what follows, we adopt the summation convention for repeated indices. Energy flux is found from the time derivative of the total energy:

$$\frac{d}{dt}(\mathcal{H}_k + \mathcal{H}_p) = \frac{d}{dt} \int \int \left( \frac{1}{2}\rho\dot{u}_\alpha\dot{u}_\alpha + \frac{1}{2}\frac{\partial u_\alpha}{\partial x_\beta}\sigma_{\alpha\beta} \right) dxdy, \tag{68}$$

where $\mathcal{H}_k$ and $\mathcal{H}_p$ are, respectively, the total kinetic and potential energies within the entire system, and $\dot{u}_\alpha = du_\alpha/dt$. Since the spatial integral in Eq. (68) is taken over a region which is static in the laboratory frame (i.e., $dx/dt = dy/dt = 0$), we have

$$\frac{d}{dt}(\mathcal{H}_k + \mathcal{H}_p) = \int \int \left( \rho\ddot{u}_\alpha\dot{u}_\alpha + \frac{\partial \dot{u}_\alpha}{\partial x_\beta}\sigma_{\alpha\beta} \right) dxdy, \tag{69}$$

where the symmetry of the stress tensor under interchange of indices has been used for the last term. Use of the equation of motion, $\rho\ddot{u}_\alpha = \partial\sigma_{\alpha\beta}/\partial x_\beta$, in Eq. (69), yields,

$$\int \int \left( \frac{\partial}{\partial x_\beta}\sigma_{\alpha\beta}\dot{u}_\alpha + \frac{\partial \dot{u}_\alpha}{\partial x_\beta}\sigma_{\alpha\beta} \right) dxdy = \int \int \frac{\partial}{\partial x_\beta}\left( \sigma_{\alpha\beta}\dot{u}_\alpha \right) dxdy$$

$$= \int_{\partial S} \dot{u}_\alpha\sigma_{\alpha\beta}n_\beta dS, \tag{70}$$

where $\partial S$ is the surface boundary of the system, and $\mathbf{n}$ is an outward unit normal with components $n_\beta$. Equation (70) is a statement of the fact that energy is transported by a flux vector $\mathbf{j}$ with components that are given by

$$j_\alpha = \sigma_{\alpha\beta}\dot{u}_\beta. \tag{71}$$

As mentioned in Chapter 5 [see Section 5.1.1 and Eq. (5.13)], the total energy flux $J$ per unit time into the fracture tip is called the $J$-integral (see Cotterell

FIGURE 7.3. Dotted lines show the most convenient contour for integrating the energy flux and calculating the energy that flows to a fracture tip. The contour runs below the fracture, closes at infinity, and comes back just above the contour (after Fineberg and Marder, 1999).

and Atkins, 1996, for a discussion of the use of the $J$-integral to ductile fracture). A convenient contour for the integration is shown in Figure 7.3. If, for a crack loaded in pure Mode I, we use the asymptotic forms, Eq. (52) for $\sigma_{yy}$ and the corresponding expression for $u_y$, we find that $J$ is given by

$$J = \frac{\alpha}{2\mu} \frac{v(1 - \beta^2)}{4\alpha\beta - (1 + \beta^2)^2} K_I^2, \tag{72}$$

where $K_I$ is the stress intensity factor defined by Eq. (55), with the subscript $I$ emphasizing that the result is specific to Mode I fracture. Thus, the energy release rate $\mathcal{H}$ in the case of pure Mode I is

$$\mathcal{H} = \frac{J}{v} = \frac{\alpha}{2\mu} \frac{1 - \beta^2}{4\alpha\beta - (1 + \beta^2)^2} K_I^2 \equiv \frac{1 - v_p^2}{Y} A_I(v) K_I^2. \tag{73}$$

The corresponding result for pure Mode II fracture is

$$\mathcal{H} = \frac{\beta}{2\mu} \frac{1 - \beta^2}{4\alpha\beta - (1 + \beta^2)^2} K_{II}^2 \equiv \frac{1 - v_p^2}{Y} A_{II}(v) K_{II}^2, \tag{74}$$

while for Mode III fracture one has

$$\mathcal{H} = \frac{v}{2\alpha\mu} K_{III}^2. \tag{75}$$

In the limit $v \to 0$, each of the functions $A_\alpha(v) \to 1$ ($\alpha =$I, II and III) and, for example, Eq. (73) simplifies to

$$\mathcal{H} = \frac{1 - v_p^2}{Y} K_I^2. \tag{76}$$

In the general case of mixed mode fracture, Eq. (10) should be used.

The functions $A_\alpha(v)$ are universal in the sense that they are independent of most details of the material's loading or geometric configuration. Assuming that there is no energy sink in the system other than the one at the tip of the fracture,

Eqs. (73)–(75) relate the total flux of energy from the entire elastic material to the tip which, when it is set to equal to the energy dissipated in the cohesive zone, yields an equation of motion for the fracture. Note that, in order to derive Eqs. (73)–(75), we have tacitly assumed that, given near-field descriptions of stress and displacement fields [Eqs. (50)–(53)], Eqs. (37) and (38) are valid. If, for example, the cohesive zone is of the order of 1 mm in a piece of a solid with dimensions that are a few centimeters, the value of the stress field on the contour $\partial S$ used in Eq. (70) will not be approximated well by the asymptotic forms of the stress and displacement fields, invalidating Eqs. (73)–(75). Since an energy balance provides no information about a fracture's path, we have assumed that the fracture travels along a straight line (see below). Although the rules for determining paths of slowly propagating fractures are known, they are not known for rapidly moving fractures.

## 7.5 The Equation of Motion for a Fracture in an Infinite Plate

As discussed above, one can derive an equation of motion for a fracture by calculating either the energy release rate $\mathcal{H}$ or, equivalently, the dynamic stress intensity factor $K$ which depends on the fracture's loading history, and its length and velocity. In what follows we derive an exact expression for $K$ for a straight semi-infinite fracture in an infinite plate with loads applied to the fracture's faces. The derivation follows closely those given by Willis (1990) and Fineberg and Marder (1999). The calculation is, in the context of linear elasticity—a boundary-value problem—and in the most general case is applicable to a system in which,

(1) the fracture is a semi-infinite straight-line branch cut in an infinite isotropic 2D elastic plate.
(2) The velocity $v(t)$ of the is not in general constant, with the position of its tip being $l(t) = \int_0^t v(t')dt'$, which, in the context of a boundary-value problem, is assumed to be known. However, $v(t)$ must be less than the relevant sound speeds at all times.
(3) The external stresses $\sigma_e$ are permitted only *along* the fracture, but are allowed arbitrary time and space dependence. This can be realized by placing wedges between the faces of the cracks in order to load them.

We derive the corresponding expressions for fracture Modes I, II and III. In the calculations that follow $u$, $\sigma$ and $c$ denote the displacement, stress, and a sound speed in each case, as listed in Table 7.2. By symmetry, $u(x, t) = 0$ for all $x > l(t)$. Due to the one-to-one relation between $K$ and the energy flux $\mathcal{H}$, we compute the latter as a function of $l(t)$ and $v(t)$, and as a functional of the external load $\sigma_e(x, t)$. We look for a Green function $G$ operating on the displacement field $\mathbf{u}$ defined by the following convolution integral,

$$G * u \equiv \int \int G(x - x', t - t')u(x', t') \, dx dt. \tag{77}$$

TABLE 7.2. Notation convention for the solution of equation of motion for a fracture in an infinite plate. The Rayleigh wave speed $c_R$ is the roots of $D = 0$ [see Eq. (54)], and is typically about 90% of the transverse wave speed $c_t$.

|  | Mode I | Mode II | Mode III |
|---|---|---|---|
| $u$ denotes | $u_y(x, y = 0^+, t)$ | $u_x(x, y = 0+, t)$ | $u_z(x, y = 0^+, t)$ |
| $\sigma$ denotes | $\sigma_{yy}(x, y = 0^+, t)$ | $\sigma_{yx}(x, y = 0^+, t)$ | $\sigma_{yz}(x, y = 0^+, t)$ |
| $c$ denotes | $c_R$ | $c_R$ | $c_t$ |

If

$$G(k, \omega) = \int \int e^{ikx' - i\omega t'} G(x', t') \, dx dt, \tag{78}$$

denotes the Fourier–Laplace transform of $G$, we require that

$$G(k, \omega) \equiv \frac{G^-(k, \omega)}{G^+(k, \omega)}, \tag{79}$$

with the properties that $G^+$ and $G^-$ vanish for $x < c_R t$ and $x > -c_R t$, respectively, where $c_R$ is the Rayleigh wave speed. Physically, this implies that $G^+$ is non-zero only for $x$ large enough that a pulse beginning at the origin at $t = 0$ could never reach it in the forward direction (with a similar condition for $G^-$). In fact, for the cases to be discussed below, we have

$$G^+ \propto \delta(x - c_R t), \quad G^- \propto \delta(x + c_R t). \tag{80}$$

While it is not yet clear that $G$ can be decomposed according to Eq. (79), or that it even exists, for the moment we simply assume these to be true.

We decompose $\sigma$ into two functions, $\sigma = \sigma^+ + \sigma,^-$ and define $u = u,^-$ where $\sigma^+$ vanishes for $x < l(t), \sigma^-$ vanishes for $x > l(t)$, while $u^-$ does so for $x > l(t)$. Therefore, $\sigma^-$ describes the stresses along the fracture faces, while $\sigma^+$ is an, as yet unknown, function. $u,^-$ on the other hand, is an unknown function along the fracture faces and vanishes ahead of its tip. Using Eqs. (77)–(79), we write, $G * u = \sigma$, which after Laplace–Fourier transforming yields,

$$G(k, \omega) \, u(k, \omega) = \sigma(k, \omega). \tag{81}$$

Using Eq. (79) we obtain

$$G^-(k, \omega) \, u(k, \omega) = G^+(k, \omega) \, \sigma(k, \omega) \tag{82}$$

which, after inverting back to real space, yields

$$G^+ * \sigma = G^- * u. \tag{83}$$

One can show that for $x < l(t)$, $G^+ * \sigma^+ = 0$. Suppose that $x > l(t)$. Since $\sigma^+$ is zero behind the fracture, the integral

$$G^+ * u^+ = \int \int G^+(x - x', t - t')\sigma^+(x', t') \, dx dt \tag{84}$$

is zero for $x' < l(t')$. The only case for the integrand to be non-zero is for $x' > l(t')$,

in which case, $x' - x > l(t') - l(t) = \dot{l}(t^*)(t' - t)$, where $t < t^* < t'$. However, this means, by the mean-value theorem, that

$$x - x' < c_R(t - t'), \tag{85}$$

since $c_R$ is the largest value that $v(t)$ can take on. On the other hand, (85) is precisely the condition under which $G^+(x - x', t - t')$ vanishes. Therefore,

$$\int\int G^+(x - x', t - t')\sigma^+(x', t')\,dxdt = 0, \quad x < l(t), \tag{86}$$

and similarly

$$\int\int G^-(x - x', t - t')u^-(x', t')\,dxdt = 0, \quad x > l(t). \tag{87}$$

From Eq. (83) one can show that

$$G^+ * \sigma^+ = -(G^+ * \sigma^-)\,H(x, t), \tag{88}$$

where $H(x, t) = \Theta[x - l(t)]$, and $\Theta$ is the Heaviside step function. Equation (88), which has now been shown to be true both for $x > l(t)$ and $x < l(t)$, yields, after inverting it back to real space,

$$\sigma^+ = -(G^+)^{-1} * [(G^+ * \sigma^-)H]. \tag{89}$$

Since $\sigma^-$ is the (known) stress exerted at the back of the fracture tip, Eq. (89) provides a formal solution to the problem. The stress intensity factor is given by

$$K = \lim_{\varepsilon\to 0^+} \sqrt{2\pi\varepsilon}\,\sigma(\varepsilon + l,\ t), \tag{90}$$

which requires identifying the terms that lead to a divergence of the form $1/\sqrt{\varepsilon}$ as $x = l(t) + \varepsilon$ approaches $l(t)$ from above.

We now show that $(G^+)^{-1}$ has a singularity for $\varepsilon \to 0$, behaving as $1/\sqrt{\varepsilon^3}$, while $G^+ * \sigma^+$ is finite. To find the singularity of Eq. (89), $G^+ * \sigma^-$ is evaluated at $x = l(t)$ and pulled outside the convolution as a multiplicative factor. The stress intensity factor can therefore be written as

$$K = \tilde{K}[l(t), \sigma] \cdot \mathcal{K}(v), \tag{91}$$

with

$$\tilde{K}(l, \sigma) \equiv -(\sqrt{2}G^+ * \sigma^-)_{(l,t)}, \tag{92}$$

and

$$\mathcal{K}(v) \equiv \lim_{\varepsilon\to 0^-} [\sqrt{\pi\varepsilon}(G^+)^{-1} * H]_{(l+\varepsilon,t)} \tag{93}$$

Physically, $\tilde{K}(l, \sigma)$, which is independent of the fracture velocity, is the stress intensity factor that would emerge at the tip of a static fracture sitting at all times at $l$ [the tip is exposed to the load $\sigma^-(t)$]. On the other hand, although $\mathcal{K}$ depends on the instantaneous velocity $v(t)$ of the fracture, it is independent of the crack's *history*, i.e., how it arrived at a particular position at time $t$.

## 7.5.1  Mode III

We now apply the general results, Eqs. (91)–(93), to the particular case of anti-plane shear, which will also allow us to verify the general structure of the Green function $G$ used so far. By calculating the stress intensity factor using Eqs. (73)–(75), and hence the energy release rate $\mathcal{H}$, we derive the equation of motion for a Mode III fracture by equating the energy release rate to the fracture energy. The starting point is the wave equation for $u_z$, Eq. (14), which after Fourier transforming in both space and time yields

$$\frac{\partial^2 u_z}{\partial y^2} = (k^2 - \omega^2/c^2 - 2ib\omega)u_z, \tag{94}$$

where a small damping $b$ has been added to help us overcome some convergence problems that will arise later. In an infinite plane, the only allowed solution is one that decays as a function of $y$, and therefore Eq. (94) is solved by

$$u_z(k, y, \omega) = \exp\left[-y\sqrt{k^2 - \omega^2/c^2 - 2ib\omega}\right] u(k, \omega). \tag{95}$$

By taking $u = u_z(y = 0)$ and $\sigma = \sigma_{yz}(y = 0)$, one has

$$G(k, \omega) = \frac{\sigma}{u} = -\mu\sqrt{k^2 - \omega^2/c^2 - 2ib\omega}. \tag{96}$$

Using Eq. (79) we can write

$$G^- = -\mu\sqrt{ik - i\omega/c + b}, \tag{97}$$

and

$$G^+ = 1/\sqrt{-ik - i\omega/c + b}. \tag{98}$$

The decomposition, $G = G^-/G^+$, satisfies the conditions of the preceding section if we write

$$G^+(x, t) = \frac{1}{(2\pi)^2} \int \int \frac{e^{-ikx - i\omega t}}{\sqrt{-ik - i\omega/c + b}} dk \, d\omega$$
$$= \frac{1}{(2\pi)^2} \int \int \frac{e^{-ipx - i\omega(t - x/c)}}{\sqrt{-ip + b}} dp \, d\omega, \tag{99}$$

with $p = k + \omega/c$, and therefore

$$G^+(x, t) = \frac{\delta(t - x/c)}{2\pi} \int \frac{e^{-ipx}}{\sqrt{-ip + b}} dp. \tag{100}$$

When $x < 0$, one must close the contour in the upper half plane, and since the branch cut is in the lower half plane, the integral vanishes. When $x > 0$, we deform the contour to surround the branch cut to obtain

$$\frac{1}{2\pi} \int_0^\infty \frac{2e^{-px}}{\sqrt{p + b}} dp = \frac{1}{\sqrt{\pi x}}. \tag{101}$$

Therefore

$$G^+(x, t) = \frac{1}{\sqrt{\pi x}}\delta(t - x/c)\Theta(x). \tag{102}$$

By a largely similar analysis we find that

$$(G^+)^{-1}(x, t) = \delta(t - x/c)\frac{d}{dx}\left[\frac{\Theta(x)}{\sqrt{\pi x}}\right]. \tag{103}$$

Having calculated $(G^+)^{-1}(x, t)$, we can now find the stress intensity factor $K_{III}(l, t)$. From Eq. (93) we find that

$$\mathcal{K}(v) = \sqrt{\pi\varepsilon}\int\int\delta(t_1 - x/c)\left[\frac{d}{dx_1}\frac{\Theta(x_1)}{\sqrt{\pi x_1}}\right]\Theta[l(t) + \varepsilon - x_1 - l(t - t_1)]\,dx_1 dt_1$$

$$= \sqrt{\varepsilon}\int\left[\frac{d}{dx_1}\frac{\Theta(x_1)}{\sqrt{\pi x_1}}\right]\Theta[\varepsilon/(1 - v/c) - x_1]\,dx_1. \tag{104}$$

Since only very small $x_1$ are important, we find that

$$\mathcal{K}(v) = \sqrt{\varepsilon}\int\frac{\Theta(x_1)}{\sqrt{\pi x_1}}\delta[\varepsilon/(1 - v/c) - x_1]\,dx_1 = \sqrt{1 - v/c}. \tag{105}$$

Similarly,

$$\tilde{K}(l, t) = -\sqrt{2}\int\int\delta(t_1 - x_1/c)\frac{\Theta(x_1)}{\sqrt{\pi x_1}}\sigma[-l(t) - x_1, t - t_1]\,dx_1 dt_1$$

$$= -\sqrt{2}\int\frac{\Theta(x_1)}{\sqrt{\pi x_1}}\sigma[l(t) - x_1, t - x_1/c]dx_1. \tag{106}$$

In particular, when $\sigma^-$ does not depend on the time and, $\sigma(x) = \sigma_0\Theta(x)$, one obtains

$$\tilde{K} = -(4/\sqrt{2\pi})\sigma_0\sqrt{l}. \tag{107}$$

The minus sign arises because the stresses ahead of the fracture tip always act against those applied on the fracture faces. Note that Eq. (104) reduces to unity when $v \to 0$, implying that in the case of time-independent loading, $\tilde{K}$ is indeed the stress intensity factor one would have had if the fracture had been static at $l$ for all times. For the propagating fracture, we obtain

$$K_{III} = \sqrt{1 - v/c}\,\tilde{K}[l(t), \sigma_0]. \tag{108}$$

We now compute the stress singularity that would have developed had we had a static fracture of length $l(t)$ at time $t$, and multiply the result by a function of the instantaneous velocity. We should emphasize that all details of the history of the crack motion are irrelevant; only the velocity and loading configuration are needed for determining the stress fields sufficiently close to the tip. As a consequence, one can use Eq. (75) to determine the energy flow to the tip of the crack:

$$\mathcal{H} = v(1 - v/c)\frac{\tilde{K}^2}{2\alpha\mu}. \tag{109}$$

The rate at which energy enters the tip of the fracture must be equal to $v\Gamma(v)$. There is nothing to prevent the fracture energy $\Gamma$ from being a function of the velocity, but the notion of local equilibrium, which has prevailed until now, strongly suggests that $\Gamma$ should be a function of $v$ alone. Therefore

$$\Gamma(v) = (1 - v/c)\frac{\tilde{K}^2(l)}{2\alpha\mu}, \tag{110}$$

which, after rearranging and using Eq. (106), yields

$$\frac{\pi\mu\Gamma}{4l\sigma_0^2} = \sqrt{(1 - v/c)(1 + v/c)}. \tag{111}$$

If we define

$$l_0 = \frac{\pi\mu\Gamma}{4\sigma_0^2}, \tag{112}$$

Eq. (111) is rewritten as

$$\frac{l_0}{l} = \sqrt{(1 - v/c)(1 + v/c)}. \tag{113}$$

## 7.5.2   Mode I

The preceding analysis can also be carried out for thin plates under tension. Although all steps of the analysis proceed as before, it is not possible to derive simple analytical expressions. This case has been discussed in detail by Freund (1990) who finds that the energy flux per unit length extension of the fracture, to an accurate approximation, is given by

$$\mathcal{H}(v) = \Gamma(v) = \frac{(1 - v/c_R)\tilde{K}^2(l)}{2\tilde{\lambda}}, \tag{114}$$

where $\tilde{\lambda}$ is a Lamé constant defined by Eq. (20). Rearranging Eq. (114) yields

$$\frac{Y\Gamma(v)}{\tilde{K}^2(l)(1 - vv_p^2)} = 1 - \frac{v}{c_R}, \tag{115}$$

where $\tilde{K}$ is still given by Eq. (106), using $\sigma_{yy}$ on the $x$-axis for $\sigma$. In the case of time-independent loading described by $\sigma(x) = \sigma_0\Theta(x)$, one obtains

$$\frac{l_0}{l} = 1 - \frac{v}{c_R}, \tag{116}$$

with

$$l_0 = \frac{\pi\Gamma\tilde{\lambda}}{4\sigma_0^2}. \tag{117}$$

Equation (116) is now written in the following form

$$v = c_R(1 - l_0/l), \tag{118}$$

which is nothing but Eq. (7) obtained by the scaling analysis, with the difference that $\Gamma$ and hence $l_0$ can depend strongly upon the crack velocity $v$. Hence, seemingly large differences between the predictions of the theory and experimental data are due to nothing more than assuming that $l_0$ is a constant!

What are the practical implications of Eq. (10) for the design of experiments? As discussed by Fineberg and Marder (1999), one may consider three experimental situations. (1) One for which the assumptions of the theory hold well. (2) A second experiment in which the theoretical assumptions are satisfied in an approximate way, while (3) in the third experiment the assumptions clearly fail. The three cases are as follows.

(1) A thin plate has a fracture running half-way through, and driven by wedging action in the middle. For times less than that needed for sound to travel from the point of loading to the material's boundaries and back to the tip of the fracture, all the assumptions of the theory are satisfied.

(2) A thin plate has a long fracture as before, but with uniform static stresses $\sigma_\infty$ applied at the outer boundaries, and the faces of the fracture being stress-free. This problem is equivalent to one in which the upper and lower outer boundaries are stress-free, but uniform stresses $-\sigma_\infty$ are applied along the fracture faces. The equivalence is due to the fact that an uncracked plate under uniform tension $\sigma_\infty$ is a solution of the equations of elasticity, so this trivial static solution can be subtracted from the first problem to obtain the second equivalent one. However, in the new problem, stresses are applied to the fracture faces all the way back to the left-hand boundary of the material. Therefore, the problem must be mapped onto one in which stresses are applied to the faces of a semi-infinite fracture in an infinite plate, but the correspondence is only approximate.

(3) Consider now a semi-infinite fracture in an infinitely long strip, shown in Figure 7.4, which is loaded by displacing each of its boundaries at $y = \pm w/2$ by a constant amount $\delta$. Far behind the tip (as $x \to -\infty$), the fracture relieves all the stresses within the strip. Far ahead of the tip (as $x \to +\infty$), the material is unaffected by the fracture with the stress field being linear in $y$. Thus, the energy per unit extension far ahead of the fracture has a constant value, $2Y\delta^2/[w(1 - v_p^2)]$, where $Y$ and $v_p$ are, respectively, the usual Young's modulus and Poisson's ratio of the material. The translational invariance of the system along the $x$-direction implies that, for a given $\delta$, the fracture should



Semi-infinite strip

FIGURE 7.4. A fracture in a semi-infinite strip.

eventually propagate at a constant velocity $v$. Writing an energy balance yields

$$\mathcal{H} = \Gamma = \frac{2y\delta^2}{w(1 - v_p^2)}. \tag{119}$$

If we now assume that we can still use Eq. (114), then, since the stress intensity factor $\tilde{K}$ of a static fracture in a strip loaded with constant displacements $\delta$ cannot depend upon where the fracture is located, $\tilde{K}$ must be a constant and Eq. (114) would predict that

$$\mathcal{H} = \Gamma = \frac{(1 - v_p^2)\tilde{K}^2}{Y}(1 - v/c_R). \tag{120}$$

However, the velocity term of Eq. (120) contradicts Eq. (119), implying that Eq. (114) has failed. The reason for this failure lies in the assumption that the fracture tip does not feel the presence of the system's boundaries, which is clearly invalid. In fact, the translational invariance of the system depends crucially on the presence of its vertical boundaries. Energy flows continuously into the material as the amount of kinetic energy reaches a steady state. In contrast, the kinetic energy within a system of infinite extent increases without bounds as ever farther reaches of the material feel propagation of the fracture, since elastic waves that carry this information propagate outward.

## 7.6    The Path of a Fracture

We now discuss briefly the path travelled by a propagating fracture. As discussed above, energy balance provides an equation of motion for the tip of a fracture only when its path or direction of its propagation is *assumed*. Although criteria for the path of a slowly propagating fracture have been established, no such criteria have been proven to exist for a rapidly moving fracture. We will discuss this issue later in this chapter.

### 7.6.1    Planar Quasi-static Fractures: Principle of Local Symmetry

A fracture is considered to propagate slowly if the velocity $v$ of its tip is much less than the Rayleigh wave speed $c_R$. Goldstein and Salganik (1974) proposed that the path taken by slow cracks satisfies the *principle of local symmetry*, according to which a crack propagates so as to set the component of Mode II loading to zero. An immediate consequence of this proposal is that if a stationary fracture is loaded in such a way as to experience Mode II loading, it forms, upon extension, a sharp kink and moves at a new angle. This rule means that the fracture moves perpendicular to the direction in which tensile stresses are *maximum*. Cotterell and Rice (1980) showed that a fracture satisfying the principle of local symmetry also chooses a direction so as to maximize the rate of energy release. The distance over which a fracture must move so as to set $K_{II}$ to zero is of the order of the size of the cohesive zone (Hodgdon and Sethna, 1993). Cotterell and Rice (1980) also

showed that the condition $K_{II} = 0$ has the following consequences for fracture propagation. Consider an initially straight fracture, propagating along the $x$-axis. The components $\sigma_{xx}$ and $\sigma_{yy}$ of the stress field are given by

$$\sigma_{xx} = \frac{K_I}{2\pi r^{1/2}} + \sigma + O(r^{1/2}), \tag{121}$$

$$\sigma_{yy} = \frac{K_I}{2\pi r^{1/2}} + O(r^{1/2}). \tag{122}$$

The constant stress $\sigma$ is parallel to the fracture at its tip. If $\sigma > 0$, any small deviations from straightness cause the fracture to diverge from the $x$-direction, whereas if $\sigma < 0$ the fracture is stable and continues to propagate along the $x$-axis. Yuse and Sano (1993) and Ronsin *et al.* (1995) conducted experiments described in Section 6.12 by slowly pulling a glass plate from a hot region to a cold one across a constant thermal gradient. The velocity of the fracture, driven by the stresses induced by the non-uniform thermal expansion of the material, follows that of the glass plate. At a critical pulling velocity, the fracture's path deviated from straight-line propagation and developed transverse oscillations. This instability is completely consistent with the principle of local symmetry: The crack deviates from a straight path if the stress $\sigma$ in Eq. (121) is positive. The wavelength of the ensuing oscillations has also been computed numerically (Adda-Bedia and Pomeau, 1995).

### 7.6.2  *Three-Dimensional Quasi-static Fractures*

Hodgdon and Sethna (1993) generalized the principle of local symmetry to 3D and showed that an equation of motion for a crack line involves, in principle, nine different constants, although we are not aware of any experimental determination of these constants. Larralde and Ball (1995) and Ball and Larralde (1995) carried out stability analysis of cracks which are *almost* planar and have a tip which is *almost* a straight line, so that the differences from a planar and straight edged crack could be considered as perturbation parameters. They showed that, in agreement with the proposal of Goldstein and Salganik (1974) (i.e., the principle of local symmetry), at least under quasi-static conditions the mechanism underlying the stability of planar fractures propagating under Mode I loading is the appearance of a Mode II loading in the vicinity of the fracture edge associated with each out-of-the plane perturbation mode. This Mode II loading, which tends to suppress the perturbation in the quasi-static propagation of the fracture, is a consequence of the global structure of the crack edge, and has nothing to do with the prior history of the fracture, nor with the local geometry around any point of the edge. As such, this stabilizing mechanism is an intrinsically *three-dimensional* effect.

Thus, the principle of local symmetry is consistent with all experimental tests that have been performed so far on slowly propagating fractures. Nevertheless, it is not based on a rigorous theoretical foundation, since there is no basic principle that predicts that a fracture must extend perpendicular to the maximum tensile stress, or that it must maximize energy release.

## 7.6.3   Dynamic Fractures: Yoffe's Criterion

In the case of rapid fractures, there is no rigorous basis for deciding the direction in which a fracture may propagate. A variety of criteria for path selection have been proposed in the literature which can be divided into two types: Those proposing that a crack propagates in the direction of a maximal stress, and those that are based on a maximum dissipation of energy. In contrast to quasi-static fractures, however, these criteria are not equivalent and, more importantly, none of them is strongly supported by experiment.

An important work is that of Yoffe (1951), already mentioned in Chapter 6, who proposed that one should check the stability of a rapidly propagating fracture by examining the dynamic stress fields given by Eqs. (51)–(53), approaching the tip of the fracture along a line at an angle $\theta$ relative to the $x$-axis, and computing the stress perpendicular to that line. If we choose, $z_\alpha = r \cos \theta + i r \alpha \sin \theta$ and $z_\beta = r \cos \theta + i r \beta \sin \beta$, and evaluate the stress in the polar coordinates, $\sigma_{\theta\theta} = \sigma_{xx} \sin^2 \theta + \sigma_{yy} \cos^2 \theta - \sigma_{xy} \sin 2\theta$, we find that below a velocity $v_c \simeq 0.61 c_R$ (which depends on the Poisson's ratio), the maximum tensile stress occurs for $\theta = 0$. Above $v_c$, the tensile stress $\sigma_{\theta\theta}$ develops a maximum in a direction $\theta > 0$, and the angle of maximum tensile stress increases smoothly until it finally develops a maximum at about $\pm 60°$ relative to the $x$-axis, implying that above the critical velocity a fracture might propagate *off-axis*. As pointed out by Fineberg and Marder (1999), this spontaneous breaking of the axial symmetry of the phenomenon is due to purely kinetic effects. Recall that in an elastic medium information is propagated at the speed of sound, and that the stress field at the tip of a rapidly moving fracture is analogous to the electric field surrounding a point charge moving at relativistic velocities. The stress field then experiences a Lorentz contraction in the direction of propagation as the fracture's velocity approaches the speed of sound, resulting in the formation of symmetric lobes around the $x$-axis of maximal tensile stress (above the critical velocity).

Fracture branching stemming from the approach of the velocity of a crack tip to Yoffe's critical velocity was first thought to provide a rigorous criterion for crack instability. However, many experiments have shown that large-scale branching occurs in a variety of materials at velocities much less than $0.61 c_R$, and that branching angles of about $10° - 15°$ (see Section 7.8.10 below), instead of Yoffe's predicted value of $60°$, are generally observed. To overcome the failure of the Yoffe's criterion for fracture branching, a number of other criteria have been proposed (see, for example, Ramulu and Kobayashi, 1985, 1986) in which the form of the stress field at the boundary of the cohesive zone near the tip is used for deriving criteria for fracture branching and its angle (see, for example, Theocaris and Georgiadis, 1985; Ramulu and Kobayashi, 1985, 1986). To obtain the angle of branching one determines the direction in which the local energy density, evaluated at the edge of the cohesive zone, is maximum. The theoretical justification for this criterion was originally suggested by Sig (1973) who proposed that fracture propagation occurs in the radial direction along which the local energy density possesses a stationary value. Experimentally-measured crack branching angles are consistent

with those predicted by variants of this criterion, although the same criteria predict critical velocities for crack branching that are nearly identical to Yoffe's prediction, namely, $0.61c_R$. Adda-Bedia and Ben Amar (1996), for example, proposed that one should draw contours of constant principal stress and search for points where these contours are perpendicular to lines drawn from the crack tip, along which the fracture travels. This criterion predicts the existence of *two* critical speeds. The first is the velocity at which the fracture must choose between *three* possible directions, whereas the second critical velocity is one at which the fracture must choose between *five* possible directions. Although this criterion is plausible, there is no experimental evidence indicating that this is in fact the preferred criterion. The branching angles that are predicted by such criteria are *not* significantly different from those determined by the following condition, which is a type of static condition. Consider the stress field formed ahead of a single propagating fracture, from which one can compute the trajectories that satisfy the quasi-static condition, $K_{II} = 0$ (Kalthoff, 1972; Parleton, 1979). The angle that is determined by this trajectory at a distance $r_c$ from the fracture tip, where $r_c$ is the typical size of the cohesive zone, is in good quantitative agreement with the experimental observations.

## 7.7   Comparison with the Experimental Data

It is instructive at this point to compare the predictions of linear continuum fracture mechanics with the experimental data. A close inspection indicates that, as long as the basic assumptions of linear fracture mechanics hold, the theory is quite successful in predicting both fracture propagation and the behavior of the stress field throughout the material. However, if one or more of these assumptions break down, the linear theory loses its predictive power. For example, continuum fracture mechanics has been successful in predicting the value of the stress intensity factor at the tip of both static and dynamic fractures for both static and dynamically-applied loads. Kim (1985) measured transient behavior of the stress intensity factor and made a quantitative comparison with the predictions of Eqs. (91)–(93). In his experiment, a step function loading was applied to the crack faces in a sheet of Homalite-100 that was large enough to be approximated as an infinite system. Homalite-100 is a thermoset polyester resin that, at room temperature, can be accurately represented as a linearly elastic material with brittle fracture behavior. Of particular importance is its property of birefringence that permits the use of optical techniques, such as photoelasticity described in Section 6.11.2.2, for mapping the stress field. Due to such desirable properties, Homalite-100 has been used in many studies of dynamic fracture. The stress intensity factor was measured optically (see Chapter 6), using a method developed by Kim himself in which the relation of the transmitted light through the fracture tip with the stress intensity factor was used. Kim's data agreed well with the calculated time dependence of the stress intensity factor. Similar agreement between the theory and experiments for PMMA was reported by Vu and Kinra (1981) who measured the transient relaxation of the stress

FIGURE 7.5. Comparison of the predictions of linear continuum fracture mechanics (the curve) with the experimental data (symbols) of Kobayashi *et al.* (1974).

field within the material. In their experiments strain gauges, with a temporal resolution of about 1 $\mu$s, were placed throughout the sample to measure the temporal behavior of the stress field surrounding a fracture at times immediately following its arrest. Their data agreed with a prediction of Freund (1990), that the stress field at a point directly ahead (behind) the fracture should reach its equilibrium value (to within a few percent) as soon as the shear (Rayleigh) wave front passes.

However, the same type of favorable comparison between the theory and experiment does not exist at high fracture velocities, and in fact experiments often seem to disagree with Eq. (14). As an example we show in Figure 7.5 the data of Kobayashi *et al.* (1974) with PMMA and compare them with the theoretical predictions. Although the theory predicts that if the fracture energy is not a strong function of the velocity, the fracture would smoothly accelerate from rest to the Rayleigh wave speed $c_R$, Kobayashi *et al.*'s data do not confirm this prediction: After the fracture initially accelerates rapidly, it becomes increasingly sluggish and eventually reaches a final velocity well below $c_R$. However, if we suppose that the fracture energy $\Gamma$ is a function of the velocity, and specify in Eq. (118) [using Eqs. (114)–(116)] that $l_0$ is defined in terms of the minimum energy $\Gamma(0)$ at which fracture propagation first happens, one obtains instead of Eq. (118):

$$v = c_R \left[ 1 - \frac{\Gamma(v)}{\Gamma(0)} \frac{l_0}{l} \right]. \tag{123}$$

Therefore, if the fracture energy $\Gamma(v)$ increases rapidly with the velocity $v$, one can obtain practically *any* functional dependence of the velocity on the fracture length. One can also interpret Eq. (123) as a way of extracting the velocity dependence of fracture energy from measurements of $v$. However, validation of the theory cannot be accomplished without an independent measurement of the fracture energy $\Gamma$, although even such validation would provide no fundamental explanation of the origin of any measured velocity dependence of the fracture energy.

FIGURE 7.6. Experimental data (triangles) of Sharon and Fineberg (1999) and their comparison with the theoretical predictions (rectangles).

Bergkvist's (1974) beautiful experiments on crack arrest in PMMA provided the first comparison of the theory and experiment where the velocity dependence of fracture energy was explicitly taken into account. His experiments allowed direct comparison of the calculated energy release rates with experimental data for fracture velocities below $0.2c_R$ (which are less than 200 m/s). He obtained a continuous distribution of the fracture tip locations with a temporal resolution of about 1 $\mu$s, and used independent measurements of the fracture energy of PMMA as a function of the fracture's velocity. Values of the fracture velocity were computed by equating the measured value of the fracture energy to the calculated value of the energy release rate. The predicted and measured velocities were in agreement to within 10%.

A similar comparison between the theoretical predictions and experimental data for PMMA was reported by Sharon and Fineberg (1999). They first carried out an independent measurement of the fracture energy of a crack by the use of a strip geometry. An additional series of experiments, which was carried out in $40 \times 40$ cm samples, yielded the velocity values which were then inserted into Eq. (123) to yield values of $\Gamma(v)$ which were then compared to the direct measurements. The results are shown in Figure 7.6. Their data agree with Eq. (123) for velocities less than about 400 m/s $\simeq 0.4c_R$. However, above $0.4c_R$ there is a large difference between the data and the predicted values of $\Gamma(v)$, which is due to the growth of the cohesive zone around the crack tip to a length scale where the assumptions of linear continuum fracture mechanics are no longer valid (see also below).

## 7.7.1   The Limiting Velocity of a Fracture

As the derivation of Eq. (118) indicated, an important prediction of linear continuum fracture mechanics is that, disallowing divergent behavior of $\Gamma(v)$, a fracture

should accelerate until it arrives asymptotically at the Rayleigh wave speed $c_R$. However, in amorphous materials, such as PMMA and glass, the maximum measured velocity of a fracture hardly exceeds a value of about $0.5c_R$, whereas in strongly anisotropic materials, such as LiF (Gilman *et al.*, 1958), tungsten (Hull and Beardmore, 1966; Field, 1971), and MgO (Field, 1971), a propagating fracture attains a speed of up to $0.9c_R$, as cleavage through a weak plane takes place, hinting that strong anisotropy in materials may be necessary for the fracture to attain the limiting velocity $c_R$. An interesting experiment by Washabaugh and Knauss (1994) indicated that this may indeed be the case. In their experiment, plates of PMMA were first fractured and then rehealed to form a preferred plane in the material that was substantially weaker than the material on either side of it. Although the interface did weaken the PMMA, the rehealed material still had between 40% and 70% of the strength of the original material. Using an interferometer together with a high-speed rotating mirror camera, interferograms of the fracture tip were recorded at equal time intervals. In this way a fracture velocity of up to $0.9c_R$ was measured. Washabaugh and Knauss (1994) also noted that none of the fractures propagating along the weakened interfaces produced branches beyond the point of fracture initiation. The same type of behavior takes place in strongly anisotropic crystalline materials. Field (1971) noted that in experiments on MgO and rolled tungsten (rolling in the preparation of tungsten induces a preferred orientation in the material, hence making it anisotropic) branching of a fracture is suppressed until very high velocities. Thus, in strongly anisotropic materials, where microscopic crack branching is inhibited, fractures approach the predicted limiting velocity of $c_R$.

Let us mention here that there have been some continuum models of dynamic fracture that predict that a fracture tip may propagate with a speed even *larger than* $c_R$. For example, Langer (1992) investigated three 1D and 2D unsteady-state models of fracture propagation. His 1D models had the following general form

$$\frac{\partial^2 u}{\partial t^2} = \frac{\partial^2 u}{\partial x^2} - \alpha_f^2(u - \delta) - F_c(u) - F, \tag{124}$$

where $u(x, t)$ is the displacement of the material at time $t$ and position $x$ along the face of the fracture, and $\alpha_f^2$ is some sort of a force constant representing a linear elastic coupling between the fracturing material and a fixed substrate. Here, $F_c(u)$ is the cohesive force, i.e., $\int_0^\infty F_c(u)du = \Gamma$ is the fracture energy, and $F$ is a function that depends on $\dot{u} = \partial u/\partial t$. A fully relaxed configuration of the system corresponds to $u = \delta$. In model 1, $F$ was a friction force given by $F = c_1 + 2c_2 u_t$, where $c_1$ and $c_2$ are constant. For this model Langer (1992) [see also Langer and Tang (1991)] found that, if $c_2 < 0$, then slipping friction causes the fracture to accelerate to the limiting wave velocity *independent of loading strength*. The second model was a 2D version of model 1 with the same qualitative behavior. In model 3, $F = -\eta\partial^2 u_t/\partial x^2$, where $\eta$ is a viscosity-like coefficient. It was found that the steady-state solutions at large applied stresses exhibit oscillating fracture-opening displacements which propagate at speeds that are comparable to, or higher than, the nominal wave speed $c_R$, i.e., the fracture propagates *supersonically*. We will come back to this interesting prediction in Section 7.8.15.

## 7.8 Beyond Linear Continuum Fracture Mechanics

Let us now discuss some of the phenomena in rapid fracture that are outside the realm of linear continuum fracture mechanics.

### 7.8.1 The Dissipated Heat

We already described in Sections 6.11.4 and 6.11.6 the technique for measuring the heat generated by a propagating crack in PMMA (Döll, 1973; Zimmerman *et al.* 1984), in glass (Weichert and Schonert, 1974), and in steel (Zimmerman *et al.*, 1984). These experiments indicate that heating accounts for most of the elastic energy that drives the fracture. In particular, for fracture velocities ranging from $0.1c_R$ to $0.6c_R$, the measured heat flux accounts for 50-60% of the energy release, whereas for fractures velocities in the range $0.1c_R - 0.3c_R$ the measured heat flux accounts for virtually the entire energy release.

Although these experiments tell us that almost all of the elastic energy is converted into dissipated heat, a central question is where this dissipation takes place within the material. Is it in, for example, the cohesive zone, or does it occur as elastic waves that propagate away from the fracture are attenuated within the material? Fuller *et al.* (1983) provided the answer to this central question by real-time infrared visualization of the fracture tip during its propagation. Their experiments were carried out on PMMA and polystyrene, and indicated that, in both materials, temperatures at the tip were approximately constant, as a function of the fracture's velocity, with a temperature rise of the order of 500 K (see also Zehnder and Rosakis, 1991; Kallivayalil and Zehnder, 1994, for similar data for AISI 4340 carbon steel and $\beta$-C titanium). These experiments also established that, in addition to the large temperature rise (in PMMA and polystyrene the temperatures at the crack tip were well above the equilibrium melting temperature), the source of the heating was within a few $\mu$m of the crack's path, or well within the cohesive zone, as defined by the material's yield stress, implying that nearly all of the heat dissipation in the material takes place in the vicinity of a fracture. The heat release appears to be caused by the extreme plastic deformation induced by the fracture process in the vicinity of the tip. This result is also supported by the experiments of Kusy and Turner (1975) who investigated the fracture energy of PMMA. They found that the fracture energy of high ($> 10^5$) molecular weight PMMA can be over two orders of magnitude larger than the surface energy (i.e., the energy needed to break a unit area of atomic bonds), which they explained it in terms of plastic deformation of the polymer chains, whereas below a molecular weight of about $10^5$, no significant plastic deformation occurred in the fracture, and its energy was comparable with the surface energy.

The dissipated heat and the associated temperature rise in the material can be calculated. For example, Langer (1993) and Langer and Nakanishi (1993) considered a 2D dynamic fracture model defined by

$$\frac{\partial^2 u}{\partial t^2} = c\nabla^2 u - m^2(u - \delta), \tag{125}$$

where $c$ is a wave speed, $m$ is the mass, and $m^2\delta$ is the applied force. The fracture is essentially moving along the center line of a strip of finite width. The traction applied to the fracture surface was assumed to be given by

$$\mu \frac{\partial u}{\partial y}\bigg|_{y=0} = \sigma_c(u) - \eta \frac{\partial^2 \dot{u}}{\partial x^2}\bigg|_{x=0}, \tag{126}$$

where $\mu$ is an elastic modulus (for example, the shear modulus), and $\sigma_c$ is the cohesive stress acting between the open fracture faces. The cohesive stress was taken to be $\sigma_c = \sigma_y$ for $0 \le u(x, 0, t) \le \delta_c$, and $\sigma_c = 0$, otherwise, where $\sigma_y$ is the yield stress, and $\delta_c$ is obviously the range of the cohesive force; note that $u(x, 0, t)$ is just the fracture-opening displacement. The second term on the right-hand side of (126) is a viscous damping stress which acts on the fracture surface. The two spatial derivatives preserve reflection and translational symmetries, and the time derivative in (126) breaks time-reversal symmetry in order to produce energy dissipation. The most interesting prediction of the model was a relation between the velocity of fracture propagation and the externally applied stress, given by

$$\frac{v/c}{[1 - (v/c)^2]^{3/2}} \simeq \left(\frac{K}{K_e}\right)^{12}, \tag{127}$$

which is valid for $1 \ll K/K_G \ll (w/\delta_c)^{1/6}(\sigma_y/\mu)^{1/6}$, where $w$ is the width of the system, and $K_G = \sqrt{2\sigma_y\delta_c/\mu}$. Here $K$ is the stress-intensity (more precisely, the strain-intensity) factor associated with the applied force, and

$$K_e = (6\delta_c)^{1/3} \left(\frac{\eta c_R}{\mu}\right)^{1/12} \left(\frac{\sigma_y}{\mu}\right)^{2/3}. \tag{128}$$

The surprising aspect of these predictions is the unusual exponent $1/12$. If $K_e > K_G$, then the fracture velocity $v$ jumps from very small values to values near $c_R$ as $K$ passes through $K_e$, and therefore $K_e$ plays the role of an effective Griffith threshold at which the fracture makes a sharp transition from slow motion to rapid propagation. Whether such predictions can be observed in an actual experiment remains an open question. The dissipated energy is $\eta(\partial \dot{u}/\partial x)^2$, and assuming that this energy is converted to heat, then the corresponding temperature rise $\Delta T$ will be

$$\Delta T \simeq \frac{K^2}{3C_p}\sqrt{\frac{\mu^3 v}{\eta}}, \tag{129}$$

where $C_p$ is the specific heat of the material.

How do the thermal effects in the cohesive zone influence our basic understanding of the fracture process? Since the fracture energy is an input into the theory of linear continuum fracture mechanics, neither the large temperature rise observed within the cohesive zone nor its cause(s) have any effect on the predictions of the theory. This is true so long as the heat dissipation is localized within the cohesive zone and does not spread out throughout the material. Otherwise, the entire rationale behind Eqs. (91)–(93) would be invalid. Thus, the total fracture energy

is related to the amount of microscopic surface actually generated by the fracture process, which in turn is related to instabilities that occur to a single fracture as a function of the energy that it dissipates.

## 7.8.2   The Structure of Fracture Surface

Studies of fracture surfaces of amorphous brittle materials indicates that the surface (which is generated by dynamic fracture) has a characteristic structure called *mirror*, *mist*, and *hackle*. These characteristics, which were already described in Section 6.13, have been reported to occur in materials as diverse as glasses and ceramics, non-crosslinked glassy polymers such as PMMA, and crosslinked glassy polymers, such as Homalite-100, polystyrene and epoxies. To summarize our description of the structure of a fracture surface given in Chapter 6, near the location of the fracture onset, the fracture surface appears smooth and shiny, and is thereby called the mirror region. As a crack propagates further, the fracture surface becomes cloudy in appearance, and is referred to as mist. When the fracture surface becomes extremely rough, it is said to be in the hackle region.

## 7.8.3   Topography of Fracture Surface

It is often useful to make quantitative measurements of the topography of a fracture surface, for which there are several techniques, each of which is appropriate for a specific length scale. For length scales ranging from 1 to 100 $\mu$m, commercial contact-type scanning profilometers is used for measuring such properties as the root mean-square of roughness of a surface. However, the tip size of the contact probe limits its resolution in resolving surface features that are under 10 $\mu$m, in which case optical profilometers have been used (Boudet *et al.*, 1995). Fracture surfaces at submicron length scales have recently been studied (Milman, 1994; Daguier *et al.*, 1997), using both scanning tunneling and atomic force microscopy.

## 7.8.4   Properties of Fracture Surface

Analysis of fracture surfaces, usually called *fractography* (Hull, 1999), is concerned with the determination of the location of the onset of fracture of a given material together with the probable cause for its failure. Although every material has its own fracture surface which is different from that of any other material, the proven usefulness of analysis of fracture surface in the determination of different fracture processes stems from the fact that, a close empirical relation exists between the deterministic dynamics of a fracture and the surface that it creates. The mechanisms that give rise to characteristic surface features are, in many cases, not known, but the fact that these features are at all general is strong evidence that they are generated by a deterministic process, independent of details of the loading or the initial conditions of the material under study.

FIGURE 7.7. Typical parabolic markings formed on the fracture surface of PMMA (after Ravi-Chandar and Yang, 1997; courtesy of Professor K. Ravi-Chandar).



## 7.8.5   Conic Markings on Fracture Surface

Fracture surfaces of amorphous materials also contain small conic (or parabolic on a surface) markings in the mist region; see Figure 7.7. They appear in all three fracture regimes, namely, mirror, mist, and hackle, and are the result of microscopic defects opening up ahead of the main fracture front. To see the origin of these markings, suppose that a microscopic void is placed directly ahead of a fracture. The large stress field, generated at the fracture's tip, causes the void to propagate some distance before the main fracture catches up with it. Smekal (1952) postulated that in the large stress field of the main crack, heterogeneities trigger the initiation of a secondary fracture ahead of the primary crack. The secondary fracture may not be in the same plane as the primary front. When these two fronts intersect in space and time, the ligament separating the two fractures breaks up, leaving a conic marking on the fracture surface. Therefore, the conic marking indicates a level difference boundary, marking the common space time interaction of the two fracture fronts, with the focus of the conic identifying the origin of the secondary fracture front. The existence of the conic markings indicates an increase in the number of voids activated into growing along the fracture path, and an increase in the nucleation distance at which the secondary microcracks

begin to grow (see Sections 7.8.9–7.8.11 for a discussion of fracture branching). Carlsson *et al.* (1972) observed that the number of the markings increases with the fracture velocity, which is consistent with the fact that an increasing number of voids is nucleated ahead of the fracture tip as the stress at the tip increases. We should, however, mention that Ravi-Chandar and Yang (1997) reported that there is no one-to-one correlation between the number or density of the markings and the mean velocity of fracture. This is similar to lack of a one-to-one correspondence between the stress intensity factor and the fracture velocity, which will be discussed below (see Section 7.8.11). Shioya and Ishida (1991) found the depth of the conic markings in PMMA to be approximately 1 $\mu$m. Ravi-Chandar and Yang (1997) carried out a comprehensive study of the development of the conic markings as a function of the velocity of a fracture for four polymeric materials which were PMMA, Homalite-100, Solithane-113 and polycarbonate. Solithane-113 is a polyurethane elastomer which exhibits brittle fracture behavior. Polycarbonate is a non-crosslinked thermoplastic polymer which is capable of inelastic deformation, since the mobility of the carbonate segments of its structure is relatively high. However, at large rates of loading, it does exhibit brittle dynamic fracture. Ravi-Chandar and Yang (1997) found that the markings in all of these materials increase in density with increasing values of the stress intensity factor.

### 7.8.6   Riblike Patterns on Fracture Surface

In the mist and hackle regions of many brittle polymers, such as polystyrene (Hull, 1970), PMMA (Fineberg *et al.*, 1992), and Solithane-113 and polycarbonate (Ravi-Chandar and Yang, 1997), rib-like patterns on the fracture surface are commonly observed. In such materials, the typical distance between the markings is of the order of 1 mm, so that they can easily be seen by naked eye. In PMMA, for example, on which extensive work has been carried out for characterizing such patterns, the rib-like patterns have been found to initiate within the mist regime. The initial width of these patterns is usually much less than the sample's thickness, but it increases with the fracture velocity and eventually, within the hackle zone, extends across the entire thickness of the sample (Sharon and Fineberg, 1996). These patterns, rather than being smooth undulations along the fracture surface, are discrete bands of jagged cliff-like structures. Their height increases with the fracture velocity, and they exist up to the point where a fracture undergoes macroscopic branching. The spacing between the ribs is also strongly related to the molecular weight of the monomers used to form PMMA (Kusy and Turner, 1975), with the typical spacing increasing by over two orders of magnitude as the molecular weight was varied between $10^4$ and $10^6$. Moreover, the fracture energy was found to be a strongly increasing function of the rib spacing.

### 7.8.7   Roughness of Fracture Surface

We already described in Chapter 6 the roughness of fracture surface of materials, and how the associated roughness exponent is measured. In effect, the fracture

surface is a self-affine fractal (see Section 1.3), and studies of aluminum alloys, steel, ceramics and concrete indicated (Bouchaut *et al.*, 1990, 1991; Måløy *et al.*, 1992) that the local width $w$ of the fracture surface scales as

$$w \sim \ell^{\alpha}, \tag{130}$$

where $\ell$ is the scale of observation within the fracture plane, and $\alpha$ (which is usually the same as the Hurst exponent $H$ defined in Chapter 1) is the roughness exponent. Characterization of rough surfaces and measurement of the associated roughness exponent $\alpha$ were discussed in Chapters 1 and 6. As discussed in Section 6.14.1, it appears that for both quasi-static and dynamic fracture a universal roughness exponent, $\alpha \simeq 0.8$, is obtained for $\ell > \xi_c$, where $\xi_c$ is a material-dependent length scale (Daguier *et al.*, 1996, 1997). For $\ell < \xi_c$ a different roughness exponent, $\alpha \simeq 0.5$, has been measured (Milman, 1994). Narayan and Fisher (1992) interpreted $\alpha \simeq 0.5$ as being the result of a crack front pinned by microscopic material inhomogeneities in very slow fracture.

As already explained in Chapter 6, the apparent length-scale dependence of the roughness exponent $\alpha$ may also be explained in another way based on the velocity of fracture propagation. According to Bouchaud and Navéos (1995) (and somewhat similar to the argument of Narayan and Fisher, 1992), one must distinguish between quasi-static (slow) and rapid fracture. In the former case, corresponding to small length scales, one may obtain a roughness exponent close to 0.5, whereas rapid fracture, which corresponds to large length scales, leads to $\alpha \simeq 0.8$. Bouchaud and Navéos (1995) thus argued for the existence of a length scale $\xi_{qs}$, such that for $\ell < \xi_{qs}$ one is in the quasi-static fracture regime and thus a low roughness exponent, while at length scales $\ell \gg \xi_{qs}$ rapid fracture is dominant and therefore one should obtain $\alpha \simeq 0.8$. As shown by Daguier *et al.* (1997), $\xi_{qs}$ depends on the velocity of fracture propagation, and thus should decrease as the velocity increases. If this picture of fracture is correct, then models that are based on minimum energy surfaces are in the quasi-static class. Bouchaud and Navéos (1995) also showed that the data for both cases can be expressed by the following equation

$$\frac{h_{max}}{r^{\alpha_{qs}}} = A_1 + A_2 r^{\alpha - \alpha_{ms}}, \tag{131}$$

where $h_{max}$ is the same as before, $\alpha_{qs}$ is the roughness exponent corresponding to the quasi-static limit, $\alpha$ is the universal roughness exponent corresponding to rapid fracture, $\alpha_{ms}$ is the roughness exponent of minimum energy surfaces, and $A_1$ and $A_2$ are two constants.

We should point out that, despite the considerable effort that had gone into understanding the properties of self-affine fracture surfaces, up until recently, there was little discussion of the fact that, in many of the experiments in which a nontrivial roughness exponent had been measured for a fracture surface, the typical length scales where the scaling behavior had been observed were several orders of magnitude smaller that the typical sample size. For example, the largest length scale observed in measurements performed on soda-lime glass (Daguier *et al.*, 1997) was of the order of 0.1 $\mu$m, which is well within the mirror regime. Thus, in the context

of continuum models of dynamic fracture, the roughness at such length scales does not constitute a departure from straightline fracture propagation, although it is conceivable that the observed scaling structure may affect the value of the fracture energy. Although it is known that the root mean-square surface roughness increases with the velocity of a crack within the mist and hackle regions in PMMA (Fineberg *et al.*, 1991; Boudet *et al.*, 1995), Homalite-100 (Ravi-Chandar and Knauss, 1984a), and crystals that are cleaved at high velocities (Field, 1971; Reidle *et al.*, 1994), we are not aware of any systematic measurements of the dependence of the roughness on the velocity of a crack, at velocities that are of interest to dynamic fracture.

Thus, as pointed out by Fineberg and Marder (1999), the length scales at which the fracture surfaces have been found to be self-affine are, in general, well within the cohesive zone. As a crack accelerates, however, the surface structure within the mist and hackle regimes may, depending on the overall sample size, become larger than the length scales at which the singular contribution to the stress field in the medium is dominant. At this point the structure within the fracture surface may no longer be swallowed up within the cohesive zone, and the description of the dynamics of a crack will be beyond the realm of linear continuum fracture mechanics.

More recent work by López and Schmittbuhl (1998) and Morel *et al.* (1998) has addressed the scale dependence of the roughness of fracture surfaces, and the associated roughness exponent $\alpha$. It has been suggested that the apparently-universal roughness exponent $\alpha \simeq 0.8$ represents a *local exponent* (even though it supposedly corresponds to rapid fracture at larger length scales). Moreover, even if the local roughness exponent, which we now denote it by $\alpha_{\mathrm{loc}}$, is universal, i.e., independent of the material, the range of length scales within which the scaling of the width of the rough surface is observed depends strongly on the material morphology. It has been shown that the scaling laws that govern the crack development in the longitudinal and transverse directions are different and material dependent. Consider, for example, the development of a fracture surface from a flat notch of length $L$ with no roughness. The mean plane of the fracture surface is marked by the coordinates $(x, y)$ where the $x$-axis is perpendicular to the direction of crack propagation, while the $y$-axis is parallel to the crack propagation direction. It has been found that the height fluctuations $\Delta h$ of the fracture surfaces of two heterogeneous brittle materials—granite (López and Schmittbuhl, 1998) and wood (Morel *et al.*, 1998)—estimated over a window of size $\ell$ along the $x$-axis and at a distance $y$ from the initial position exhibits scaling properties that are much more complex than what is predicted by Eq. (130) for the *transverse direction*, and are described by the following *anomalous* scaling properties,

$$\Delta h(\ell, y) \simeq A \begin{cases} \ell^{\alpha_{\mathrm{loc}}}\xi(y)^{\alpha-\alpha_{\mathrm{loc}}}, & \text{if } \ell \ll \xi(y), \\ \xi(y)^{\alpha_{\mathrm{loc}}}, & \text{if } \ell \gg \xi(y), \end{cases} \tag{132}$$

where $\xi(y) = By^{1/z}$ depends on the distance to the initial notch $y$ and corresponds to the crossover length along the $x$-axis below which the fracture surface is self-

affine with a local roughness exponent $\alpha_{loc}$. The quantity $z$ is the dynamic exponent for rough surfaces that was already introduced in Section 1.5.

The scaling laws (132) indicate that along the $y$-axis the roughness develops according to two different regimes: For large length scales $[\ell \gg \xi(y)]$, the roughness grows as $\Delta h \sim y^{\alpha/z}$, where $\alpha$ is called the *global roughness exponent*, whereas for small length scales $[\ell \ll \xi(y)]$ the roughness growth is characterized by the exponent $(\alpha - \alpha_{loc})/z$. Unlike the local roughness exponent, the global exponent $\alpha$, as well as the dynamic exponent $z$ and the prefactors $A$ and $B$ are material dependent, and hence non-universal. Thus, despite exhibiting universality in the transverse direction, roughening in the longitudinal direction is material dependent.

An important consequence of scaling laws (132) is that, when the global saturation occurs, i.e., far from the notch for $y \gg y_{sat}$ [where $y_{sat} = (L/B)^z$], the magnitude of the roughness is not only a function of the window size $\ell$ but also of the system size $L$, since in this case, $\Delta h(\ell, y \gg y_{sat}) \simeq A\ell^{\alpha_{loc}} L^{\alpha - \alpha_{loc}}$. It is for this reason that scaling laws (132) are viewed as *anomalous* because in the conventional scaling of rough surfaces that were described in Section 1.5 one has

$$\Delta h(\ell, y) \simeq A \begin{cases} \ell^{\alpha_{loc}}, & \text{if } \ell \ll \xi(y), \\ \xi(y)^{\alpha_{loc}}, & \text{if } \ell \gg \xi(y). \end{cases} \quad (133)$$

If fact, the scaling laws (132) and (133) become equivalent only if we take the global roughness exponent $\alpha$ to be equal to the local exponent $\alpha_{loc}$.

These anisotropic scaling laws have important implications for the Griffith criterion which will be described shortly.

## 7.8.8  *Modeling Rough Fracture Surfaces*

Although, in addition to their experimental realization, self-affine fracture surfaces have been clearly produced in molecular dynamics simulations of dynamic fracture (see Chapter 9), an important unsolved problem, which is outside the realm of linear continuum fracture mechanics, is a proper model that can generate self-affine fracture surfaces with the roughness exponents that have been measured in many experiments. As usual, this problem has been attacked by many, employing many different ideas. For example, J.P. Bouchaud *et al.* (1993) proposed a model based on directed percolation. In directed percolation (see, for example, Kinzel, 1983; Duarte, 1986, 1990, 1992; Duarte *et al.*, 1992), the bonds of a lattice are directed and diode-like. Transport along such bonds is allowed only in one direction. If the direction of the external potential is reversed, then, there may be no macroscopic transport in the new direction of the external potential. Unlike the regular percolation, there are *two* correlation lengths in directed percolation that characterize the shape of the percolation clusters. One is the longitudinal correlation length (in the direction of the external potential), while the second one is the transverse correlation length, in the direction perpendicular to the direction of the external potential. As a result, one must also have two critical exponents that characterize the scaling of the correlation lengths near the percolation threshold. One is $\nu_L$ which is associated with the longitudinal correlation length, while the

second exponent is $\nu_T$, associated with the transverse correlation length. However, although the directed percolation model does provide a prediction for the roughness exponent, namely, $\alpha = \nu_T/\nu_L$, its numerical value in 2D, $\alpha \simeq 0.63$, or in 3D, $\alpha \simeq 0.57$, is not in good agreement with the data discussed above. J.P. Bouchaud *et al.* (1993) also proposed a set of coupled equations which do have some of the required symmetries and properties appropriate to this phenomenon. Their equations are given by

$$\frac{\partial x}{\partial t} = v + \Sigma\frac{\partial^2 x}{\partial y^2} + \frac{\lambda_{xx}}{2}\left(\frac{\partial x}{\partial y}\right)^2 + \frac{\lambda_{xz}}{2}\left(\frac{\partial z}{\partial y}\right)^2 + \mathcal{N}_x(y,t), \qquad (134)$$

$$\frac{\partial z}{\partial t} = \Sigma\frac{\partial^2 z}{\partial y^2} + \lambda_z\frac{\partial x}{\partial y}\frac{\partial z}{\partial y} + \mathcal{N}_z(y,t). \qquad (135)$$

Here, $x$ is the direction of fracture propagation, $y$ is along the fracture front, $z$ is the tensile axis, $v$ is the nominal fracture velocity, $\Sigma$ is the line tension, $\mathcal{N}$ represents noise or disorder in the material, and the $\lambda$s are coupling constants. The nonlinear terms signify the fact that the local velocity of the fracture depends on its local direction. They are designed to satisfy the required symmetries, namely, $y \rightarrow -y$ and $z \rightarrow -z$. The same type of equations were discussed by Ertas and Kardar (1992, 1993, 1994, 1996) in the context of driven vortex lines in superconductors, and the morphology of polymers in shear flows. In their model the flux lines are pulled away by a constant force. Their equations are nonlinear, with the nonlinearity accounting for the variations of the local propagation speed with the local orientation of the front. Depending on the values of the parameters, many distinct scaling regimes are predicted by these equations. In particular, in a certain limit and for a finite velocity, Ertas and Kardar found that $\alpha \simeq 0.75$ at large length scales, and $\alpha \simeq 0.5$ at short length scales, quite close to the experimental values of $\alpha$ discussed above and in Chapter 6. However, the exact correspondence between the problem discussed by Ertas and Kardar and self-affine fracture surfaces is not clear.

J.P. Bouchaud *et al.* (1993) and E. Bouchaud *et al.* (1993a) also suggested that a fracture surface may be modeled as the trace that is left by the fracture front propagating in a medium with randomly-distributed obstacles. The model proposed by Hansen *et al.* (1991), based on an analogy with directed polymers in random media first proposed by Kardar *et al.* (1986)—the KPZ equation described in Section 1.6—also does not produce the experimentally-measured value of the roughness exponent, since it predicts that $\alpha = 2/3$. Schmittbuhl *et al.* (1995) proposed a perturbative approach to describe the evolution of a fracture between two elastic solids, in which the driving force was the stress intensity factor along the fracture front. The resulting fracture surface was rough and self-affine, but the roughness exponent was only $\alpha \simeq 0.35$, which does not agree with any of the experimental data described above. We will come back to this issue later in this chapter and also in Chapter 8.

A completely different approach was suggested by Räisänen *et al.* (1998). They suggested an analogy between quasi-static fracture surfaces and minimal energy surfaces. Although both types of surfaces are rough, it may seem surprising that

the two can be related, since the minimal energy surfaces, such as those obtained in the random-bond Ising model, seem to have little, if anything, to do with fracture of a material. Nevertheless, Hansen *et al.* (1991) suggested, and Räisänen *et al.* (1998) confirmed by extensive numerical simulations, that the roughness exponent of the two types of surfaces in 2D are the same. In particular, Räisänen *et al.* (1998) used a scalar approximation to model fracture of a brittle material—the random fuse model described in Section 5.2—to provide strong numerical evidence for this equality. However, in 3D the scalar quasi-static fracture model was found to be rougher than the minimal energy surfaces.

### 7.8.9  Fracture Branching at Microscopic Scales

As described in Chapter 6, in an early study of fracture of glass rods, Johnson and Holloway (1968) demonstrated, by progressive etching of the fracture surface in the mist region, the existence of microscopic cracks that branch away from the main fracture. Similar microscopic branched cracks were later observed by Hull (1970) in polystyrene, by Ravi-Chandar and Knauss (1984b) in Homalite-100, and by Anthony *et al.* (1970) during rapid fracture of tool steel. In fact, as we discuss later in this chapter, formation and evolution of micro-branches strongly influence the dynamics of a fracture.

### 7.8.10  Multiple Fractures Due to Formation and Coalescence of Microscopic Voids

Experiments carried out on Homalite-100 by Ravi-Chandar and Knauss (1984a) suggest that, one should not view dynamic fracture as the propagation of a single fracture, but as the coalescence of microscopic voids that are formed ahead of a fracture front. In their experiments fracture was generated via the electromagnetic loading method described in Chapter 6 in which a trapezoidal pressure profile with a 25 $\mu$s rise time and 150 $\mu$s duration was applied to the faces of a seed microcrack. The sample material was large enough that the first reflected waves from its boundaries would not interact with the fracture throughout the experiment. It was observed that within the mist and hackle regions, a front of multiple microscopic parallel cracks, instead of a single fracture, was formed. The cracks in the mirror region tended to propagate within a single plane, whereas in the mist region caustics due to the formation of multiple fractures tips (which were seen in high speed photographs) were observed, the intensity of which increased within the hackle regime as the secondary fractures increased in size. Ravi-Chandar and Knauss (1984a) proposed that formation of the multiple micro-cracks was due to the nucleation of microscopic material flaws or voids, the traces of which were indicated by the conic markings left on the fracture surface. Earlier, Broberg (1979) had in fact proposed that, these voids are nucleated by the large stresses ahead of the fracture front, so that the dynamics of fracture propagation is dictated by the interactions between these growing flaws and the fracture front. We have already discussed this phenomenon, and therefore do not elaborate further.

## 7.8.11   *Microscopic Versus Macroscopic Fracture Branching*

If, relative to the size of a sample material, the crack branches remain small, then they can be considered as part of the cohesive zone. In materials such as Homalite-100, above a certain energy flux, fractures are made of many microscopic cracks propagating in unison. Such microscopic multiple fractures are observed in a variety of materials within the mist and hackle zones. However, in sample materials of any given size, an increase in size of microbranches with the energy release rate $\mathcal{H}$ will eventually make the size of the cohesive zone large enough that the assumptions of continuum fracture mechanics break down. As soon as a crack begins branching, single fracture models are, of course, no longer valid. Therefore, theories that are based on formation of a single fracture can, at best, provide a criterion for when fracture branching may begin. We already discussed a few of such branching criteria, such as that of Yoffe (1951) and those that are based on extremal energy density. However, as discussed above, the same criteria also predict fracture velocities for the onset of branching that are much too large. Other criteria, such as those that postulate a critical value of the stress intensity factor, are not consistent with experiments (Arakawa and Takahashi, 1991; Adda-Bedia and Ben Amar, 1996) since they indicate that there is considerable variation of the stress intensity factor $K_I$ at the point of branching. Another criterion was suggested by Eshelby (1971) according to which a fracture branches when the energy $\Gamma$ that creates a single propagating fracture is large enough to support two single cracks. However, this criterion suffers from the fact that if $\Gamma$ were not a strongly increasing function of $v$, then once branching began, one should observe a large decrease in the velocities of the branches relative to that of the single fracture that preceded the branching event. In glass, however, the post branching velocities either do not decrease at all (Schardin, 1959), or decrease at most by about 10% (Kerkhof, 1973). It should, however, be clear that the Eshelby criterion is a necessary, but not sufficient, condition for fracture branching.

Yoffe's proposal that there exists a universal critical velocity for macroscopic fracture branching is not supported by experimental observations. For example, branching velocities in glass are between $0.18c_R$ and $0.35c_R$ (Schardin, 1959), in PMMA are consistently about $0.78c_R$ (Cotterell, 1965), and in Homalite are between $0.34c_R$ and $0.53c_R$ (Arakawa and Takahashi, 1991). In any experiment on fracture branching, one must ensure that branching occurs at locations that are far from the lateral boundaries so that the system can be considered as effectively infinitely large. Otherwise, experiments have indicated (Ravi-Chandar and Knauss, 1984c) that branching can be created by the arrival of waves generated at the onset of fracture and reflected at the lateral boundaries of the system back into the fracture tip. Despite such difficulties, the consistent values of the measured branching angles in many different materials indicate that there may be a degree of universality in the macroscopic branching process. The branching angles have been typically determined by measurement of the tangent of a branched fracture at distances of the order of a fraction of a millimeter from the fracture tip. They range from 10° in PMMA (Cotterell, 1965) and glass (Johnson and Holloway, 1968) to

14° in Homalite, 15° in polycarbonate (Ramulu and Kobayashi, 1985) and about 18° in steel (Anthony *et al.*, 1970), all of which were measured for materials that were under pure uniaxial tension.

### 7.8.12   Nonuniqueness of the Stress Intensity Factor

Another discrepancy between the theory and experiment was discovered by Ravi-Chandar and Knauss (1984a) in experiments on Homalite-100. They took high speed photographs of the caustic formed at the tip of a fracture initiated by electro-magnetic loading at high loading rates. The velocity of the fracture was estimated from the position of its tip in the photographs, and was compared with the instan-taneous value of the stress intensity factor, which had been estimated from the size of the caustic. In agreement with the theory, at low velocities (below about 300 m/s $= 0.3c_R$) a change in the value of the stress intensity factor resulted in an instantaneous change in the fracture's velocity. However, at higher velocities significant changes in the stress intensity factor produced no measurable change in the fracture's velocity, indicating that the stress intensity factor is *not* a unique function of fracture velocity.

### 7.8.13   Dependence of the Fracture Energy on Crack Velocity

Due to its fundamental importance, the fracture energy Γ—the energy needed for generating a unit fracture surface—and its dependence on the crack velocity have been measured for many different materials, for which the most common technique is the method of caustics described in Section 6.11.2.1. Measured values of the fracture energy Γ in single crystals, which are necessary for initiating crystal cleavage, agree well with the theoretical predictions (see, for example, Lawn, 1993, for a review). In amorphous or polycrystalline materials, however, experiments indicate that $\Gamma(v)$ is a strongly increasing function of a fracture's velocity, the form of which is known only empirically. Most of the fracture energy is dissipated as heat within the cohesive zone, or is radiated from the crack as acoustic energy, or is lost as the emission of photons from excited molecules along the fracture surface—the so-called fracto-emission (Dickinson, 1991).

Figure 7.8 presents some typical measurements of fracture energy Γ versus fracture velocity $v$ for PMMA, Homalite-100 and AISI 4340 steel. Also shown are the dimensionless velocity $v/c_R$ versus $\Delta = K_I/K_{Ic} = \sqrt{\mathcal{H}/\mathcal{H}_c}$, a dimensionless measure of loading which is the ratio of the stress intensity factor $K_I$ and the critical value $K_{Ic}$ of $K_I$ at which fracture first begins. Although these materials are quite different, a common feature among them is the steep rise in Γ as the fracture velocity $v$ increases. For steels, the increase in Γ is due to the fact that the cohesive zone acts as a plastically deforming region (Freund, 1990). However, in the case of PMMA and Homalite-100, which are brittle amorphous materials, there is no reason to expect the classical theory of plasticity to describe deformations near the fracture tip.

FIGURE 7.8. The dependence of the fracture energy $\Gamma$ on the fracture velocity $v$, for (top row, left to right) AISI 4340 steel (Rosakis *et al.*, 1984), and PMMA (Sharon *et al.*, 1996). The bottom row shows the rescaled data, where $\Delta = K_I/K_{Ic}$, and $c_R$ is the Rayleigh wave speed.

Figure 7.8 does in fact reflect the view of Dally (1979) who studied extensively dynamic fracture in amorphous polymers, and in steels. According to him,

(1) the proper way to characterize a dynamic fracture experiment is through presenting the data by two dimensionless numbers which are $v/c$, the ratio of the fracture velocity and a wave speed, and $\Delta = K_I/K_{Ic}$, the ratio of the dynamic stress intensity factor and its critical value at the fracture onset. The relation $v/c = f(\Delta)$ contains most of the information about the dynamics of fracture.
(2) The energy needed for fracture of brittle amorphous materials increases steeply past a critical velocity, where the straight-line fracture becomes unstable to frustrated branching events.

We will come back to these points later in this chapter.

### 7.8.14 Generalized Griffith Criterion for Fractures with Self-Affine Surfaces

If fracture surfaces are self-affine fractals, then one must think about modifying the Griffith criterion in order to accommodate this fact. Such a generalization was first suggested by Mosolov (1993). Bouchaud and Bouchaud (1994), considered the case in which no distinction was made between the growth of the fracture surface in the longitudinal and transverse directions, and the local and global

roughness exponents were assumed to be the same. This case, as described in Section 7.8.7, corresponds to an isotropic fracture surface at small length scales, which we consider first. Thus, consider the case of non-fractal fracture surfaces and derivation of, for example, Eq. (66). We assume quite generally that $K_I \sim r^{-\zeta}$, where $K_I$ is the stress intensity factor [$\zeta = 1/2$ yields Eq. (66)]. If the fracture path is smooth, then the surface energy is simply

$$\mathcal{H} = 2\Sigma w \ell, \tag{136}$$

where $\Sigma$ is the surface tension, $w$ is the width, and $\ell$ is the fracture length increment. The released elastic energy $\Gamma$ is estimated by noting that, since the stress field is relaxed on length scales $r < \ell$ and unperturbed on larger scales, then

$$\Gamma \simeq \frac{K_I^2}{2Y} \int_{r_c}^{\ell} r^{-2\zeta} w r \, dr \simeq \frac{w K_I^2}{4Y(1-\alpha)} \ell^{2-2\zeta}, \tag{137}$$

where $r_c$ is a microscopic cutoff length scale below which the stress saturates, and $Y$ is the Young's modulus. According to the Griffith's criterion, at the onset of fracture one must have $\mathcal{H} = \Gamma$, which results in $\zeta = 1/2$, as expected.

We now suppose that the fracture surface is self-affine at length scales $\xi$ and is represented by a height profile $h(r)$ given by

$$h(r) = \Lambda(r) h_{max} \left( \frac{r}{\xi} \right)^{\alpha}, \quad r \ll \xi \tag{138}$$

where $\Lambda(r)$ is a random variables of order 1. For $r \gg \xi$ we must have $h(r) = \Lambda(r) h_{max}$. Following Griffith's method, one must calculate the surface energy corresponding to opening of the fracture along a distance $\ell \ll \xi$, which is given by

$$\mathcal{H} \simeq 2\Sigma w \int_0^{\ell} \sqrt{1 + \left( \frac{dh}{dr} \right)^2} \, dr. \tag{139}$$

Equation (139) indicates that there is a new length scale $\xi^*$ at $r = \xi^*$ such that one has $dh/dr \simeq 1$; for $r \ll \xi^*$ one has $dh/dr \gg 1$. Bouchaud and Bouchaud (1994) argued that

$$\frac{\xi^*}{\xi} \simeq \left( \frac{h_{max}}{\xi} \right)^{1/(1-\alpha)} \tag{140}$$

One must distinguish between two distinct cases:

(1) If $h_{max} \ll \xi$ or $\xi^* \ll \xi$, which is the regime in which the surface is a self-affine fractal but *shallow*, i.e., it has a mean local angle of the crack profile smaller than 45°, and there is no sample size effect. Then

$$\mathcal{H} \simeq 2\Sigma w \xi \left( \frac{h_{max}}{\xi} \right) \left( \frac{R}{\xi} \right)^{\alpha}, \quad \ell < \xi^*. \tag{141}$$

However, as soon as $\ell > \xi^*$ one has $\mathcal{H} \simeq 2\Sigma w \ell$, even if $R < \xi$, so that the surface energy is similar to that needed to create flat surfaces, even though

the surface is rough, and thus the stress-field singularity is the usual Griffith's singularity, $\ell^{-1/2}$. Equating (137) and (141) leads to $\alpha = 2 - 2\zeta$ (yielding the Griffith's result, $\zeta = 1/2$, when $\alpha = 1$, i.e., when the fracture surface is smooth). Thus, rougher fractures, i.e., those with smaller $\alpha$, lead to a more singular stress field.

(2) In the second regime, $h_{max} \gg \xi$ or $\xi \ll \xi^*$. In this case the slope of the surface over the entire fractal domain is larger than one, resulting in a *spiky* regime, and $h_{max}/\xi$ is a measure of this spikiness. Near the tip of the fracture ($r < \xi$) the stress field is characterized by the exponent $\zeta = \frac{1}{2}(2 - \alpha)$.

However, the above considerations are valid when the anisotropy in the growth of rough fracture surfaces is not taken into account. As described in Section 7.8.7, the height fluctuations in the longitudinal and transverse directions exhibit distinct scaling properties that are characterized by Eqs. (130) and (132). In particular, one has an anomalous, size-dependent scaling in the saturation regime, which must be taken into account if one is to generalize the Griffith criterion for the onset of fracture. Based on these scaling laws, Morel *et al.* (2000) proposed a modified form of the Griffith criterion. To understand their proposal, consider a semi-infinite linear elastic material of thickness $L$ that contains an initial crack at position $\Delta a$ and in Mode I (i.e., under a uniaxial stable and low tension). In the zone where the roughness of the fracture surface grows, i.e., for $\Delta a \ll y_{sat}$ [where $y_{sat} = (L/B)^z$ defined in Section 7.8.7], the critical energy release rate $\mathcal{H}_c$ during fracture propagation (which, in Griffith's approach, is set to be equal to the energy required for generating the corresponding free surfaces at the microscale; see above and Section 6.7) is given by

$$\Gamma_c = 2\Gamma_s \sqrt{1 + \left(\frac{AB^{\alpha-\alpha_{loc}}}{\ell_0^{1-\alpha_{loc}}}\right)^2 \Delta a^{2(\alpha-\alpha_{loc})/z}}, \quad \Delta a \ll y_{sat}, \tag{142}$$

where $\ell_0$ is the lower cut-off for the length scale over which the fracture surface is a self-affine fractal, i.e., $\ell_0$ is the characteristic size of the smaller microstructural element which is relevant for the fracture process, and $\Gamma_s$ is the specific surface energy that characterizes the resistance of the material to fracturing. The quantities $A$ and $B$ and the exponents $\alpha$ and $\alpha_{loc}$ were already introduced and discussed in Section 7.8.7.

On the other hand, when the crack increment is large (i.e., $\Delta a \gg y_{sat}$), which corresponds to the saturation state of the roughness, one has

$$\Gamma_c = 2\Gamma_s \sqrt{1 + \left(\frac{A}{\ell_0^{1-\alpha_{loc}}}\right)^2 L^{2(\alpha-\alpha_{loc})}}, \quad \Delta a \gg y_{sat}, \tag{143}$$

implying that the energy $\Gamma_c$ is independent of $\Delta a$, but depends on the linear size $L$ of the sample, an important characteristic of brittle fracture of heterogeneous materials. Equation (143) indicates that the size effect gives rise to two asymptotic behaviors which are, $\Gamma_c \sim 2\Gamma_s$ and $\Gamma_c \sim L^{\alpha-\alpha_{loc}}$. The crossover between the two

occurs at a length $L_{co} = (\ell_0^{1-\alpha_{\mathrm{loc}}}/A)^{1/(\alpha-\alpha_{\mathrm{loc}})}$. Hence, for $L \ll L_{co}$ the fracture surface is shallow, and there is no size effect, $\Gamma_c \simeq 2\Gamma_s$. In this case, the classical results of linear continuum fracture mechanics are applicable to fracturing of the material. However, for $L \gg L_{co}$ one has a power law

$$\Gamma_c \sim L^{\alpha-\alpha_{\mathrm{loc}}} > 2\Gamma_s. \tag{144}$$

Equation (144) was found to agree with the experimental data for wood (Morel *et al.*, 1998). Note that, if the anomalous scaling is neglected, and the fracture surface is described by scaling laws (133), then

$$\Gamma_c(\Delta a) \simeq 2\Gamma_s \sqrt{1 + \left(\frac{A}{\ell_0^{1-\alpha_{\mathrm{loc}}}}\right)^2}, \tag{145}$$

that is, there is no dependence on the size of the material, which is the case for purely elastic brittle materials.

## 7.8.15  Crack Propagation Faster Than the Rayleigh Wave Speed

Our discussions so far should have made it clear that linear continuum fracture mechanics predicts that a crack cannot propagate with a speed larger than the Rayleigh wave speed $c_R$. Briefly, continuum mechanics predicts that for Mode I tensile loading there is a *forbidden velocity zone* (FVZ) for fracture propagation which is a zone in which the speed of the propagation cannot be larger than $c_R$. For Mode II shear loading, the FVZ exists only for speeds between $c_R$ and shear wave speed $c_t$. Therefore, in Mode I a crack's limiting speed is also $c_R$ because its FVZ between $c_R$ and $c_t$ acts as an impenetrable barrier for the shear cracks to go beyond $c_R$.

However, several experiments have been reported in which the cracks propagated with a speed *larger than $c_R$*. Winkler *et al.* (1970) reported supersonic crack propagation along weak crystallographic planes in anisotropic single crystals of potassium chloride, where the fracture tip was loaded by laser-induced expanding plasma. Supersonic crack tip speeds are those that are larger than the dilatational wave speed $c_l$ which itself is larger than $c_t$. At much larger length scales, indirect observations of intersonic (i.e., one with a speed $v$ between $c_t$ and $c_l$) shear ruptures have been reported for shallow crustal earthquakes (Archuleta, 1982; Olsen *et al.*, 1997). In this case, the fault motion is primarily shear dominated, and the material is not strictly monolithic because preferred weak rupture propagation paths exist in the form of fault lines.

Rosakis *et al.* (1999) carried out interesting laboratory experiments to determine whether in-plane shear intersonic crack growth can be obtained in materials that are under remote shear loading conditions. They utilized two identical plates of Homalite-100 polymer, and introduced a weak plane ahead of the notch tip (used for initiating crack propagation) in the form of a bond between the two identical samples of the materials. The bonding process was done carefully so that the constitutive properties of the bond were close to those of the bulk material. In this way, fracture toughness along the line was lower. Dynamic photoelasticity described

FIGURE 7.9. Supersonic crack propagation velocity in Homalite-100 (after Rosakis *et al.*, 1999).

in Section 6.11.2.3 was used for recording the stress field near the propagating fracture. The sample was subjected to asymmetric impact loading with a projectile at 25 m/s, and sequences of isochromatic fringe patterns were recorded around a shear fracture as it propagated along the interface between the two Homalite halves. Crack tip speeds were measured independently from crack length history. Figure 7.9 shows the speed of the propagating crack versus the crack length. Initially, the crack tip speed is close to the shear wave speed of Homalite-100, beyond which it accelerates and becomes intersonic. Thereafter, it continues to accelerate up to the plane stress dilatational wave speed of the material, then decelerates and approaches a steady-state value of about $\sqrt{2}c_t$. As mentioned above, the speeds between $c_R$ and $c_t$ are in the FVZ.

Observations of fast shear rupture during earthquakes have also provided the impetus for a considerable amount of theoretical work. We already mentioned in Section 7.7.1.1 the theoretical work of Langer (1992) which predicted the possibility of supersonic fracture propagation. Theoretical analysis of Andrews (1976) had already shown that a shear fracture can have a terminal velocity either less than $c_R$ or slightly greater than $\sqrt{2}c_t$, depending on the cohesive strength of the fault plane ahead of the fracture. Burridge *et al.* (1979) concluded from their theoretical analysis that the crack speed regime $c_t < v < \sqrt{2}c_t$ is inherently unstable for dynamic shear crack growth. Broberg (1989) showed that the crack speed regime $c_R < v < c_t$ is forbidden for both opening and shear mode cracks, a result that was mentioned above. He also showed that the regime $c_t < v < c_l$ is forbidden for opening mode cracks only. Finally, Freund (1979) showed that $\sqrt{2}c_t$ is the only speed permissible for a stable intersonic shear crack.

The existence of crack growth with a speed larger than $c_R$ has also been confirmed by large-scale molecular dynamics simulations of dynamic fracture. These simulations will be discussed in Chapter 9.

## 7.9    Shortcomings of Linear Continuum Fracture Mechanics

Our discussion so far has been an attempt for providing an overview of linear continuum fracture mechanics. As discussed above, the general principle is that by balancing the energy flowing into the vicinity of a fracture's tip with what is required for creating new fracture surface one can predict the motion of a straight, smooth fracture. In addition, continuum fracture mechanics can predict both the strength and functional form of the near-field stresses, and its predictions agree well with the experimental data (see below). However, as Fineberg and Marder (1999) pointed out, there still remain several issues that linear continuum fracture mechanics cannot resolve:

(1)  How does the fracture energy in brittle material vary with its velocity?
(2)  What are the main processes happening in the cohesive zone?
(3)  What controls a non-straight path of a rapidly propagating fracture?
(4)  What controls branching of a crack into two macroscopic fractures?
(5)  As discussed above (see also Chapter 6), fractures can develop rough, self-affine surfaces. What is the controlling factor in the transition from a smooth fracture surface to a rough one?

Many of these questions have been answered by the beautiful experimental and theoretical work of the past decade by a few research groups, most notably by Fineberg, Marder, and co-workers, published in a series of papers in the 1990s. Therefore, we first discuss in the next section the essence of these experimental results and the definitive conclusions that one may draw from them. We then describe in the next chapter the recent theoretical and computational work, the predictions of which have turned out to be in excellent agreement with the experimental observations. These developments have helped the emergence of a coherent picture of dynamic fracture in which instabilities caused by fracture branching play a key role. Our discussion of the experimental results follows closely that presented in the review by Fineberg and Marder (1999), while the discussion of the theoretical and computational approaches is patterned after Sahimi (1998) and Fineberg and Marder (1999).

## 7.10    Instability in Dynamic Fracture of Isotropic Amorphous Materials

The experiments of Fineberg, Marder, and co-workers (Fineberg *et al.*, 1991, 1992, 1997; Gross *et al.*, 1993; Sharon *et al.*, 1995; Marder and Gross, 1995; Sharon and Fineberg, 1996, 1998, 1999; Hauch and Marder, 1999) used the conductive

strip method described in Section 6.11.5.2 with high resolutions—up to $\pm 5$m/s for the velocities and 0.2 mm for the spatial resolution. They used PMMA and Homalite-100 and were able to measure the fracture's velocity at $1/20 \, \mu$s intervals for about $10^4$ points throughout the duration of an experiment, which allowed them to follow the long-time dynamics of a fracture in considerable detail. What follows is a discussion of their results as well as those of others. These experiments have helped us understand and resolve a few of the outstanding issues in dynamic fracture.

### 7.10.1   The Onset of Velocity Oscillations

Typical data for fracture propagation in PMMA are shown in Figure 7.10. The fracture was initially at rest. Its tip had ample time to become slightly blunted, hence making it difficult for the fracture to begin propagating. Note that the crack first accelerates abruptly, over a very short a time ($< 1 \mu$s), to a velocity of the order of $v_c = 100 - 200$ m/s, beyond which the dynamics of the fracture is no longer
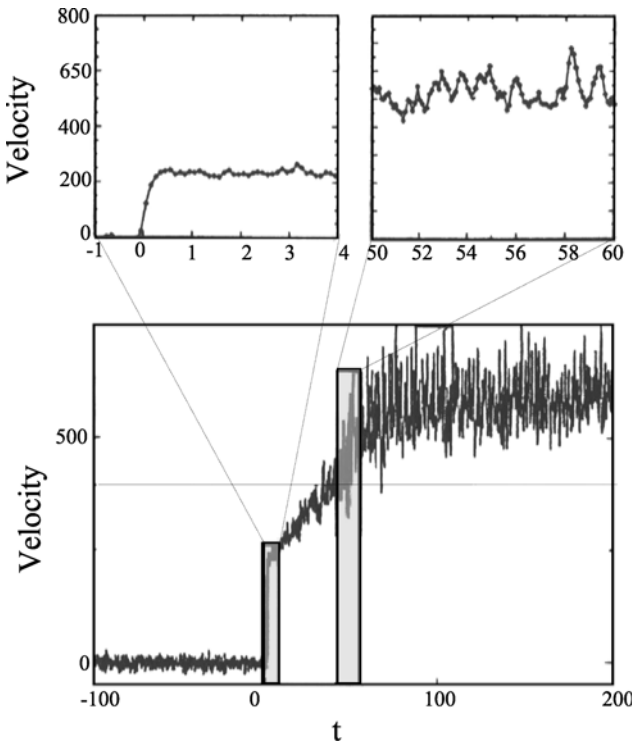


FIGURE 7.10. Typical measurements of velocity (in m/s) of a fracture tip as a function of its length in PMMA. The fracture velocity initially jumps to 150 m/s, and then accelerates smoothly to the critical velocity $v_c$ (dotted line), beyond which strong oscillations set in. The times are in $\mu$sec (after Fineberg and Marder, 1999).

smooth. Instead, one has rapid oscillations in the fracture's velocity which increase in amplitude as $v$ does. On the other hand, Hauch and Marder (1999) carried out experiments in which the energy available per unit length decreased slowly through the length of the sample. In both PMMA and Homalite-100, fractures decelerated gradually to zero velocity, supporting strongly the notion that initial trapping, rather than any intrinsic dynamical effect, is responsible for the velocity jumps, such as those in Figure 7.10, which are always seen when fractures begin to propagate. Indeed, in the case of glass, it is possible to prepare very sharp initial cracks so that their propagation can begin gradually and then continue steadily at velocities that are only a small fraction of the Rayleigh wave speed $c_R$.

The next question is whether the velocity oscillations are random fluctuations or are periodic in time. A careful examination of the oscillations indicate that, although they are not completely periodic, a well-defined time scale does exist with a value that, in the case of PMMA, is typically between 2 and $3\mu$s. Moreover, in experiments in which the fracture accelerates continuously, the location of the peak of the power spectrum of the data in the frequency domain is *constant*, although the velocity varies by as much as 60% of its mean value. As Figure 7.10 also indicates, there is a critical velocity $v_c$ beyond which the fracture velocity begins to oscillate. Many experimental observations indicate that $v_c$ is independent of the sample geometry and thickness, and the applied stress. The value of the critical velocity for PMMA is about $0.36c_R$ which, when surpassed, results in oscillations in the fracture velocity and an increase in the fracture surface area.

## 7.10.2   *Relation Between Surface Structure and Dynamical Instability*

We already described in Sections 7.8.2–7.8.7 the various features that appear in the structure of a fracture surface. How are these features related to the dynamics of fracture propagation? Experiments by Fineberg *et al.* (1997) indicate that the surface structure appears in the close vicinity of $v_c$. The initial surface structure is apparent on only a relatively small amount of the fracture surface. To characterize the amplitude of this structure obtained for PMMA, Fineberg *et al.* (1992) [see also Boudet *et al.* (1995, 1996)] plotted the average height of the points not found in the mirror-like regions within the fracture surface as a function of the mean velocity of the fracture. The results are shown in Figure 7.11. This figure indicates that a well-defined transition occurs where surface structure is created. This happens when the fracture velocity has reached $v = v_c = 0.36c_R$. Moreover, the surface structure is a well-defined and monotonically increasing function of the mean velocity of the fracture. Finally, both the transition point and functional form of the graph are independent of such details as the initial and boundary conditions utilized in the experiment. They are, therefore, intrinsic to the fracture process. Thus, the existence of a well-defined critical velocity $v_c$ for the onset of oscillatory behavior of the fracture and the monotonic dependence of the surface structure created by

FIGURE 7.11. The root mean square values of the surface heights (in $\mu$m) as a function of the mean fracture velocity (in m/s) in PMMA. Different symbols are for various stresses and sample geometries (after Fineberg *et al.*, 1992).

the fracture for $v > v_c$ demonstrate the existence of a dynamical instability in propagation of fracture beyond $v_c$. The dynamical instability is not influenced by either the boundary or initial conditions, and is only a function of the mean velocity of the fracture or, equivalently, the energy release rate, and thus is intrinsic to the system. Moreover, this dynamical instability is a general feature of brittle fracture.

### 7.10.3   Mechanism of the Dynamical Instability

Although there is little, if any, doubt about the existence of an intrinsic dynamical instability during fracture propagation in brittle amorphous materials, the mechanism that gives rise to this instability must be identified, a task that was accomplished by Sharon *et al.* (1995). As already discussed above, experiments indicate that microscopic branches appear within the mist region in a variety of brittle materials, ranging from PMMA to hardened steels. The morphology of these branches was analyzed by Sharon *et al.* as a function of fracture velocity. Their analysis indicated that below the critical velocity $v_c$ no microbranches appear. They begin to emerge at $v_c$, and as the mean velocity of the fracture increases, they become both longer and more numerous. Figure 7.12 presents the mean length of a microbranch as a function of the mean velocity of the fracture, indicating that this quantity is a smooth and well-defined function of the mean velocity. Moreover,

FIGURE 7.12. Mean branch length (in $\mu$m) as a a function of the mean fracture velocity in PMMA. The critical velocity is about 340 m/s (the data are from Sharon and Fineberg, 1996, and Sharon *et al.*, 1996).

similar to Figure 7.11, at $v = v_c$ there is a sharp transition from a state which has no branches to one in which both the main fracture and its daughter cracks are observed. This feature is independent of the initial state of the material. At the same time, a single value of $v_c$ describes both the transition to formation of microbranches and the emergence of the surface structure. Indeed, the surface structure is a result of the crack branching process, and in fact the structure observed on the fracture surface is, essentially, the initial stage of a microbranch which subsequently continues in the material in a direction transverse to the fracture plane.

The microbranching instability is also responsible for the increase in the size of the velocity fluctuations. As a fracture accelerates, the energy released from the potential energy stored in the surrounding material is utilized for generating new fracture surface (i.e., the two new faces created by the fracture). At $v_c$, the energy flowing into the fracture tip is divided between the main fracture and its daughters, resulting in less energy for each crack and a decrease in velocity of the crack ensemble. However, the daughter cracks cannot win their competition with the main fracture, and thus have a finite lifetime. This is presumably because the daughter cracks are screened by the main fracture which, due to its straight-line propagation, outruns them. Thus, after some time the growth of the daughter cracks stops and the energy that was being diverted from the main fracture now returns to it, causing it to accelerate again until the scenario repeats itself.

### 7.10.4   Universality of Microbranch Profiles

For a given mean velocity both the lengths and distances between consecutive microbranches are broadly distributed. Sharon and Fineberg (1996) showed that in PMMA log-normal distributions characterize these quantities with a mean and standard deviation that increase linearly with increasing mean fracture velocity. However, although a given branch may select its length from a broad distribution, such as a log-normal distribution, all microbranches propagate along a highly well-defined trajectory. Indeed, Sharon and Fineberg (1996, 1998) found that these trajectories in *both* PMMA and glass, when considered at the same mean velocity, follow a power law of the form

$$y = 0.2x^{0.7}, \tag{146}$$

where $x$ and $y$ are, respectively, the directions parallel and perpendicular to the direction of propagation of the main crack, with the origin being the point at which the microbranch begins. Much earlier, Hull (1970) had obtained the same result for fracture of polystyrene. These identical trajectories in highly different materials suggest that the microbranch profiles in brittle materials are universal, caused by the universal behavior of the stress field surrounding the fracture. Hull had also proposed that the branch profiles follow the trajectory of maximum tangential stress of the singular field created at the tip of the main fracture (see also the numerical calculations of the stress field of a single static fracture by Parleton, 1979). Moreover, recall (see Section 7.8.11) that value of the branching angle for macroscopic branching in various materials ranges from 11° to 15°, which suggest that a smooth transition between microscopic and macroscopic fracture branches takes place in brittle materials, and that the characteristic features of fracture branches exhibit a high degree of universality. If this is true, then the criterion for the formation of macroscopic fracture branches is identical with the onset of the microbranching instability.

### 7.10.5   Crossover from Three-Dimensional to Two-Dimensional Behavior

The next question to be taken up is the following. What are the circumstances under which a fracture branch survives and continues to propagate away from the main crack? Sharon and Fineberg (1996) proposed that a necessary condition for a microbranch to develop into a full-fledged fracture is the *coherence of the microbranch over the entire thickness of the sample material*. They showed in their experiments on PMMA that, near the onset of the instability, the width of a microbranch is quite small, but as the fracture velocity surpasses the critical velocity $v_c$, both the branch width and length increase. Sharon and Fineberg (1996) used two methods to quantify the increase in the coherence width of the branches. One method was based on a study of the velocity-dependence of the width of the fracture patterns formed by the branches along the fracture surface, which indicated that, beginning with fracture velocities that are close to $v_c$, the width of the pattern

increases sharply with the mean velocity of the crack. When the velocity reaches a value of about $1.7v_c$, the pattern of the growing branches becomes coherent across the entire thickness of the sample. At still higher velocities, macroscopic branching occurs. This phenomenon represents a crossover from a 3D behavior to a 2D one.

The second method for quantifying the coherence of the microbranches, which also helps to further quantify the crossover between the 3D and 2D behavior, is based on measuring the ratio of the total amount of fracture surface produced by the crack and its branches located at the sample faces, and that produced at the center of the sample. Sharon and Fineberg (1996) found that the difference in surface production between the outer and center planes decreases continuously until the fracture velocity is about $1.65v_c$, at which the ratio approaches 1, indicating that microbranch production across the sample is homogeneous. These results are also supported by the experiments of Boudet *et al.* (1996) on PMMA that indicated that both the sound emissions and surface roughness diverge as the mean fracture velocity approached $1.7v_c$, hence suggesting that a *second* transition may occur at $v \sim 1.7v_c$. As the divergence of surface roughness is an indication of macroscopic branching, the crossover from 3D to 2D may be considered to be a *sufficient* condition for macroscopic branching to occur.

### 7.10.6    Energy Dissipation

As we discussed above, the fracture energy $\Gamma$ increases sharply with the fracture velocity. In PMMA, for example (see Figure 7.8), the energy release rate increases by nearly an order of magnitude as the mean fracture velocity exceeds $v_c$. Since for $v > v_c$ the microbranching instability occurs, the total amount of fracture surface created by the fracture front must also increase, thereby leading to an increase in $\Gamma$. Sharon and Fineberg (1996) and Sharon *et al.* (1996) measured, for PMMA and as a function of the mean crack velocity, the relative surface area, defined as the ratio of the total area per unit crack width created by both the main fracture and microbranches, and that which would be formed by a single crack. They also measured the energy release rate. Their data indicate that the amount of surface area formed is a *linear* function of the energy release rate, implying that, both before and after the onset of the instability, the fracture energy is nearly constant. Thus, the fracture energy "increase" shown in Figure 7.11 is entirely a direct result of the microbranching instability. The rise in the fracture energy is due to the formation of more surface by the microbranches. The *cost* of creating a unit fracture surface remains, however, constant with a value which is close to the fracture energy immediately preceding the onset of the instability.

In light of these results, the long-standing question of why in isotropic materials a propagating fracture never seems to approach the Rayleigh wave speed $c_R$ can be answered. A propagating fracture does not *have* to dissipate increasing amounts of energy by accelerating, thereby increasing the amount of kinetic energy. Beyond the critical velocity $v_c$ a fracture has the option of dissipating energy by generating an increased amount of fracture surface at the expense of a reduction in the total kinetic energy. As the amount of energy to its tip increases, a fracture forms a

corresponding amount of surface via *microscopic branching*, the mean length of which also increases with increasing the energy flux to the tip. If this energy increases further, a *second* generation of microbranches may also form (Sharon and Fineberg, 1996) which are the daughters of the daughter cracks. The process of formation of the second, third, $\cdots$, generation of the microbranches may very well be the mechanism for the generation of a fractal structure.

### 7.10.7  *Universality of the Dynamical Instability*

An important question is whether the microbranching instability is a universal feature of dynamic fracture, or is limited to certain types of brittle materials. Much of the experimental data (and also the theoretical work to be discussed later) indicate that the instability is indeed a general feature of brittle fracture. We already mentioned that patterns on the fracture surface have been observed within the mist region in a variety of brittle polymers. In addition to PMMA that was used by Sharon and Fineberg, microscopic branches have also been observed in polycarbonate, polystyrene, hardened steels, glass, as well as in brittle polymers. Additional evidence for this universality is supplied by the fact that, as discussed in Section 7.10.4, microbranches in glass and PMMA develop nearly *identical* trajectories. Moreover, Irwin *et al.* (1979), Ravi-Chandar and Knauss (1984a,b,c) and Hauch and Marder (1999) reported that microbranches are initiated in Homalite beyond $v_c = 0.37c_R$, which is within 2% of the critical velocity observed in PMMA, although the critical velocity for glass seems to be slightly higher (Gross *et al.*, 1993), $v_c \simeq 0.42c_R$, which is still within 20% of the critical velocity for PMMA and Homalite. These results all point to the universal nature of the dynamical microbranching instability in a wide variety of materials, which also makes it possible to describe dynamic fracture of many heterogeneous materials by a unified theory.

## 7.11   Models of the Cohesive Zone

Having described the experimental facts that have helped us understand the nature and characteristics of the microbranching instability in dynamic fracture of amorphous materials, we are now in a position to discuss the theoretical developments, the predictions of many of which agree with the experimental data. An important task is development of a reasonable model of the cohesive zone. Although there has been considerable work devoted to modeling of metals' cohesive zone, our focus in this chapter is on brittle materials. We describe in this chapter the progress that has been made based on the continuum models. Chapter 8 will discuss the lattice models and the insights that they have provided.

As discussed earlier in this chapter, linear continuum fracture mechanics predicts that, as one approaches the tip of a fracture, the stress field diverges as $r^{-1/2}$. However, a divergent stress field is not tenable in a real material. This has motivated the development of many models, both simple and complex, of the cohesive zone

in order to explain how the apparent stress singularity actually joins smoothly a region around the fracture tip where all the fields are finite.

### 7.11.1   The Barenblatt–Dugdale Model

One of the simplest models of the cohesive zone was proposed by Barenblatt (1959a,b) and apparently independently by Dugdale (1960) (see also Langer, 1992). In their model, one assumes that, up to a certain distance $L$ from the tip of the fracture—the length of the cohesive zone—the faces of the fracture are pulled together by a uniform stress $\sigma_c$, which then drops abruptly to zero when the separation between the surfaces reaches a critical separation of $l_c$, as shown in Figure 7.13. The energy absorbed by the cohesive zone can be determined easily, if the fracture propagates in a steady state so that the cohesive zone and all the elastic fields translate in the $x$-direction without changing their form, since in this case translating the fracture by a distance $\Delta x$ increases the length of the material by $\Delta x$ that has passed through the cohesive zone. The energy cost $\Delta \mathcal{H}$ for bringing a length $\Delta x$ of the material through the cohesive zone, per unit length along $z$, is given by $\Delta \mathcal{H} = \Delta x \int_0^{l_0} \sigma_c dy = \Delta x l_0 \sigma_c$. If all the energy that flows into the fracture tip is dissipated by the cohesive forces, then the energy release rate $\mathcal{H}$ equals $l_0 \sigma_c$. The main idea of this model of the cohesive zone is to select $l_0$ and $\sigma_c$ in such a way that the singularities from the linear elastic problem are removed. Therefore, the condition

$$\mathcal{H} = l_0 \sigma_c, \tag{147}$$

must coincide *exactly* with the condition for eliminating the stress singularities. With the aid of Eq. (66), one can then determine the length $L$ of the cohesive
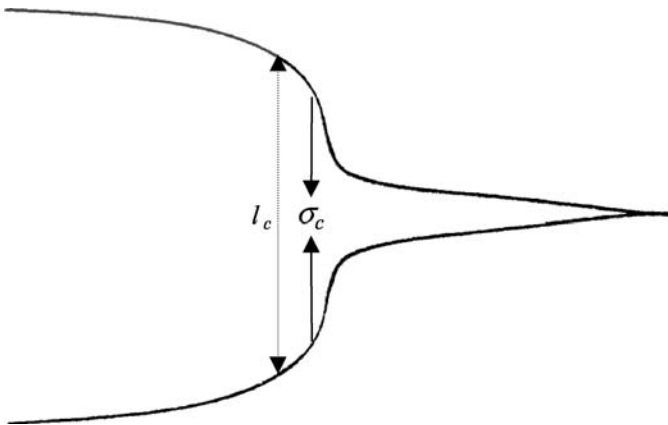


FIGURE 7.13. Schematics of the cohesive zone model of Barenblatt and Dugdale. The faces of the fracture are pulled apart by a cohesive stress $\sigma_c$ until the faces are separated by a critical distance $l_c$. The fracture moves from left to right.

zone, since this zone can be viewed as a superposition of delta-function stresses of the type considered in Section 7.3.2, but with tensile stresses $\sigma_c$, rather than the compressive stresses used there. From Eq. (66), the stress intensity factor is given by

$$K_I = - \int_{-L}^{0} \sigma_c \sqrt{2/\pi l_0} \, dl_0 = -\sigma_c \sqrt{8L/\pi}, \qquad (148)$$

where the negative sign is due to the fact that the cohesive zone is pulling the fracture faces together and cancelling out the positive stress intensity factor which is being generated by other forces outside the fracture. Substituting Eq. (148) into (73) and using Eq. (147) yield

$$\mathcal{H} = \left[ \frac{1 - v_p^2}{Y} A_I(v) \right] \left( \sigma_c^2 \frac{8L}{\pi} \right) = l_c \sigma_c, \qquad (149)$$

from which $L$ is determined. Since as discussed earlier, $A_I(v)$ diverges as the fracture velocity $v$ approaches the Rayleigh wave speed $c_R$, the length $L$ of the cohesive zone must vanish, because the fracture opens more and more steeply as $v$ increases, and therefore it reaches the critical separation $l_c$ sooner and sooner. This type of cohesive zone is frequently observed in fracture of polymers, since behind the fracture tip of such materials, there still are polymers that are arrayed in the craze zone (see Chapter 6) which stretch between the two fracture faces and pull them together. The cohesive zone in metals, on the other hand, is viewed as a simple representation of plastic flow around the fracture tip.

Although this model is simple and has had some success in explaining some aspects of two physics of the fracture energy, as far as explaining the dynamical microbranching instability that we discussed in the last section is concerned, it is not useful at all because, in essence, it replaces one phenomenological parameter, $L$, by the two phenomenological parameters, $\sigma_c$ and $l_0$, and hence provides no new information or even a clear sense of how the dissipated energy varies with fracture velocity.

More realistic models of the cohesive zone have played an important role in providing a better understanding of the dynamical instabilities in brittle fracture of materials. Some of these models are first formulated based on the continuum mechanics, but are then discretized using the finite-element method and simulation, or utilize analytical or semi-analytical analysis. In what follows we discuss these models and the implications of their predictions.

## 7.11.2   Two-Field Continuum Models

An important advance in the continuum formulation of dynamic fracture has been the development of the so-called two-field models that couple the equation for elastic deformation of materials to one for the *order-parameter* of the system. The concept of an order parameter is borrowed from theories of thermodynamic phase transitions in which this parameter represents, for example, the difference between

the densities in the gas and liquid phases that are in equilibrium with each other. This concept is also well-defined for geometrical models, such as the percolation model for which the order parameter represents the fraction of the uncut bonds or sites, where the cut bonds represent the "defects." Hence, at the percolation threshold $p_c$, where the geometrical connectivity of the system is lost due to the presence of too many defects, the order parameter is zero, slightly above $p_c$ is very small, while far from $p_c$ the order parameter is nearly unity, since in this region the defects are too few. In a similar spirit, the order parameter for dynamic fracture should be related to the concentration of point defects in the material, hence characterizing *local order*. In this formulation, the order parameter is (similar to the percolation model) unity outside of the propagating fracture, but zero inside the crack where all the atomic bonds have been broken. On the crack surface, the order parameter varies continuously between 0 and 1, on length scales that are much larger than the interatomic distances. This would then justify use of a continuum formulation of dynamic fracture propagation, in which case one would need an equation for the order parameter that couples it to the equation for the elastic deformation, hence the name two-field models. The advantage of formulating the problem in terms of an order parameter and coupling it to the displacement field is that, by allowing the order parameter to vary in the cohesive zone, the stress singularity at the tip of the fracture is avoided, hence removing one main deficiency of continuum fracture mechanics.

One such two-field model was developed by Aranson *et al.* (2000). They focused on 2D materials in Mode I fracture, and represented the elastic deformation of an amorphous material by the usual wave equation, coupled to a term that represents viscous damping:

$$\rho_0 \frac{\partial^2 \mathbf{u}}{\partial t^2} = \eta \nabla^2 \left( \frac{\partial \mathbf{u}}{\partial t} \right) + \mathbf{\nabla} \cdot \boldsymbol{\sigma}, \tag{150}$$

where the first term on the right-hand side accounts for viscous damping with $\eta$ being the viscosity, and $\rho_0$ is the material's density which is taken to be unity. The stress tensor $\boldsymbol{\sigma}$ is related as usual to the strain tensor $\boldsymbol{\epsilon}$, except that their relation now contains a term involving the order parameter $\mathcal{P}$. This relation, in component form, is given by

$$\sigma_{ij} = \frac{Y}{1 + v_p} \left( \epsilon_{ij} + \frac{v_p}{1 - v_p} I_\epsilon \delta_{ij} \right) + a_1 \frac{\partial \mathcal{P}}{\partial t} \delta_{ij}, \tag{151}$$

where $a_1$ is a constant, $I_\epsilon$ is the trace of the strain tensor, and the rest of the notations are as before. One must take into account the effect of the material's weakening by fracture which reduces the Young's modulus $Y$. Therefore, Aranson *et al.* (2000) assumed that, $Y \sim Y_0 \mathcal{P}$, where $Y_0$ is the initial Young's modulus. In Eq. (151) the term that couples the stress and strain tensors to the order parameter accounts for the hydrostatic pressure that one must apply to the material in order to generate new defects. Although one might be tempted to interpret this term as being due to the material's thermal expansion, such identification would be erroneous since, as discussed in Chapter 6 and earlier in this chapter, during fracture propagation

the temperature at the tip of the crack will be high and therefore it is unlikely that the tip will be in thermal equilibrium. Note that for $\mathcal{P} = 1$—the area outside the fracture surface—one has the usual equations of elasticity, while for $\mathcal{P} = 0$ the dynamics is trivial as nothing is happening inside the fracture.

The next step is to develop an expression for the order parameter. Aranson *et al.* (2000) assumed that $\mathcal{P}$ is governed by purely dissipative dynamics. As such, the order parameter may be derived from an free-energy functional $\mathcal{H}$,

$$\frac{\partial \mathcal{P}}{\partial t} = -\frac{\delta \mathcal{H}}{\delta \mathcal{P}}, \tag{152}$$

which is the standard practice in thermodynamics. In the theory of phase transitions one has (see, for example, Landau and Lifshitz, 1980)

$$\mathcal{H} = \int \left[ a_2 |\nabla \mathcal{P}|^2 + \mathcal{H}_p(\mathcal{P}) \right] dx dy, \tag{153}$$

where $\mathcal{H}_p$ is a local potential energy that has minima at $\mathcal{P} = 0$ and 1. If we choose $\mathcal{H}_p$ to be a polynomial in $\mathcal{P}$, we arrive at (Aranson *et al.*, 2000)

$$\frac{\partial \mathcal{P}}{\partial t} = a_2 \nabla^2 \mathcal{P} - a_3 \mathcal{P}(1 - \mathcal{P}) F(\mathcal{P}, I_\epsilon) + f(\mathcal{P}) \frac{\partial \mathcal{P}}{\partial x_l} \frac{\partial u_l}{\partial t}. \tag{154}$$

Therefore, the order parameter is coupled to the displacement field through Eqs. (151) and (154) and the function $F(\mathcal{P}, I_\epsilon)$, which is subjected to the constraint that it must have one zero in the interval $0 < \mathcal{P} < 1$, so that $F(\mathcal{P}_c, I_\epsilon) = 0$ for $0 < \mathcal{P}_c < 1$, and $\partial F(\mathcal{P}_c, I_\epsilon)/\partial \mathcal{P} = 0$. The simplest functional form for $F$ that satisfies these constraints is given by $F = 1 - (a_4 - a_5 I_\epsilon)\mathcal{P}$, where $a_4$ and $a_5$ are material-dependent constants that can be set to 1 by rescaling of the time, $t \rightarrow a_3 t$ and the spatial coordinates $x_i \rightarrow a_1 x_i$ with $a_1^2 = a_2/a_3$.

The last term on the right-hand side of Eq. (154) couples the order parameter to the speed $d\mathbf{u}/dt$, and represents the localized shrinkage of the fracture caused by the motion of the material. Aranson *et al.* stated that the precise form of the function $f(\mathcal{P})$ is immaterial, and therefore they used a simple form, $f = a_6 \mathcal{P}(1 - \mathcal{P})$, where $a_6$ is a dimensionless material constant (taken to be 1). This completes the formulation of the problem.

However, we must point out that although these functional forms for $F$ and $f$ facilitate the solution of the problem, they also lead to certain anomalies. For example, the model predicts that the crack opening depends logarithmically on the sample size, as opposed to the correct linear dependence. The root of this anomalous dependence is in the fact that in this model the strain in the bulk of the material is *not* fully relieved after passage of the fracture, and hence a more sophisticated formulation of these functions is necessary. Despite this deficiency, the model does predict dynamical instability of the type described above which we now describe.

This model predicts crack branching, with the size of the branches being dependent on the parameters of the materials. The angle of the branches with the main propagating crack is around 30°, but increases with the crack speed. Figure 7.14

FIGURE 7.14. Fracture velocity $v$, normalized by the Rayleigh wave speed $c_R = 926$ m/s in PMMA, versus dimensionless energy $\mathcal{H}/\mathcal{H}_c$. Open circles correspond to stable propagation, crosses to unstable propagation, while diamonds are experimental data of Sharon *et al.* (1996). The inset shows the curvature for unstable propagation at $\rho = 0.5$, with the arrows indicating the progression of time (after Aranson *et al.*, 2000).

compares the predictions of the model for the crack velocity $v$ in PMMA with the experimental data of Sharon *et al.* (1995, 1996). The crack velocity has been normalized by the Rayleigh speed $c_R$, and is plotted versus the fracture energy $\mathcal{H}$ normalized by its value at $v = 0.2c_R$. The parameters used in the simulations were $Y_0 = 10$, $v_p = 0.36$, and $\eta = 13/\sqrt{Y_0}$. For PMMA, the Rayleigh wave speed is $c_R = 926$ m/sec. The model predicts that, depending on the material's parameter, a crack instability develops when its speed varies anywhere from $0.32c_R$ to $0.55c_R$, with the instability manifesting itself as pronounced velocity oscillations, sound emission from the crack tip and, of course, crack branching, as mentioned above. The agreement between the predictions and the experimental data shown in Figure 7.14 is quite good, indicating the correctness of the model in having most of the essential features of dynamic fracture. For a somewhat related model see Karma *et al.* (2001).

## 7.11.3   Finite-Element Simulation

Johnson (1992,1993) and Xu and Needleman (1994) carried out extensive numerical simulations of dynamic fracture in model isotropic elastic materials. Their simulations, which were based on discretization of the governing equations with

the finite-element (FE) method, have the closest correspondence with experiments in brittle amorphous materials. In particular, similar to the experiments discussed above, these FE simulations produced frustrated crack branching, oscillations in fracture velocities, and limiting crack velocities below the Rayleigh wave speed $c_R$. Let us describe and discuss these successful efforts.

The basis of Johnson's work was the physical fact that the size of the cohesive zone is *not* predetermined, but is adaptive and changes in accordance with the fracture's behavior. Since the main purpose of the work was to investigate material weakening and the role of the cohesive zone, accurate modeling of the continuum region outside the zone was not essential, and therefore Johnson assumed the continuum to be linearly elastic. In addition, the material modeled was highly idealized in the sense that, no viscoplastic flow or other rate-dependent properties were incorporated in the cohesive zone. A planar stress model was used, and an initial crack of length $a_0 = 0.6h$ was inserted in the system, where $h$ is the length of the plane. The material in the vicinity of the crack tip was assumed to have a large number of sites where nucleation of defects, all being of the same type and having the same size, occurs. The fractures were driven by loading their faces with a number of different loads. Depending on the applied load, the FE simulations produced maximum fracture velocities of $0.29c_R$, $0.44c_R$ and $0.55c_R$. Moreover, the simulations predicted that, at the lowest velocities, a fracture would accelerate smoothly. As the external loading was increased, multiple attempts at microbranching were observed and, similar to the experiments discussed above, the length of the attempted branches increased with the loading. Moreover, the experimental observations of Ravi-Chandar and Knauss (1984a) (see Section 7.8.12) that the stress intensity factor is not a unique function of the fracture velocity $v$, once $v$ exceeds a certain limit, were also reproduced by these FE simulations. None of these results was dependent on the various parameters of the simulations.

More extensive FE simulations were carried out by Xu and Needleman (1994), although their model of the cohesive zone was different from Johnson's, and was also much more elaborate. The continuum was characterized by two constitutive relations; (1) a volumetric constitutive law that related stress and strain, and (2) a cohesive surface constitutive relation between the tractions and displacement jumps across a specified set of cohesive surfaces that were interspersed throughout the continuum. The first constitutive law was that for an isotropic hyperelastic solid:

$$\boldsymbol{\sigma}_{\text{PK}} = \frac{\partial \mathcal{H}_s}{\partial \boldsymbol{\epsilon}}, \tag{155}$$

where $\mathcal{H}_s$ is the strain energy density which is given by

$$\mathcal{H}_s = \frac{1}{2}\boldsymbol{\epsilon} : \mathbf{C} : \boldsymbol{\epsilon}, \tag{156}$$

where $\mathbf{C}$ is the tensor of the elastic moduli. Here, $\boldsymbol{\epsilon}$ is the Lagrangian strain, and $\boldsymbol{\sigma}_{\text{PK}}$ is the so-called second Piola–Kirchhoff stress, given by

$$\boldsymbol{\sigma}_{\text{PK}} = \boldsymbol{\sigma} \cdot (\mathbf{F}^{-1})^{\text{T}}, \tag{157}$$

where $\boldsymbol{\sigma}$ is the non-symmetric nominal stress tensor, $\mathbf{F}$ is the deformation gradient, and T denotes the transpose operation. If, relative to a fixed Cartesian coordinate system, a material point was initially at $\mathbf{x}_0$ and in the current position is at $\mathbf{x}$, then, $\mathbf{F} = \partial\mathbf{x}/\partial\mathbf{x}_0$. In addition,

$$\boldsymbol{\epsilon} = \frac{1}{2}(\mathbf{F}^{\mathrm{T}} \cdot \mathbf{F} - \mathbf{U}), \tag{158}$$

where $\mathbf{U}$ is the identity tensor.

The constitutive law for the cohesive surface was taken to be a phenomenological mechanical relation between the traction $\mathbf{T}$ and displacement jump $\boldsymbol{\Delta}$ across the surface. This constitutive law must be such that, as the cohesive surface separates, the magnitude of $\mathbf{T}$ first increases, reaching a maximum, and then approaches zero with increasing separation. Xu and Needleman (1994) assumed the constitutive relation for each cohesive surface to be elastic, so that any dissipation associated with the separation is neglected, in which case one has

$$\mathbf{T} = \frac{\partial\phi}{\partial\boldsymbol{\Delta}}, \tag{159}$$

where $\phi$ is a potential which in 2D is given by

$$\phi(\boldsymbol{\Delta}) = \phi_n + \phi_n \exp(-\Delta_n/\delta_n)\left\{\left(1 - r + \frac{\Delta_n}{\delta_n}\right)\frac{1-q}{r-1} - \left[q + \left(\frac{r-q}{r-1}\right)\frac{\Delta_n}{\delta_n}\right]\exp(-\Delta_t^2/\delta_t^2)\right\}. \tag{160}$$

Here, $\mathbf{n}$ and $\mathbf{t}$ are unit vectors that are normal and tangent, respectively, to the surface at a given point in the reference configuration, $\Delta_n = \mathbf{n} \cdot \boldsymbol{\Delta}$, $\Delta_t = \mathbf{t} \cdot \boldsymbol{\Delta}$, $q = \phi_t/\phi_n$, and $r = \Delta_n^*/\delta_n$, where $\phi_n$ and $\phi_t$ are the work of normal and tangential separation, respectively, $\Delta_n^*$ is the value of $\Delta_n$ after complete shear separation with $T_n = 0$, and $\delta_t$ and $\delta_n$ are characteristic lengths. The two separation works are given by, $\phi_n = e\sigma_n\delta_n$, and $\phi_t = \sqrt{e/2}\,\sigma_t\delta_t$, where $\sigma_n$ and $\sigma_t$ are the cohesive surface normal strength and tangential strength, respectively, and e= $\exp(1)$. All the physical parameters used in the simulations were made to correspond to an isotropic elastic material with the properties of PMMA.

To model the fracture tip, Xu and Needleman (1994) used a model of the cohesive zone similar to what we described in Section 7.11.1 that takes into account both tensile and shear stresses, and also allows for the creation of new fracture surface with no additional dissipation added to the system. In order to allow fractures to branch off the main fracture line, an underlying grid of lines was used on which material separation was allowed if a critical condition was reached. Therefore, this type of simulation combines features of FE models with lattice models (see Chapter 8), but is in some respect more realistic than the lattice models. The computations were carried out for a center-fractured rectangular block, and plane strain conditions were assumed to prevail. Since both the volumetric and surface constitutive relations are elastic, no dissipation mechanism was incorporated into the model. As a result, the work done by the imposed loading was partitioned into kinetic energy, strain energy stored in the material volume, and elastic energy stored

in the cohesive surfaces. The FE discretization was based on linear displacement triangular elements that were arranged in a cross-triangle quadrilateral pattern.

The results of these simulations were very much similar to the experiments in PMMA. Beyond a critical velocity of $0.45c_R$, fracture velocity oscillations together with attempted fracture branching were produced. The branching angle was $29°$, which is close to the maximum branching angle of $32°$ that has been obtained in the experiments. Moreover, when the fracture was constrained to move along a straight line, it accelerated to velocities close to $c_R$, in agreement with the experiments of Washabaugh and Knauss (1994); see Section 7.7.1. Hence, these FE computations produce results that can describe many, but not all, of the instabilities in the fracture of PMMA observed in experiments and described above, and in this regard are more successful than most approaches to dynamic fracture.

## 7.11.4  Fracture Propagation in Three Dimensions

Several investigations have explored the possibility that the instability of fracture tip arises naturally from a wiggly fracture front that propagates through a heterogeneous material. Notable among these investigations are those of Rice and co-workers (Rice *et al.*, 1994; Perrin and Rice, 1994; Morrissey and Rice, 1998), Willis and Movchan (1995, 1997), Movchan and Willis (1995) and Ramanathan and Fisher (1997,1998), which we now discuss.

Rice *et al.* (1994) studied the stability of a straight-line, half-plane fracture front propagating dynamically through an unbounded heterogeneous solid. We provide here some details of their method for studying this problem, as a good example of the type of effort that such problems require. They considered the scalar approximation,

$$\frac{\partial^2 u}{\partial t^2} = c^2 \nabla^2 u, \tag{161}$$

where $u$ is a displacement field representing tensile opening or shear slippage, and $c^2 = E_m/\rho$, with $E_m$ and $\rho$ being an elastic modulus and density, respectively. Equation (161) is easily derived by assuming that the material occupies a volume $\Omega$ with an external surface $S$ on which a load $q$ is applied. If we then form the Lagrangian $\mathcal{L}$, i.e., the difference between the kinetic and potential energies of the system,

$$\mathcal{L} = \int_\Omega \frac{1}{2}\left(\frac{\partial u}{\partial t}\right)^2 d\Omega - \left(\int_\Omega \frac{1}{2}E_m|\nabla u|^2 d\Omega - \int_S qu\, dS\right), \tag{162}$$

and use variational principles, Eq. (161) is obtained. Suppose now that $x = \ell(t)$ is the growth history of the fracture in the 2D version of the problem in which a straight line front propagates in the $x$-direction, that the loadings are such that the static solution of the problem has a stress intensity factor $K_0$ for any position of the fracture front, and that, compared to length scales of interest, all loadings are applied far from the fracture tip. Then, the 2D version of the model equations become identical to those that govern anti-plane strain in actual elastodynamics.

Eshelby (1969) derived the following equation for anti-plane solution for arbitrary fracture propagation:

$$u(x, y, t) = \sqrt{\frac{2}{\pi}} \frac{K_0}{E_m} \, \mathrm{Im} \left[ \sqrt{x - \ell(t_r) + iy} \right], \tag{163}$$

where $t_r = t_r(x, y, t)$ is a retarded time at which a signal arriving at position $(x, y)$ at time $t$ was launched at the fracture tip, and satisfies the equation, $c^2(t - t_r)^2 = [x - \ell(t_r)]^2 + y^2$. The actual analysis is for a finite body prior to the arrival back at the fracture tip of waves that are reflected from boundaries or from another fracture tip, in which case the 2D solution very near the tip is given by

$$u(x, y, t) = \sqrt{\frac{2}{\pi}} \frac{K}{\alpha E_m} \, \mathrm{Im} \left[ \sqrt{x - \ell(t) + i\alpha y} \right] + \text{higher order terms}, \tag{164}$$

where $\alpha = \sqrt{1 - v^2(t)/c^2} = \sqrt{1 - (d\ell/dt)^2/c^2}$, and $K$ is the instantaneous stress intensity factor given by,

$$K = K_0 \sqrt{1 - v(t)/c}. \tag{165}$$

The corresponding energy release rate $E$ is then given by

$$\mathcal{H} = \mathcal{H}_0 \sqrt{[1 - v(t)/c]/[1 + v(t)/c]}, \tag{166}$$

where $\mathcal{H}_0 = K_0^2/(2M)$.

Rice *et al.* (1994) derived the 3D solution as a linearized perturbation about the 2D solutions for a fracture propagating at a steady speed $v_0$ [hence, $\ell(t) = v_0 t$]. Thus, if we use polar coordinates such that, $r \exp(i\theta) = x - v_0 t + i\alpha_0 y$, the 2D solution becomes

$$u_0(x, u, t) = \sqrt{\frac{2}{\pi}} \frac{K_0}{\alpha_0 E_m} \, \mathrm{Im} \left[ \sqrt{x - v_0 t + i\alpha_0 y} \right] = \sqrt{\frac{2}{\pi}} \frac{K_0}{\alpha_0 E_m} \sqrt{r} \sin \left( \frac{1}{2}\theta \right), \tag{167}$$

which is consistent with that of actual elastodynamics for anti-plane strain, if we identify $E_m$ and $c$ with the shear modulus and shear wave speed. To develop the 3D solution, one sets $x = \ell(z, t) = v_0 t + \epsilon f(z, t)$, a first-order expansion in $\epsilon$ about the 2D results corresponding to a straight fracture ($\epsilon = 0$) propagating along the $x$ axis with a constant velocity $v_0$. Thus, the shape of the fracture front can deviate from being straight. The 3D solution is then of the form,

$$u(x, y, z, t; \epsilon) = u_0(x, y, t) + \epsilon\phi(x, y, z, t) + O(\epsilon^2), \tag{168}$$

where $\phi(z, y, z, t) = (\partial u/\partial \epsilon)_{\epsilon=0}$. The singular part of the 3D solution must be of the 2D character, but now relative to the *local* direction of fracture propagation, so that for any $\epsilon$ we must have

$$u(x, y, z, t; \epsilon) = \sqrt{\frac{2}{\pi}} \frac{K(z, t; \epsilon)}{E_m \alpha(z, t; \epsilon)} \mathrm{Im} \left[ \sqrt{x - (v_0 t + \epsilon f) \cos \gamma + i\alpha(z, t; \epsilon)y} \right] + \cdots \tag{169}$$

where $\alpha(z, t; \epsilon) = \sqrt{1 - v^2(z, t; \epsilon)/c^2}$, $v(z, t; \epsilon) = (v_0 + \epsilon \partial f / \partial t) \cos \gamma(z, t; \epsilon)$, and $\cos \gamma(z, t; \epsilon) = [1 + (\epsilon \partial f / \partial z)^2]^{-1/2}$, with $\gamma$ being the angle between the local normal to the fracture front and the $x$- axis. It is then easy to show that as $r \to 0$ [i.e., as $x \to \ell(z, t)$ and $y \to 0$] one has

$$\lim_{r \to 0} [\phi(x, y, z, t)\sqrt{r}] = \sqrt{\frac{1}{2\pi} \frac{K_0}{\alpha_0 E_m}} f(z, t) \sin\left(\frac{1}{2}\theta\right), \qquad (170)$$

so that $\phi(x, y, z, t)$ satisfies the same equation as (161), subject to the stress-free boundary condition, $\partial \phi / \partial y = 0$ at $y = 0$ if $x < v_0 t$. We also have, by symmetry, $\phi = 0$ at $y = 0$ when $x > v_0 t$. In the harmonic case, $f(z, t) = F(k, \omega) \exp(-ikz + i\omega t)$, the solution for $\phi$ is written as

$$\phi(x, y, z, t; k, \omega) =$$

$$\sqrt{\frac{1}{2\pi} \frac{K_0}{\alpha_0 E_m}} F(k, \omega) \exp[i(\omega t - kz)] \exp[-i\omega v_0 (x - v_0 t)/\alpha_0^2 c^2] \psi(x - v_0 t, y; k, \omega). \qquad (171)$$

Since $\phi$ must satisfy Eq. (161), we find that $\psi$ must satisfy the following equation

$$\left(\frac{\partial^2}{\partial x^2} + \frac{1}{\alpha_0^2} \frac{\partial^2}{\partial y^2}\right) \psi = \left(\frac{\partial^2}{\partial r^2} + \frac{1}{r} \frac{\partial}{\partial r} + \frac{1}{r^2} \frac{\partial^2}{\partial \theta^2}\right) \psi = Q^2 \psi, \qquad (172)$$

where

$$Q(k, \omega) = \frac{|k|}{\alpha_0} \left(1 - \frac{\omega^2}{\alpha_0^2 k^2 c^2}\right)^{1/2}, \quad \omega^2 < \alpha_0^2 k^2 c^2, \qquad (173)$$

$$Q(k, \omega) = \frac{i\omega}{\alpha_0^2 c} \left(1 - \frac{\alpha_0^2 k^2 c^2}{\omega^2}\right)^{1/2}, \quad \omega^2 > \alpha_0^2 k^2 c^2. \qquad (174)$$

Equation (173) corresponds to letting $k$ approach the positive real axis through $\mathrm{Im}(k) > 0$ and the negative real axis through $\mathrm{Im}(k) < 0$; these approaches are then taken to be branch-cut portions of the $\mathrm{Re}(k)$-axis where $|k| > |\omega|/(\alpha_0 c)$. Equation (174) holds for any direction of approach. Note that the combination $\alpha_0 c$, which often appears in solutions of fracture propagation problems, has a clear physical interpretation: It is the speed at which information is transmitted transversely along the propagating fracture front. That is, two points of the fracture front a distance $\Delta z$ apart do not influence each other before the time delay $\Delta z/(\alpha_0 c)$.

The solution $\psi$ must satisfy the asymptotic requirement (170) as $r \to 0$. Any more general fracture perturbation $f(z, t)$ can be represented as a Fourier superposition, so that

$$F(k, \omega) = \int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} f(z, t) \exp[-i(\omega t - kz)] \, dz dt. \qquad (175)$$

The general solution for $\phi(x, y, z, t)$ for any $f(z, t)$ is then given by

$$\phi(x, y, z, t) = \sqrt{\frac{r}{2\pi^3}} \frac{K_0 \sin(\frac{1}{2}\theta)}{\alpha_0 E_m} \int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} \frac{\partial f(z', t')}{\partial t'} \frac{c(t - t') - v_0(x - v_0 t)/\alpha_0^2 c}{(x - v_0 t)^2/(\alpha_0^2) + y^2 + (z - z')^2}$$

$$\times \frac{\Theta[c(t - t') - \sqrt{(x - v_0 t')^2 + y^2 + (z - z')^2}]}{\sqrt{c^2(t - t')^2 - (z - z')^2 - y^2 - (x - v_0 t')^2}} \, dt' dz',$$

(176)

where $\Theta[\ ]$ is the Heaviside unit-step function. Once $\psi$ is obtained, $\phi$, and hence the displacement field $u$, are also obtained.

One can now derive an expression for the stress intensity factor (and hence the energy release rate). To do this, it is convenient to replace $\epsilon f(z, t)$ by $\ell(z, t) - v_0 t$ and $\epsilon \partial f(z, t)/\partial t$ by $v(z, t) - v_0$ in all the expressions. To obtain the first order perturbation to the stress intensity factor at some location $\zeta$ along the $z$-axis, the crack front $\ell(z, t)$ is written as

$$\ell(z, t) = v_0 t + [\ell(\zeta, t) - v_0 t] + \{\ell(z, t) - \ell(\zeta, t)\},$$

(177)

where the $[\ ]$ term describes a 2D perturbation, which is solvable exactly to all orders by Eqs. (163), (165) and (166), while the $\{\ \}$ term corresponds to a 3D perturbation that vanishes at $z = \zeta$ for all $t$. The stress intensity factor at $z = \zeta$, due to small deviations from straightness in other fractures, is determined by applying to solution (176) for $\phi$ the operator $\lim_{r \to 0} E_m \sqrt{2\pi r} \partial/\partial y$. The result is given by

$$K(z, t) = K_0 \sqrt{1 - v_0/c} + \left[ K_0 \sqrt{1 - v(z, t)/c} - K_0 \sqrt{1 - v_0/c} \right] + \left\{ K_0 \sqrt{1 - v_0/c} \, I(z, t) \right\},$$

(178)

with

$$I(z, t) = \frac{1}{2\pi} \text{PV} \int_{-\infty}^{+\infty} \int_{-\infty}^{t - |z - z'|/(\alpha_0 c)} \frac{c(t - t')[v(z', t') - v(z, t')]}{(z - z')^2 \sqrt{[\alpha_0 c(t - t')]^2 - (z - z')^2}} \, dt' dz',$$

(179)

with PV denoting the principal value integral, and $v(z, t) = \partial \ell(z, t)/\partial t$ being the local velocity of the propagating fracture. Therefore, the dependence of the stress intensity factor on the shape of the fracture front and its deviations from being straight are expressed in terms of $I(z, t)$. The $[\ ]$ term of Eq. (178) is actually exact for arbitrarily large perturbations of $v(z, t)$, but the $\{\ \}$ term is exact only to first order in the deviation $v(z, t) - v_0$. The choice of $v_0$ is arbitrary so long as it is in the range of "first-order difference" from $v(z, t)$.

If we examine the expression for $I(z, t)$, we see that when a segment of the fracture front suddenly slows down relative to neighboring locations along the front, a reduction in $K$ radiates outward from that segment at speed $\alpha_0 c$. Similarly, when a segment speeds up, an increase of $K$ is radiated. Such elementary slow-down and speed-up are due to the encounters of the fracture front with regions of higher or lower resistance to fracture. Rice *et al.* also (1994) found that when a straight fracture front approaches a slightly heterogeneous strip which lies parallel to the fracture tip along an otherwise homogeneous fracture plane, it may be pinned by asperities after some advancement into the heterogeneous region, if it is propagat-

ing with a relatively small velocity. If, however, the velocity is relatively high, the asperities give way, the fracture front becomes curvy and propagates further into the bordering homogeneous region, where it recovers a straight-line configuration through slowly-damped space-time oscillations which, if they are in response to spatially-periodic heterogeneities, decay as $t^{1/2}$ with time. Such a slow decay suggests that the configuration of a straight fracture front may be sensitive to even small but sustained heterogeneity in the fracture resistance (i.e., in the material).

Using the results of Rice *et al.* (1994), Perrin and Rice (1994) showed that a fracture propagating through a heterogeneous material, in which the heterogeneities are represented as randomly-distributed asperities with which the fracture front interacts continually, will *never* reach a statistically steady state. Instead, heterogeneities in the fracture energy lead to a logarithmic divergence of the root mean-squares deviations of an initially straight fracture front. In particular, the variance $V$ of the deviation of propagation velocity from the mean, was found to be

$$V \sim \log(2\alpha_0 v_0 t). \tag{180}$$

More interestingly, if the material is uniform over the remaining part of the fracture plane, after the encounter with the heterogeneous portion of the material, the propagating fracture becomes asymptotically (i.e., in the limit $t \to \infty$) straight again. These predictions suggest that perhaps the roughness of a fracture surface may be the direct result of a continuous roughening of the surface that is driven by small heterogeneities within the material. More recently, Willis and Movchan (1995) and Movchan and Willis (1995) computed the coupling of the energy release rate to random perturbations to the fracture front in the case of planar perturbations to the crack in Mode I fracture, and in shear loading. Willis and Movchan (1997) extended the analysis to the perturbations to the stress intensity factors induced by a small 3D dynamic perturbation of a propagating, nominally planar, fracture.

Ramanathan and Fisher (1997, 1998) calculated the dynamics of planar perturbations to a tensile crack front and found that, in contrast to the case of the scalar model for which Perrin and Rice (1994) had obtained logarithmic instability of the crack front, in Mode I fracture weak heterogeneity of the material can lead to a non-decaying unstable mode that propagates along the fracture front. They predicted that this propagating mode occurs in materials having $\partial\Gamma/\partial v \le 0$, where a constant value of $\Gamma$ is a marginal case. For $\partial\Gamma/\partial v > 0$, the propagating mode was predicted to decay, with the propagation velocity of the new mode being between $0.94c_R$ and $c_R$. These predictions are supported by the numerical simulations of Mode I fracture in a 3D material with a constant $\Gamma$, carried out by Morrissey and Rice (1998), indicating that the propagating mode is highly localized in space, and indeed propagates at the predicted velocities.

Ramanathan and Fisher (1997, 1998) and Morrissey and Rice (1998) both showed that these localized modes lead to linear growth of the root mean-square deviations of an initially straight fracture with its distance of propagation. They suggested that this may provide a new mechanism for the roughness produced by a propagating fracture in materials in which the fracture energy does not increase rapidly with the velocity of a crack. Both the calculations and simulations

were performed for in-plane disturbances to a fracture front. Disturbances of this type cannot, of course, generate the out-of-plane roughness typically seen along a fracture surface.

### 7.11.5  Failure of Dynamic Models of Cohesive Zone

Langer and collaborators (Barber *et al.*, 1989; Langer, 1992, 1993; Ching, 1994; Ching *et al.*, 1996a,b,c; Langer and Lobkovsky, 1998) carried out extensive theoretical studies of dynamic models of cohesive zone. They defined the cohesive zone in a manner similar to what was described in Section 7.11, but did not assume that cracks propagate at a constant rate, or always in a straight line, and therefore the cohesive zone becomes a dynamical entity which interacts with the fracture in a complex fashion. Their goal was to understand whether fracture tip instabilities can be predicted by such models.

In a first set of calculations, Barber *et al.* (1989), Langer (1992, 1993) and Ching (1994) studied the dynamics of cracks confined to straight lines, and found that such cracks always propagate in a stable fashion, which is consistent also with the predictions of Marder (1991), although there were also tantalizing hints of instabilities. Therefore, Ching, Langer and Nakanishi (1996a,b,c) studied dynamics of fractures that are allowed to follow curvy, out of plane, paths. In its most elaborate version, their model allows the fracture to pursue an oscillating path, and the cohesive zone to contain both tensile and shear components. In most, although not all of these models, fracture propagation is violently unstable to very short-length oscillations of the tip. Their general conclusion is that these cohesive-zone models are inherently unsatisfactory for use in dynamical studies. They are extremely difficult mathematically and they seem to be highly sensitive to details that, from a physical view point, ought to be unimportant. Pathological short-wavelength instabilities of fractures also emerge from their analysis which have a simple underlying explanation, which is as follows. The logic of the principle of local symmetry (see Section 7.6.1) states that atomic bonds under the greatest tension must break first, and therefore cracks loaded in Mode I propagate straight ahead, at least until a velocity, identified by Yoffe (see Section 7.6.3), is reached when a fracture is predicted to spontaneously break the symmetry inherent in straight-line propagation. This logic has been called into question by a very simple calculation, first described by Rice (1968).

To see this, let us look at the ratio $\sigma_{xx}/\sigma_{yy}$ right on the fracture line. Using Eqs. (51) and (52), we find that

$$\frac{\sigma_{xx}}{\sigma_{yy}} = \frac{(\beta^2+1)[1+2(\alpha^2-\beta^2)-4\alpha\beta]}{4\alpha\beta-(1+\beta^2)^2} = \frac{2(\beta^2+1)(\alpha^2-\beta^2)}{4\alpha\beta-(1+\beta^2)^2} - 1, \quad (181)$$

which after a Taylor expansion for low velocities $v$ becomes

$$\frac{\sigma_{xx}}{\sigma_{yy}} = 1 + \frac{v^2(c_t^4+c_l^4)}{2c_l^2c_t^2(c_l-c_t)(c_l+c_t)} + \cdots, \quad (182)$$

indicating that $\sigma_{xx}/\sigma_{yy}$ is greater than unity for all $v$ [$c_l$ and $c_t$ are defined by Eqs. (24) and (25)]. This result is surprising because it states that, in fact, as soon

as the fracture begins to propagate, the greatest tensile forces are *perpendicular* to its tip and not parallel to it. Therefore, it is difficult to imagine how a fracture can ever propagate in a straight line.

That Langer and collaborators found their dynamic models of the cohesive zone to be unsatisfactory may imply that, such models must be replaced by those in which plastic yielding is distributed across an area, and not restricted to a line. The two-field continuum models of the type described in Section 7.11.2 represent progress in this right direction. Another possibility is that calculations of Langer and co-workers indicate a fundamental failure of the continuum formulation of the type that they employ, and that the resolution must be sought either at the atomic or molecular scale (see Chapters 9 and 10), or one should resort to two-field continuum models that take into account the variations of the order parameter in the fracture zone.

## 7.12   Brittle-to-Ductile Transition

The last topic that we would like to briefly discuss is the brittle-to-ductile (BTD) transition that occurs in materials as the temperature is lowered and the strain rate is increased. Kelly *et al.* (1967) and Rice and Thomson (1974) were probably the first to offer a fundamental perspective on the class of materials that are capable of this fracture transition. In particular, Rice and Thomson developed a theoretical criterion for establishing the intrinsic brittle behavior and distinguishing it from intrinsic ductility. According to their criterion, an atomically sharp fracture governs the behavior of a material in the absence of any other form of plastic response in the background, by either (1) nucleating dislocations from its tip, or (2) by propagating in a cleavage mode due to the presence of an energy barrier to the emission of such dislocations. In the first class are intrinsically-ductile materials which cannot undergo a fracture transition, whereas the materials in the second group are usually considered as intrinsically brittle that are capable of making a transition to ductility. The BTD transition takes place at a characteristic temperature $T_{\mathrm{BTD}}$, and one main goal of research in this area has been developing a theory for quantitative prediction of this transition temperature. A variety of factors affect $T_{\mathrm{BTD}}$, with chief among them being the rate of loading the material. Many experimental studies (see, for example, Burns and Webb, 1970) indicate that mere nucleation of some dislocations from the tip of a fracture may not ensure ductile behavior. Despite this evidence, the Rice–Thomson mechanism resembles a threshold process, somewhat similar to the threshold nonlinearities that we have been considering in this book, that triggers ductility in a class of intrinsically-brittle materials in which the mobility of the dislocation is relatively high. Examples of such materials include BCC transition metals and most alkali halides. However, completely satisfactory confirmation of the Rice–Thomson criterion is rare.

Most models that are based on the Rice–Thomson criterion have been developed based on the assumption that, while background plastic relaxation serves to lower $T_{\mathrm{BTD}}$, the most important controlling factor of the transition temperature is the ability of the fracture tip to emit dislocations that can shield the entire fracture front

and hence trigger extensive plastic deformation *before* the fracture can propagate by cleavage. However, Argon (1987) showed that the Rice–Thomson-type models, in which the activation configuration consists of a fully-developed dislocation line, greatly over-estimate the energy barriers to nucleation of dislocations. This remains true even if one considers a modified Rice–Thomson-type model developed by Cheung *et al.* (1991) in which fracture tip nonlinearity and tension softening were incorporated.

On the other hand, consider the response of silicon, and many other similar co-valent compounds and materials, that have very sluggish dislocation mobility, and hence are in contrast with high-mobility hypothesis and the nucleation-controlled response of some materials. In such materials, the transition from brittleness to toughness is governed by the mobility of groups of dislocations that are away from the tip of the fracture (see, for example, St. John, 1975; Hirsch *et al.*, 1989; George and Michot, 1993). It is now well-established for both classes of materials that, the emission of the dislocations from the tip of a fracture occurs preferentially from specific sites on the tip, and that, in order to guarantee ductile behavior, the entire fracture front must be shielded from local break-out of the cleavage fracture from unprotected parts of the fracture front. Thus, it is now widely believed that the fundamental BTD transition is governed by the behavior of a cleavage crack.

In addition to the experimental studies mentioned above, theoretical analyses of fracture behavior of Si, carried out by Rice and Beltz (1994) and Xu *et al.* (1995), indicate that the activation configuration of dislocation embryo is a double kink of dislocation core matter. Thus, one may identify two distinct types of BTD transitions:

(1) In the BCC transition metals, where barrier to kink mobility along the dislocation are low, the BTD transition is governed by the formation of dislocation embryos at the fracture tip, which then results in a nucleation-controlled transition.

(2) By contrast, experimental work (see, for example, Yonenaga and Sumino, 1989) and theoretical modeling (Bulatov *et al.*, 1995) suggest that, in semiconductors and compounds the kink mobility is hindered by substantial energy barrier, hence rendering the BTD transition controlled by dislocation mobility away from the tip of the fracture.

A complete understanding of the BTD transition can be obtained based on atomistic modeling of the formation and outward propagation of the dislocation embryo at the tip of the fracture. Such atomistic modelings are based on 'molecular dynamics simulation that will be described in Chapters 9 and 10. However, atomistic models provide quantitative predictions for this phenomenon only if accurate potentials for describing the interatomic interactions are available. Several promising interatomic potentials have been developed over the past decade or so that will be described in Chapter 9. Alternatively, one may utilize a multiscale modeling approach—one that combines continuum modeling for the region away from the fracture tip with atomistic simulations in the tip region—in order to study this phenomenon. This represents a realistic and powerful approach that is rapidly gaining

popularity; Chapter 10 will describe this method. So far as the BTD transition is concerned, Xu *et al.* (1995) have already developed a multiscale model for studying this phenomenon. They showed that the energetics of the dislocation embryo formation on inclined slip planes that contain the fracture tip, when compared with an additional surface production resistance, is quite unfavorable and cannot explain the known BTD transition temperatures. Xu *et al.* conjectured that nucleation may be more favorable on oblique slip planes, or may occur *heterogeneously* at the edges of the fracture front. However, we must realize that, although dislocation nucleation on oblique planes has often been suggested as a likely scenario, approximate analyses that were based on the Rice–Thomson criterion have led to estimates of $T_{BTD}$ that are several orders of magnitude larger than the experimental values.

We note that, although experiments have established the ability of dislocation nucleation at the fracture tip for accounting for the exceedingly sharp BTD transitions in Si and similar materials, Khanta *et al.* (1994) questioned this well-understood fact, and instead advocated an approach based on an analogy with thermal phase transitions. Specifically, they considered, unlike the more traditional methods described above, the thermally-induced instability of many small loops in the presence of an applied stress, and proposed that the creation of many atomic-size loops by thermal activation induces a temperature-dependent *cooperative screening effect* that enhances the subsequent growth of the loops. This cooperative effect is completely different from the dislocation shielding of fracture tip stress described above. To develop their theory, they extended the concept of dislocation screening, originally developed by Kosterlitz and Thouless (1973) in an entirely different context, namely, 2D phase transitions. In the Kosterlitz–Thouless (KT) theory, the generation of dislocations (which is an unstable process) is driven by only thermal fluctuations, without the aid of an applied stress. The KT transition occurs at a temperature close to the melting temperature, which then gives rise to a dislocation-mediated melting transition (Nelson and Halperin, 1979; Young, 1979). In the model developed by Khanta *et al.* (1994), both the external stress and thermal fluctuations assist the growth of dislocation loops. The model then predicts the existence of a KT-type instability, but not a phase transition in the thermodynamic sense, at a temperature well below the melting temperature, at a stress level that corresponds to the Griffith threshold that is needed for brittle fracture propagation. This temperature is then identified with $T_{BTD}$. If the transition temperature is zero and the applied load is equal to the Griffith threshold, the model reduces to the Rice–Thomson model described above. Thus, one advantage of this theory is that it is applicable to systems that are at a finite temperature, in contrast with the Rice–Thomson model that is strictly valid for zero temperature. Despite this success, there is not yet convincing evidence for the role of thermal fluctuations advocated by Khanta *et al.* (1994). Indeed, the meticulous experiments of George and Michot (1993), who used X-ray direct imaging of the stages of evolution of the fracture-tip plastic response, starting from nucleation of crack tip heterogeneities and followed by very rapid spread and multiplication of dislocation length from such sources, demonstrate clearly the vast numbers of degrees

of freedom available to dislocation for populating the highly-stressed fracture tip, but do not indicate any significant role for thermal fluctuations.

Finally, we note that there are many morphological aspects of a BTD transition in polycrystalline materials in which microcracks, nucleation and crack arrest at grain boundaries become very important, and modulate the actual $T_{BTD}$. Our understanding of such processes is still not complete, and therefore this is an active research area (see, for example, Falk and Langer, 1998; Falk, 1999).

## Summary

As stated at the beginning of this chapter, it was believed for a long time that there is a conceptual problem with the continuum mechanical formulation of brittle fracture of amorphous materials, as its prediction for the terminal velocity of propagating fractures, i.e., the Rayleigh wave speed $c_R$, had seemed to be experimentally unattainable (apart from highly anisotropic materials). However, the discussions of this chapter should have made it clear that the problem persisted not because of a fault in the continuum mechanics, but because it had not been properly posed. The correct question should have been about the nature of energy dissipation near the fracture tip. However, such a problem was not studied for several decades, because it had seemed natural to assume that, in a sufficiently brittle material, energy will be consumed mainly for breaking the atomic bonds and generating new fracture surface, a process that should depend only weakly on the fracture velocity. However, by loading fractures in differing fashions, greatly-fluctuating quantities of energy can be forced into the fracture tip. The tip must then find some mechanism for dealing with the energy not needed to break a minimum set of atomic bonds. A small fraction of the remaining energy is consumed by such minor events as phonon emission, after which the tip begins consuming energy by a sequence of dynamical instabilities, giving rise to ramified networks of fractures (or broken atomic bonds) on small length scales.

Thus, there is actually no discrepancy between the conventional continuum fracture mechanics and the experimental observations and data. In a large enough amorphous material, the fracture-tip instabilities occur within the cohesive zone where linear continuum fracture mechanics is not even an appropriate theoretical framework for analyzing the instabilities, let alone predicting them. The finite-element simulations, models of fracture propagation in 3D, the two-field continuum models, the lattice models that will be described in Chapter 8, and many precise and beautiful experiments carried out over the past decade, have now provided us with a much better understanding of the structure and dynamics of energy dissipation in the vicinity of the tip of a propagating fracture in a brittle material. It is now clear that fracture in brittle materials is governed by a dynamic instability that gives rise to repeated attempts for branching off of the main propagating fracture, hence preventing the terminal fracture velocity from reaching the Rayleigh wave speed.

# 8
# Brittle Fracture: The Discrete Approach

## 8.0  Introduction

As discussed in Chapters 6 and 7, theoretical and computer simulation studies of fracture of materials are usually based on one of the following three approaches.

(1) The first approach formulates the problem using linear continuum fracture mechanics. This approach, which was described in detail in Chapter 7, allows one, in many cases, to derive the analytical solution of the problem of fracture propagation in a given material, subject to certain initial and boundary conditions. If, however, such analytical solutions cannot be derived, then the governing equations must be discretized by, for example, a finite-difference or finite-element method and solved by numerical simulations, in which case the model reduces to a type of discrete or lattice model.

(2) The second approach is based on molecular dynamics (MD) simulation of fracture propagation which studies the phenomenon at atomic length scales. Molecular dynamics is a discrete approach in that, the system under study is represented by a discrete set of atoms connected to one another by atomic bonds. This approach will be described in Chapter 9.

(3) The third approach is based on lattice models which can be used for both quasi-static and dynamic fracture phenomena. However, we must point out that there is a major difference between lattice models of fracture that we describe and discuss in this chapter and the MD approach to fracture. The difference is due to the fact that, in MD simulation of fracture breaking of an atomic bond is a natural outcome of the simulations, whereas in the lattice models described in this chapter, how or when a bond breaks is an *input* of the models that must be specified at the outset. There are, in general, two types of lattice models.

  (i) One class of such models is intended for quasi-static fracture. Such models consist of a lattice of springs or beams, together with a criterion for nucleation of local microcracks. In these models, each node of the lattice is connected to only a finite number of other sites (which are usually the nearest-neighbor sites), and a force balance is written down for each node, resulting in a set of simultaneous equations that govern the nodal displacements. Unlike the MD method, the nodes of the lattice do not represent the material's atoms, nor do the bonds represent the atomic bonds. Instead,

the lattice models represent a material at length scales much larger than the distance between two neighboring atoms in the material, and therefore one does not have to be concerned about developing accurate interatomic potentials between the atoms, a subject that will be discussed in detail in Chapter 9.

(ii) The second class of such models are intended for dynamic fracture. This class of models is itself divided into two subclasses. (a) In one group are models that represent generalization of the lattice models of quasi-static fracture. The nodes of the lattice do not represent atoms. Some of such models contain *quenched* (fixed in space) disorder, while others have been developed for fracture of materials with *annealed* disorder (i.e., one that may change with time). (b) The lattice sites in the second group do represent atoms. However, instead of assuming interatomic potentials between the atoms, as in MD simulations, one adopts, in a manner similar to lattice models of quasi-static fracture, a simple force law between the atoms, one in which the forces rise linearly up to a critical separation between the atoms, beyond which they abruptly vanish. If the lattice contains no disorder, then *exact* calculations can be carried out (see below).

In essence, most of these models represent generalizations of the lattice models for linear transport properties of heterogeneous materials (described in detail in Volume I), and also those for the phenomena of electrical and dielectric breakdown described in Chapter 5. Aside from the fact that for certain materials, such as fibrous composites, lattice models are natural, the motivation for developing such models of brittle fracture is twofold.

(1) In most materials, either manufactured (such as composite solids) or natural (such as rock), the presence of heterogeneities in the form of either a distribution of microscopic elastic constants, or in terms of flaws or defects with various sizes, shapes and orientations, makes fracture a very complex phenomenon. Thus, as already pointed out in Chapter 6, the effect of even small initial disorder can be enormously amplified during fracture, with the result being the fact that fracture is a *collective* phenomenon which is controlled by the disorder. In fact, due to disorder, especially when it is strong, brittle materials generally exhibit large statistical fluctuations in their fracture strengths, when nominally identical samples are tested under *identical* loading. Thus, as is now well-understood, due to the fluctuations, it is inappropriate to analyze the phenomena of fracture of a disordered material by a mean-field theory or an effective-medium approximations. Incorporating the effect of disorder in a continuum model of dynamic or even quasi-static fracture is, however, a daunting task, especially when the heterogeneities are broadly distributed. In addition, such lattice models allow one to investigate, in a convenient and meaningful manner, various properties of the *morphology* of the networks of microcracks that are formed, e.g., those that are formed in rock and rock-like materials, such as concrete.

(2) Over the past fifteen years there has been considerable theoretical progress towards understanding the dynamics of elastic manifolds moving through disordered media, such as charge density waves (see, for example, Narayan and Fisher, 1992), fluid-surface contact lines (see, for example, Ertas and Kardar, 1992), and interfaces between two phases, such as those that are encountered in multiphase flow in a disordered porous medium (see, for example, Sahimi, 1993b, 1995b), all of which exhibit a sort of *non-equilibrium* critical phenomenon close to the onset of motion. Fracture of materials does have similarities with these phenomena (although it has important differences too) which have provided the impetus for developing some of the models that were described in Chapter 7, and those that will be described in the present chapter. In particular, one is interested to understand the extent of the similarities between these seemingly different phenomena, so that the possibility of a unified approach to most, if not all, of them can be explored. Moreover, if such similarities do exist, then the knowledge that already exists about some of such phenomena can be immediately "transferred" into new insight about fracture phenomena.

To make this point clearer, let us go back to Chapter 7 and recall the essentials of brittle fracture phenomenon. Suppose that there exists a crack front in a material and that an external load $\sigma$ is applied to it. If $\sigma$ is small, there is no steady-state motion and the crack front is *pinned* by the heterogeneities of the material in one of the many locally-stable configurations. As the external load increases, there are a series of local instabilities that become larger as $\sigma$ increases further. At a critical load (stress) $\sigma_c$ the crack front *depins* and begins to move. In a large enough system, the transition from the stationary to the moving state exhibits features of a non-equilibrium dynamic critical phenomenon which, to some extent, are similar to those of second-order phase transitions, such as the percolation transition emphasized in this book. For example, the mean velocity $v$ of the moving fracture just above $\sigma_c$ obeys the following power law (Ramanathan and Fisher, 1997):

$$v \sim (\sigma - \sigma_c)^{\zeta}, \tag{1}$$

where $\zeta$ is a critical exponent which is, hopefully, independent of many microscopic properties of the material. Moreover, in the quasi-static case, as $\sigma$ increases, segments of the crack front overcome the local toughness caused by the heterogeneities and move forward, causing other segments to jump, thereby triggering an avalanche which will eventually be stopped by tougher regions. It has been found that, up to a characteristic length $\xi^-$, the avalanches exhibit a power-law size distribution, where by size we mean roughly the extent $l$ along the crack front of an avalanche. This size distribution is given by

$$P(\text{size} > l) \sim l^{-\kappa} f(l/\xi^-), \tag{2}$$

where $\kappa$ is a characteristic critical exponent. The cutoff length scale $\xi^-$ itself obeys the following power law near $\sigma_c$:

$$\xi^- \sim (\sigma_c - \sigma)^{-\nu^-}, \tag{3}$$

where $\nu^-$ is the critical exponent associated with $\xi^-$. Note that the cutoff length scale $\xi^-$ plays a role similar to $\xi_p$, the correlation length of percolation which, as has been emphasized throughout this book, plays a fundamental role in determining the length scale over which materials with percolation heterogeneity can be considered as homogeneous. Moreover, we expect that

$$\int_0^{\sigma_c} l^{-\kappa} f(l/\xi^-) \, d\sigma_c \sim l^{-1}. \tag{4}$$

Just above $\sigma_c$, the fluctuations in the crack velocity are correlated up to a length scale $\xi^+$ which follows another power law given by

$$\xi^+ \sim (\sigma - \sigma_c)^{-\nu^+}. \tag{5}$$

In general, we expect $\nu^- = \nu^+ = \nu$ (see Chapter 3 for examples for which this is not true). As discussed in Chapters 6 and 7, at the threshold $\sigma_c$ the fracture surface has a self-affine structure with a roughness exponent $\alpha$, so that the correlation function $C(r)$ scales as,

$$C(r) \sim r^{2\alpha}. \tag{6}$$

Finally, the time scale $t_l$ that an avalanche of size $l$ lasts is characterized by a dynamic exponent $z$, similar to what was defined in Chapter 2:

$$t_l \sim l^z. \tag{7}$$

Not only are these exponents well-defined, but also satisfy certain scaling relations. In fact, Ramanathan and Fisher (1997) showed that

$$\zeta = (z - \alpha)\nu, \quad \nu = (1 - \alpha)^{-1}, \tag{8}$$

so that, similar to percolation and other second-order phase transitions, there are only two independent exponents that characterize this transition. Two-dimensional (2D) numerical simulations of Ramanathan and Fisher (1997) yielded, $z \simeq 0.74$, $\alpha \simeq 0.34$, $\nu \simeq 1.52$, and $\zeta \simeq 0.34$. The estimated $\alpha$ is smaller than the typical value of the roughness exponent, $\alpha \simeq 0.8$, that, as discussed in Section 7.8.7, has been reported for several classes of materials. However, MD simulations of fracture by Nakano *et al.* (1995), to be described in Chapter 9, indicate that, in agreement with our discussion in Chapter 7, there may be two regimes of fracture propagation, characterized by different roughness exponents. Nakano *et al.* found that at the initial stages of fracture propagation, when the crack tip moves slowly, $\alpha \simeq 0.44$, which is reasonably close to the estimate of Ramanathan and Fisher (1997), while at latter stages when fracture propagation proceeds at relatively high speeds, $\alpha \simeq 0.8$.

In addition, the lattice models that are described in this chapter have enabled us to resolve the conflicts between the predictions of linear continuum fracture mechanics and the experimental observations. In particular, the phenomena of fracture instabilities, microbranching, and the inability of a propagating fracture for reaching the Rayleigh wave speed $c_R$ (the experimental aspects of which were

described in detail in Chapter 7) have been explained in a satisfactory manner by such lattice models of dynamic fracture.

We begin this chapter by discussing important aspects of models of fibrous materials and the predictions that they have provided. We then describe in detail lattice models of quasi-static brittle fracture, and the considerable insight that they have provided into the fracture of *heterogeneous* materials, after which lattice models of dynamic fracture are described and discussed. As usual and whenever possible, we compare the predictions of the models with the relevant experimental observations and data.

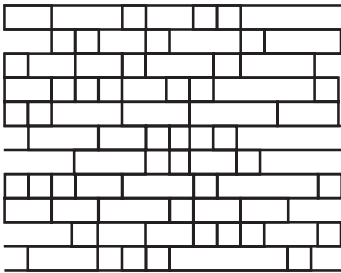## 8.1    Quasi-static Fracture of Fibrous Materials

As our discussions in Chapter 7 indicated, despite *decades* of effort, there are very few exact results for fracture dynamics of disordered materials. Exact analytical analysis of fracture of any type of material, regardless of whether a discrete model is used or linear or nonlinear continuum mechanics is employed, is a complex task. Moreover, quasi-static fracture processes are sensitive to the sample size, but the approach to their asymptotic (large sample size) behavior is slow. At the same time, numerical simulation of quasi-static fracture (of the type that is discussed in this chapter) in very large systems is currently very difficult, if not impossible. Thus, an exact solution of the fracture problem in *any* physically viable system would be very valuable, as it would shed light on a very complex process.

Some of the early work on fracture phenomena concentrated on tensile failure of continuous-fiber composites using relatively simple models (see, for example, Daniels, 1945; Coleman, 1958). The reason for this was twofold. One was the wide applications that such materials have, ranging from paper to glass-fiber mats. In addition, many composite materials of industrial importance are reinforced by rigid fibers. The second reason for these early studies was that some of the relatively simple models developed for such materials, which could provide insight into their fracture process, are amenable to analytical analysis. Hence, study of fracture of such materials has remained an active research field (see, for example, Harlow and Phoenix, 1978, 1991; Smith *et al.*, 1983; Phoenix and Smith, 1983; Curtin, 1991; Phoenix and Raj, 1992; Åström *et al.*, 1994, 2000; Kellomäki *et al.*, 1996; Räisänen *et al.*, 1997). Some of these studies involved analytical computations of mechanical and fracture properties of fibrous materials, while others, which also used more realistic models of such materials, utilized large-scale computer simulations. We aim to describe the important results that have emerged from such studies, starting with the analytical results.

One of the rare models for which an *exact* analysis can be carried out is the fiber-bundle model, the simplest example of which is shown in Figure 8.1. The tensile stress is applied vertically. Suppose that $p$ and $q = 1 - p$ are the fractions of the bonds that are present (unbroken) and absent (broken or failed) in the bundle, respectively, and that each bond is characterized by a failure stress $\sigma_f$. One can construct a 2D model of such fibers by putting together $L$ of such bundles,

FIGURE 8.1. Fiber bundle (top) and chain-of-bundles model (bottom).



which is also shown in Figure 8.1. The survival probability $p_s$ (the probability that the bundle does not fail macroscopically) is then, (survival probability of a 1D bundle)$^L$.

The model is physically viable only if the applied stress or strain is shared by the bonds in a meaningful manner, and thus the issue of *load sharing* is critical. As discussed by Duxbury and Leath (1994a), there are two classes of such load-sharing models which we now describe and analyze.

## 8.1.1  Equal-Load-Sharing (Democratic) Models

In this class of models, also called the *democratic* models, the load carried previously by a failed bond is shared equally by all the remaining bonds in the system (Daniels, 1945; Harlow and Phoenix, 1978). As simple as it may seem, this model might be applicable to a variety of materials, such as cables or ropes made of numerous fibers, and even geological faults that are locked by asperity barriers sharing the total stress. The democratic model of failure of the material is a type of an effective-medium or a mean-field approximation, and has been used in a variety of situation, such as modeling of ceramic-matrix continuous-composites. Because of their mean-field nature, such models can often be solved exactly. Here, we briefly describe the solution for such models which is due to Sornette (1989).

Consider $n$ independent vertical fibers with identical spring constant $\kappa^{-1}$ but random failure threshold $X_i, i = 1, 2, \cdots, n$. Suppose that the total stress exerted on this system is $\sigma$, and that the strengths $X_1, X_2, \cdots$ of the individual links are independent and randomly distributed variables with the cumulative distribution $P(X_j < x) = F(x)$. Under a total load $\sigma$, a fraction $F(\sigma/n)$ of the threads will be submitted to more than their rated strength, and therefore will fail (break) immediately, after which the total load will be redistributed by the transfer of stress from the broken links to the unbroken ones, which will then induce secondary failures, and so on. Thus, one has a cascade of induced failure which we would like to describe. An important question to be answered is: Does the cascade stop at some point or propagate until the entire system fails? The answer does, of course, depend on the way the total stress is redistributed each time a link or bond fails. Although

the democratic model may appear to be difficult but, as pointed out by Sornette (1989), it can in fact be solved by using the theory of extreme order statistics which was also used in our discussion of models of electrical and dielectric breakdown of materials in Chapter 5. The key idea is that, the bundle will not break under an external load $\sigma$ if there are $k$ links in it, each of which can withstand a load $\sigma/k$. In other words, if $X_{1;n} \leq X_{2;n} \leq \cdots \leq X_{n;n}$ is the way in which the strengths of the individual links are ordered, then, if the first $k-1$ weakest links fail, the bundle will resist macroscopic failure under a stress $\sigma_n \leq (n-k+1)X_{k;n}$, because of the remaining $(n-k+1)$ links of breaking strength $\geq X_{k;n}$. Therefore, the strength $\sigma_n$ of the bundle is given by

$$\sigma_n = \max\{(n-k+1)X_{k;n}; \quad 1 \leq k \leq n\}. \tag{9}$$

We now search for the *strongest subgroup* of the bonds. The variables $X_{k;n}$ are strongly dependent since they are correlated. However, regardless of the specific form of $F(x)$, there is a very general result for $\sigma_n$ due to Galambos (1978) which is as follows.

**Theorem:**    Suppose that $F(x)$ is an absolutely continuous function with finite second moment, and that $x[1-F(x)]$ has a unique maximum at $x = x_0 > 0$ such that $y_0 = x_0[1-F(x_0)]$. If $F(x)$ has a positive second derivative in the neighborhood of $x_0$, then as $n \to +\infty$, one has

$$\lim_{n \to \infty} P(\sigma_n < ny + x\sqrt{n}) = (2\pi)^{-1/2} \int_{-\infty}^{x} \exp\left(-\frac{1}{2}z^2\right) dz, \tag{10}$$

which is essentially a central-limit theorem. Equation (10) implies that

$$P(\sigma_n = \sigma) \sim (2\pi n x_0)^{-1/2} \exp[-(\sigma - ny)^2/2nx_0^2]. \tag{11}$$

Equation (11) states that the density distribution of the global failure threshold is Gaussian around the maximum $\sigma = ny$ with a variance that scales as $n$, hence implying that the typical strength of the system increases as $\sigma_n \sim n$, if $n$ is large. Although by a naive argument one may predict that $\sigma_n = n\langle x \rangle$, where $\langle x \rangle$ is the mean one-link threshold, Eq. (11) shows that $\sigma_n = ny$, with $y$ being in fact significantly smaller than $\langle x \rangle$, and therefore the naive argument greatly overestimates the global failure threshold.

The mechanical characteristics of the system under a given applied stress $\sigma < \sigma_n$ depend upon the history of the system, i.e., on the number and the way the links have failed as the stress was increased from zero to $\sigma$. With each value of $\sigma < \sigma_n$ we associate an integer $m(\sigma)$ with $1 \leq m(\sigma) \leq n$ such that

$$[n-m(\sigma)+2]X_{m-1;n} \leq \sigma \leq [n-m(\sigma)+1]X_{m;n}, \tag{12}$$

which can be rearranged to

$$\{1-[m(\sigma)-2]/n\}X_{m-1;n} \leq \sigma/n \leq \{1-[m(\sigma)-1]/n\}X_{m;n}. \tag{13}$$

Note that $m(\sigma) - 2$ is the number of links which have failed under a stress $\leq$ $\sigma/[n - m(\sigma) + 2]$. Moreover, by definition of $F(x)$,

$$(m - 2)/n \leq F[\sigma/(n - m + 2)] \leq (m - 1)/n, \tag{14}$$

which follows from the fact that, for large $n$, counting the number of links with failure threshold less than $\sigma/[n - m(\sigma) + 2]$ amounts to computing the cumulative failure distribution $F(x)$ at $x = \sigma/[n - m(\sigma) + 2]$. Relations (13) and (14) indicate, roughly speaking, that, as $n \to \infty$, $\sigma/n$ is increasingly better approximated by $x[1 - F(x)]$ with

$$\frac{\sigma}{n} = x(\sigma)\{1 - F[x(\sigma)]\}. \tag{15}$$

Note that Eq. (15), in the limit $n \to \infty$, is a continuous function. It is then not difficult to show that, for large $n$, the number of links which have failed under $\sigma$ is given by

$$k(\sigma) = nF[x(\sigma)]. \tag{16}$$

For large but finite $n$, $\sigma(x)$ or $x(\sigma)$ is a staircase with plateaux of width decreasing to zero as $n \to \infty$. The width of each plateau, for a given $\sigma$, can be obtained from (13), since the interval in $\sigma$ is such that (13) holds with the same integer $m(\sigma) = m$.

Just before complete failure of the bundle, the total number of failed links is given by

$$k_n = k(\sigma_n) = nF(x_0), \tag{17}$$

implying that a finite fraction of the links fail before global rupture occurs. If we consider, for example, the (cumulative) Weibull distribution (WD) (see also Chapter 5),

$$F(x) = 1 - \exp[-(x/\lambda)^m], \tag{18}$$

where $\lambda$ and $m$ are the parameters of the distribution, then

$$\frac{k_n}{n} = 1 - \exp(-1/m), \tag{19}$$

which for $m = 2$ yields $k_n/n = 0.393$. For $\sigma \leq \sigma_n$, $x(\sigma)$ is in neighborhood of $x_0$ and may be expressed as

$$x(\sigma) = x_0 - A(y - \sigma/n)^{1/2}, \tag{20}$$

where $A$ is a constant with a value that depends on the shape of $F(x)$. For example, for the WD, $A = [x_0 \exp(1/m)/m]^{1/2}$. Then, the number of links that have failed under the stress $\sigma$ is given by

$$\frac{k(\sigma)}{n} = F(x_0) - B(y - \sigma/n)^{1/2}, \tag{21}$$

where $B$ is another constant. For example, for the WD, $B$ is given by $B = (mx_0)^{-1/2} \exp(-1/2m)$. Equation (21) indicates that $k$ increases rapidly as $\sigma \to \sigma_n$, approaching $nF(x_0)$ with a square-root singularity.

We can thus predict the strain-stress characteristics of the bundle of the fibers. Suppose that each individual link is made of a brittle material, so that its strain-stress relation is given by, $\epsilon_l = \kappa \sigma_l$ up to its failure point, where $\epsilon_l$ is the strain. Then,

(1) for $\sigma \le \sigma_1$, where $\sigma_1$ is the strength of the weakest link (the first to fail), all links are intact and the system has a linear stress-strain characteristic with slope $\kappa^{-1}$. Note that for the WD, $\sigma_1 \sim \lambda n^{-1/m}$.

(2) For $\sigma_1 \le \sigma \le \sigma_n$, some of the links have failed, and the system is elastic but nonlinear, which can be established by the following argument. We see from Eq. (16) that $n\{1 - F[x(\sigma)]\}$ links support the total external stress $\sigma$, which means that the stress *per remaining link* is given by

$$\sigma_r = \frac{\sigma}{n\{1 - F[x(\sigma)]\}} = x(\sigma). \tag{22}$$

Thus, for every $\sigma_r$ there is a corresponding strain per link $\epsilon_r$, which is equal to the strain of the entire bundle of links associated in parallel, and is given by

$$\epsilon_r = \kappa x(\sigma), \tag{23}$$

and therefore we have a strain-stress characteristic which becomes flat with zero slope as the global failure threshold is approached, $\sigma \to \sigma_n$. Hence, the effective elastic modulus of the system decreases as $\sigma$ increases. This nonlinearity is due to the fact that as $\sigma \to \sigma_n$, more and more links fail and therefore the total external stress is transferred to fewer and fewer links. The stress transfer is of course a nonlinear process. The nonlinear behavior of the system is characteristic of an irreversible process, with the irreversibility in the present problem being the deterioration of the bundle as $\sigma \to \sigma_n$.

Note that the failure transition in the democratic model is abrupt and hence it represents a first-order phase transition. There is a rapid increase in the number of the failed links as the global failure point is approached. If we assume $F(x)$ to be a WD with $m = 2$, then the value of $F$ at the failure threshold is $F = 0.168$, implying that, before the global failure threshold, few precursory failures have taken place. Thus, in a sense, the system fails without any "warning."

## 8.1.2 Local-Load-Sharing Models

In this class of models the stress carried previously by a failed bond is shared *locally* by the remaining bonds in its vicinity, which is of course what happens in most real materials. Suppose that the total number of bonds in a bundle, the sum of the intact and failed ones, is $L$. Since a defect or vacant cluster grows as bonds at its ends fail, catastrophic failure occurs as soon as a bond fails. Therefore, all one must do is finding the bond that suffers the largest stress enhancement, and adjusting the external stress until this bond fails. The adjusted stress is then the fracture stress $\sigma_f$ of the bundle as a whole. In practice, this is easier said than

done, because failure depends on the largest vacant cluster the statistics of which are difficult to analyze.

An elegant solution of this problem was developed by Duxbury and Leath (1994a) (for the solution of the problem in which the stress carried previously by a failed fiber is shared by its nearest and next-nearest neighbors, see Phoenix and Beyerlein, 2000). We present a brief description of their solution. With the cluster-end-load-sharing rule, the bond which suffers the largest stress enhancement is one at the end of the largest cluster of the absent bonds. Under this scenario then, the survival probability is related to the probability $P_L(n)$ that there is no cluster of vacant bonds of size greater than some prescribed value $n$. An important load sharing rule is that, $\sigma_t = \sigma(1 + \frac{1}{2}n)$, where $\sigma_t$ is the stress at the tip of the failed bond. Duxbury and Leath (1994a) calculated $P_L(n)$ following a method proposed by Harlow (1991) in which one identifies the possible endings of a fiber bundle of length $L + 1$, and the way by which these endings may be generated from a bundle of length $L$. In essence, this method is similar to the transfer-matrix technique described in Section 5.14.2 of Volume I. Suppose that {1} stands for a present (unbroken) bond and {0} for an absent (failed) one. If the size of the vacant sites is restricted to be $n$, then the bundle endings that are allowed are (1), (10), (100), (1000$\cdots$), where the number of zeros in the last probability is $n$. One now constructs a transition probability matrix for going from each of these possible configurations at the end of a bundle of length $L$ to the same endings in a bundle of length $L + 1$, by considering the probability of their occurrence. For example, the probability of going from ending (1) to ending (10) is $q$, since the probability that the next bond added is vacant is just $q$. We define $\mathbf{P}_L^T = [p_{(1)}, p_{(10)}, p_{(100)}, \cdots, p_{(100\cdots0)}]$ as the probability vector of having the set of possible endings on a fiber bundle of length $L$. Then $\mathbf{P}_{L+1}$ is obtained from $\mathbf{MP}_L = \mathbf{P}_{L+1} = \mathbf{M}^L\mathbf{P}_1$, where

$$
\mathbf{M} = \begin{bmatrix}
p & p & p & \cdots & p \\
q & 0 & 0 & \cdots & 0 \\
0 & q & 0 & \cdots & 0 \\
\cdot & \cdot & \cdot & 0 & 0 \\
0 & 0 & \cdots & q & 0
\end{bmatrix}
\tag{24}
$$

is called the *transition matrix*. Then, the probability $P_L(n)$ that there are no vacant clusters of size larger than $n$ is found from

$$
P_L(n) = \sum_l (p_l)_L.
\tag{25}
$$

One may use a variety of boundary conditions, the simplest of which is perhaps the periodic conditions which require that the first and the last site of the bundle to be equivalent, in which case

$$
P_L(n) = \text{tr}(\mathbf{M}^L),
\tag{26}
$$

where tr denotes the trace of the matrix. Thus, all one must do is studying the eigenvalues of $\mathbf{M}$. Let $a_1 = p/\lambda$ and $a_2 = q/\lambda$, where $\lambda$ is the eigenvalue of $\mathbf{M}$.

If we define a determinant $D_n$ by

$$D_n = \begin{vmatrix} a_1 - 1 & a_1 & \cdots & a_1 & a_1 \\ a_2 & -1 & 0 & \cdots & 0 \\ 0 & \cdots & & & \\ \cdots & \cdots & a_2 & -1 & 0 \\ 0 & \cdots & 0 & a_2 & -1 \end{vmatrix} \tag{27}$$

then

$$D_n = -D_{n-1} + (-1)^n a_1 a_2^n, \tag{28}$$

with $D_0 = a_1 - 1$. The solution to the recursion relation (28) is

$$(-1)^n D_n = a_1 - 1 + a_1 a_2 + a_1 a_2^2 + \cdots + a_1 a_2^n = 0. \tag{29}$$

It is then easy to see that

$$\lambda^{n+2} - \lambda^{n+1} + pq^{n+1} = 0. \tag{30}$$

Because $\mathbf{M}$ is non-negative, then according to the Perron–Frobenious theorem (see, for example, Noble and Daniel, 1977) its largest eigenvalue $\lambda_\ell$ is real and unique. Moreover, it is not difficult to see that $\lambda_\ell \to 1$ as $n$ becomes large. Therefore, setting $\lambda_\ell = 1 - \delta$, Eq. (30) yields

$$\lambda_\ell \simeq 1 - pq^{n+1} + O(q^{2n}). \tag{31}$$

and hence for periodic boundary conditions

$$P_L(n) = \text{tr}(\mathbf{M}^L) = \lambda_1^L + \lambda_2^L + \cdots + \lambda_n^L \simeq \lambda_\ell^L + O(|\lambda_{s\ell}|^L), \tag{32}$$

where $\lambda_{s\ell}$ is the second largest eigenvalue of $\mathbf{M}$. We thus obtain

$$P_L(n) = [1 - pq^{n+1} + O(q^{2n})]^L + O(|\lambda_{s\ell}|^L). \tag{33}$$

This result agrees with what Duxbury *et al.* (1986) derived for the electrical breakdown problem discussed in Section 5.2.5. We can now find the failure probability $p_f$ when a stress $\sigma$ is applied to the bundle by noting that, since failure of the bond that carries the largest stress causes catastrophic failure, we must have

$$p_f = \frac{\sigma_f}{\sigma} = 1 + \frac{1}{2}n, \tag{34}$$

where $\sigma_f$ is the failure stress. Therefore, the probability $p_s$ that the fiber bundle will survive is

$$p_s(\sigma) = \left(1 - pq^{2\sigma_f/\sigma - 1}\right)^L. \tag{35}$$

If $L$ and $n$ are large, then Eq. (35) is essentially equivalent to a double exponential form, also called a Gumbel distribution, a result that was also obtained for electrical and dielectric breakdown phenomena described in Section 5.2.5.

A more complex situation arises when an intact bond is between two clusters of vacant bonds, in which case the bond suffers a large stress enhancement. Thus, for a more complete analysis one must also consider this situation. The same

technique that was described above can be used to analyze this case, except that some modifications must be made. For example, the distinct endings that must be considered are (11), (110), (1100),$\cdots$,(110$\cdots$0); (101), (1010), (10100),$\cdots$, (10100$\cdots$0); (1001), (10010), $\cdots$, and (10$\cdots$0), each of which occurs with a certain probability analogous to $p_{(10)}$, $p_{(100)}$, and so on. Duxbury and Leath (1994a) then showed that these more complex configurations do not change the essence of their analysis described above. After some algebra one obtains

$$P_L(n) \simeq \left\{ 1 - [(n+1)p^2 - pq]q^{n+1} + O(q^{3n/2}) \right\}^L, \tag{36}$$

and the probability of survival is given by

$$p_s(\sigma) = \left[ 1 - \left( \frac{2\sigma_f}{\sigma}p^2 - p \right) q^{2\sigma_f/\sigma - 1} \right]^L. \tag{37}$$

Observe that, compared to (33) and (35), only some prefactors are different in (37). The average strength of the fiber bundle can then be calculated as

$$
\begin{aligned}
\frac{\langle \sigma \rangle}{\sigma_f} &= \sum_{n=0}^{L-1} \frac{2[P_L(n) - P_L(n-1)]}{n+2} \\
&= \frac{2P_L(n)}{L+1} - \frac{L^2 pq^{L+1}}{(L+1)(L+2)} + \sum_{n=1}^{L-1} \frac{2P_L(n-1)}{(n+1)(n+2)},
\end{aligned}
\tag{38}
$$

where the second term on the right side of the second equation represents a correction term for preventing (38) from having unphysical behavior as $L$ becomes large.

   In two other papers, Duxbury and Leath (1994b) and Leath and Duxbury (1994) developed interesting recursion relations for calculating the failure probability and average strength of the fiber-bundle model, so that one can numerically study the behavior of the model [for a different approach, based on calculating the Green functions, see Zhou and Curtin (1995); for a Green function analysis of fracture in more general systems see also Zhou *et al.* (1993)]. As usual, suppose that {1} denotes an intact (unbroken) bond and {0} a failed one. Then for $L = 2$ the surviving configurations are {11, 10, 01}, while for arbitrary $L$ there are $2^L - 1$ surviving configurations and one failure configuration {0$\cdots$00}. The probability $p_{sn}$ that a bond with $n$ failed neighbors survives is $p_{sn} = 1 - \int_0^{(1+n/2)} q(x)dx$, where $q(x)$ is the differential failure probability of a bond. Duxbury and Leath (1994b) separated the full set of $2^L - 1$ survival configurations into judiciously selected subsets. Suppose that a lone surviving fiber is surrounded by failed fibers, and let {A} be the set of all survival configurations which contain only failed fibers, *and* lone fibers, *and* which are bracketed at both ends by lone fibers. Some of such configurations are {101, 1001, 10001, 1010, $\cdots$}. From {A} construct {B}, the set of the configurations one specified end of which must be failed. The failed configuration at the end can be on the left or the right end, but no distinction is made between them. A third set {C} is also constructed out of {A} in which both

ends of a configuration have failed, e.g., $\{010, 0100, \cdots\}$. Finally, suppose that $\{P\}$ is the set of configurations with no failed bond, e.g., $\{1, 11, 111, \cdots\}$. One then defines generating functions

$$A(z) = \sum_{L=3}^{\infty} A_L z^L, \quad B(z) = \sum_{L=2}^{\infty} B_L z^L, \quad C(z) = \sum_{L=3}^{\infty} C_L z^L, \quad (39)$$

where $A_L$, $B_L$, and $C_L$ are the sums, respectively, of the survival probabilities of the sets $\{A\}$, $\{B\}$, and $\{C\}$ for a fixed $L$. Likewise, a generating function for $\{P\}$ is also defined

$$P(z) = \sum_{L=0}^{\infty} (p_{s0})^L z^L = \frac{1}{1 - p_{s0} z}, \quad (40)$$

where $p_{s0}$ is the probability that a bond with no failed neighbors survives. Leath and Duxbury (1994) showed that the generating function for the survival configurations, $S(z) = \sum_L p_{sL} z^L$, is given by ($p_{sL}$ is the survival probability for a fixed $L$)

$$S(z) = C(z) + \frac{P(z)[1 + B(z)]^2}{1 - P(z)A(z)}. \quad (41)$$

Since $p_{fL} = 1 - p_{sL}$, where $p_{fL}$ is the failure probability for a fixed $L$, then

$$f(z) = \frac{1}{1 - z} - S(z) \quad (42)$$

where $f(z) = \sum_L p_{fL} z^L$, with $p_{f0} = 0$ and $p_{sL} = 1$. We thus obtain

$$(1 - z)[1 + B(z)]^2 - [1 - p_{s0} z - A(z)]\{1 - (1 - z)[f(z) + C(z)]\} = 0. \quad (43)$$

Expanding identity (43) in powers of $z^L$ and setting the coefficient of the $z^L$ term to zero, one finds the following recurrence relation

$$X_L = X_{L-1} + p_{s0}D_{L-1}X - 2D_L B - A_L + p_{f1}A_{L-1} - B_2 B_{L-2}$$
$$+ \sum_{i=1}^{L-4}(A_{i+2}D_{L-i-2}X - B_{i+1}D_{L-i-1}B), \quad (44)$$

in which $X_L = p_{fL} + C_L$, and $D_L Y = Y_L - Y_{L-1}$. Thus one needs $A_L$, $B_L$, and $C_L$ to use recursion relation (44). These are found by defining new subsets $\{a_{L,l}\}$, $\{b_{L,l}\}$, and $\{c_{L,l}\}$, where, e.g., $\{c_{L,l}\}$ is the set of survival configurations of length $L$ which end with exactly $l$ failed bonds. Recursion relations are also found for these new quantities. For example, $a_{L,l} = b_{L-1,l} \, p_{sl}$, and

$$b_{L,l} = p_{fl}p_{sl}\delta_{L-l-1} + \sum_{i=1}^{L-l-2} b_{L-l-1,i} p_{s,L+i} p_{fl}. \quad (45)$$

These recursion relations can then be used efficiently for calculating various quantities of interest. Because of their efficiency, the behavior of the system for large $L$, of the order of several thousands, can be studied.
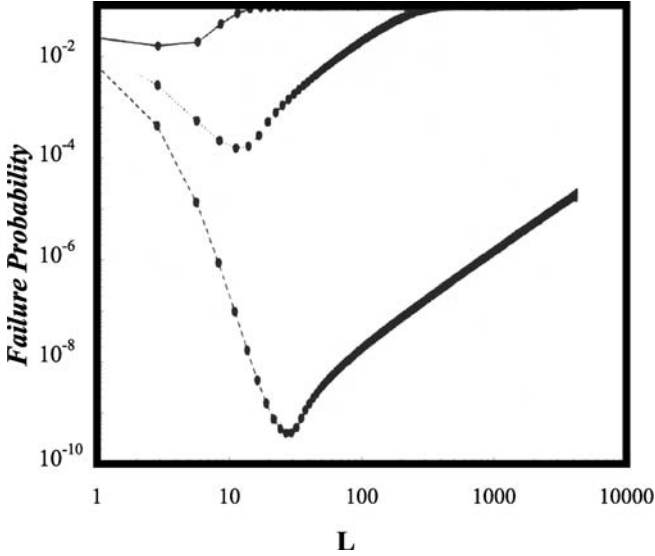
FIGURE 8.2. Dependence of the failure probability of the chain-of-bundles model on the linear size $L$ of the system (after Duxbury and Leath, 1994b).

An interesting and unexpected result of these calculations is that, the failure probability possesses a *deep minimum* with respect to $L$. Figure 8.2 presents a sample of the results (Duxbury and Leath, 1994b). For a large applied stress, the failure probability increases monotonically with $L$. However, if the applied stress is small, then the failure probability possesses a deep minimum at an optimal size $L_o$, hence pointing to the intriguing possibility of designing fibrous materials that operate near their minimum failure probability.

A similar, but simpler, exact recursive method was developed by Wu and Leath (1999). They considered a bundle of parallel fibers in which the local fiber strengths were distributed according to a statistical distribution $f(\sigma)$. Periodic boundary conditions were imposed on the system. Their analysis indicated that there is a critical size $n_c$ (measured in units of the number of fibers) at which there is a transition from a *tough* material to a brittle-like one. More specifically, one has one of the three following scenarios.

(1) If the size $n$ of the system is less than $n_c$, then the material is in the tough region which is characterized by very small stresses and small system sizes. The probability of failure of the material is a superposition of a very large number of local distributions $f(\sigma)$. Since the failure of the material is path dependent, the number of such local distributions can be as large as $2^{n-1}(n!)$. In this case, if the statistical distribution of the local strengths is given by a Weibull distribution, Eq. (18), then the cumulative failure probability $F_n(\sigma) = 1 - P_n(\sigma)$ is given by

$$F_n(\sigma) = 1 - \exp\left[-(n!)^{\gamma(m)m}\sigma^{mn}\right], \tag{46}$$

where $0 < \gamma(m) < 1$ is a parameter that depends on $m$. Equation (46) has the general form of a Weibull distribution. Thus, the optimal sample size $n_{\min}$ that corresponds to the minimum failure probability is obtained from $F_{n-1} = F_n$, yielding

$$n_{\min} \sim \sigma^{-1/\gamma}. \tag{47}$$

(2) If $n \gg n_c \gg 1$, then the material is in the brittle region, where it is macroscopically brittle but microscopically tough. Roughly speaking, the failure of the material depends on whether the size of the weakest region exceeds $n_c$. Since, as discussed in Chapter 5, the probability of finding a weak region of size larger than $n_c$ decays exponentially (because in this case the statistics of the weak or failed regions is described well by percolation statistics), the cumulative failure probability is of the Gumbel type:

$$F_n(\sigma) = 1 - \exp\left[-an\exp(b\ln\sigma/\sigma^m)\right], \tag{48}$$

where $a$, $b$ and $m$ are fitting parameters. The size dependence of the mean failure stress $\langle\sigma_f\rangle$ can then be obtained by neglecting the slow-varying factor $\ln\sigma$ and taking the median as the average, which then yield

$$\langle\sigma_f\rangle \sim (\ln n)^{-1/m}. \tag{49}$$

(3) For $n_c \sim O(1)$ the material is in the *super-brittle* regime. This situation arises when the applied stress is so large that the critical nuclei exist almost everywhere, and thus almost all the fibers fail simultaneously. The cumulative failure probability is then simply

$$F_n(\sigma) = 1 - [1 - f(\sigma)]^n, \tag{50}$$

where $f(\sigma)$ is the local strength distribution.

For related work on this problem see Wu and Leath (2000) and Kun *et al.* (2000).

## 8.1.3  *Computer Simulation*

Simulation of more realistic models of fiber networks (with interconnected fibers) have also been undertaken by, for example, Åström *et al.* (1994, 2000) who used a realistic model in which the fibers were linearly elastic beams, described in Section 8.13 of Volume I, up to a threshold to be defined below, so that the fibrous material can be considered as being brittle. Consider, as an example, a 2D system of such fibers, each of which has a length $l_f$. The network is constructed within a rectangular surface of size $L_x \times L_y$. The $(x, y)$ coordinates of the fibers' centers are selected from uniform distributions in the intervals $[-l_f, L_x + l_f]$ and $[0, L_y]$, respectively, while their orientations are chosen from a uniform distribution in the interval $[-\pi/2, \pi/2]$. The cross sections of the beams are assumed to be squares of width $w$ with $w \ll l_f$. The beams can be stretched and bent and are made of a material with a Young's modulus $Y_f$. Two crossing fibers are rigidly bounded together at their intersection, meaning that all the elastic energies are stored in the
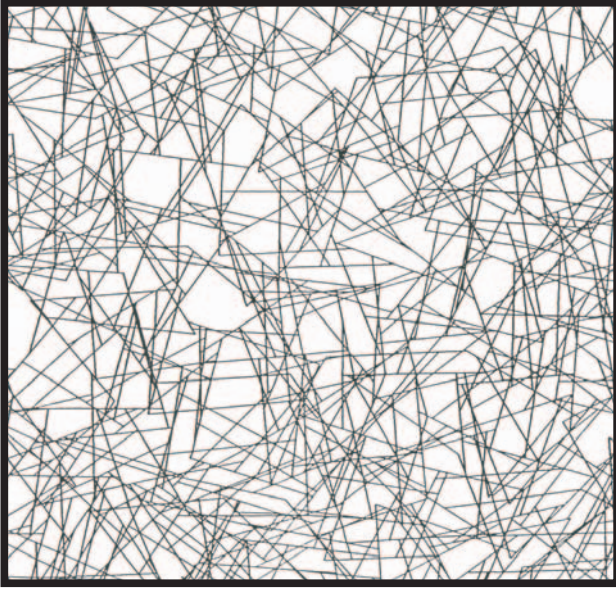
FIGURE 8.3. A typical realization of a 2D model of a fibrous material with randomly distributed and intersecting fibers.

beams and not at the intersections, and that when the network is deformed, the angle between the crossing fibers will remain constant. Each fiber-fiber bond has three degrees of freedom: Horizontal and vertical displacements, and rotations. An example of a typical realization of such a model is shown in Figure 8.3. Two distinct cases can be considered. (1) The beams are embedded by a background material with specific elastic properties, as in, for example, a sheet of paper. (2) Alternatively, the system consists of a network of the beams alone, as in, for example, a polymer network.

The elastic properties of the model depend on the aspect ratio $w/l_f$, as well as the density $p$ of the fibers, defined as the average total length of fibers in an area of $l_f^2$. The percolation threshold, or the critical density of the fibers, is given by

$$p_c \simeq 5.71 l_f. \tag{51}$$

Each fiber contains a segment of length $l_s$ which is that part of the fiber that is between the two intersections that the fiber has with two other fibers. Clearly, the length of the segments is a random variable, as the fibers are distributed randomly in the system. The average segment length is given by

$$\frac{\langle l_s \rangle}{l_f} = \frac{\pi}{11.42(p/p_c)} \simeq \frac{p_c}{3.6p}. \tag{52}$$

The elastic interaction between two connected bonds is characterized by a stiffness matrix $\mathbf{C}$. If the moment of inertia of the cross section is $\mathcal{M} = w^4/12$, then

the stiffness matrix for $w \ll l_s$ is given by

$$
\mathbf{C} = \begin{pmatrix}
Y_f w^2/l_s & 0 & 0 & -Y_f w^2/l_s & 0 & 0 \\
0 & 12Y_f\mathcal{M}/l_s^3 & 6Y_f\mathcal{M}/l_s^2 & 0 & -12Y_f\mathcal{M}/l_s^3 & 6Y_f\mathcal{M}/l_s^2 \\
0 & 6Y_f\mathcal{M}/l_s^2 & 4Y_f\mathcal{M}/l_s & 0 & -6Y_f\mathcal{M}/l_s^2 & 2Y_f\mathcal{M}/l_s \\
-Y_f w^2/l_s & 0 & 0 & Y_f w^2/l_s & 0 & 0 \\
0 & -12Y_f\mathcal{M}/l_s^3 & -6Y_f\mathcal{M}/l_s^2 & 0 & 12Y_f\mathcal{M}/l_s^3 & -6Y_f\mathcal{M}/l_s^2 \\
0 & 6Y_f\mathcal{M}/l_s^2 & 2Y_f\mathcal{M}/l_s & 0 & -6Y_f\mathcal{M}/l_s^2 & 4Y_f\mathcal{M}/l_s
\end{pmatrix}
$$

$$(53)$$

The forces acting on the bonds at the segment ends are obtained by multiplying $\mathbf{C}$ by the vector $(u_{x1}, u_{y1}, \varphi_1, u_{x2}, u_{y2}, \varphi_2)$, where $\mathbf{u} = (u_x, u_y)$ is the displacement vector and $\boldsymbol{\varphi} = (\varphi_1, \varphi_2)$ is the rotation vector. If $l_s$ is short, the bending stiffness $12Y_f\mathcal{M}/l_s^3 = Y_f w^4/l_s^3$ should, as a first approximation, be replaced by the shear modulus $Y_f w^2/[2(1 + v_p)l_s]$, where $v_p$ is the Poisson's ratio of the material.

The fiber network is deformed by, for example, stretching it uniformly in the $x$-direction, which means, for example, fixing the edge at $x = 0$ and pulling in the positive $x$-direction the edge which is initially at $x = L_x$, with the fibers crossing these edges rigidly tied to them. Periodic boundary condition is used in the $y$-direction. Computations of the system deformation, when there is a background matrix, is not straightforward. Typically, a finite-element method, of the type described in Section 7.11.3, is used. Several commercial computer programs that are capable of performing such computations are available. If the system consists only of the fiber network (with no background material), then the computations proceed in the same manner that was described for elastic percolation networks (see Chapter 8 of Volume I). If the fiber density is too low, the system is not rigid and the elastic stiffness is zero.

To study brittle fracture of the material, a failure criterion must be defined. Although such criteria will be described in the next section where we discuss more general discrete models of brittle fracture, we mention a few of them here. One can, for example, consider a fiber as broken or failed if the axial tension or bending of its corresponding beam exceeds a pre-set threshold. Alternatively, failure of a fiber can be defined based on the shear-lag strain, defined as the magnitude of the jump in the axial strain on a fiber across a bond. A combination of all such criteria can also be considered, and in fact Åström *et al.* (1994) studied the case in which fracture occurred by segment breaking due to axial tension and failure at a critical value of shear-lag strain. Once the failure criterion is set, the fracture simulations begin. Each time a fiber fails, the stress and strain distributions in the network must be recomputed, as the network's configuration changes dynamically. As such, the computations are very intensive. Computer simulations indicated that, the failure of the system at the initial stages occurs more or less randomly, and thus the fracture process is similar to percolation. Figure 8.4 shows the stress-strain diagram of the system when relatively few fibers have failed. As expected, up to a certain strain, the stress-strain relation is perfectly linear which is what is expected
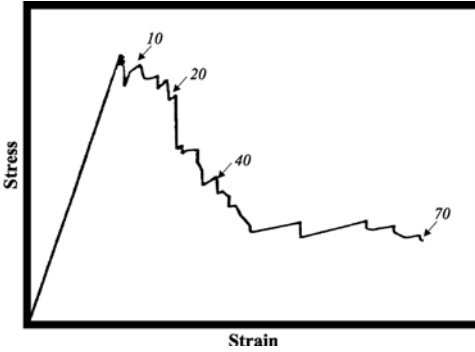
FIGURE 8.4. Stress-strain curve of the fiber network shown in Figure 8.3. The numbers refer to the iterations (after Åström *et al.*, 1994).

of brittle materials. Beyond the critical strain, however, the stress shows a generally downward trend with increasing strains, accompanied by fluctuations that are the result of having fibers failing at essentially random locations in the system.

However, as the number of the microcracks increases, the facture zone becomes quasi-1D, populated mainly by such microcracks with no dominating fracture that can propagate. The absence of a dominating fracture is presumably because of the random orientations of the fibers that help distribute the applied stress in the network more evenly than in a regular network where a dominating fracture usually forms (see below). This behavior is also different from what is usually observed during fracture of composite, but non-fibrous, materials described in Chapters 6 and 7.

### 8.1.4   Mean-Field and Effective-Medium Approximations

One may develop a mean-field or an effective-medium approximation for estimating the elastic *and* fracture properties of such fiber networks. The oldest of such approximations for fiber networks appears to have been developed by Cox (1952). The development of this type of approximation parallels those previously discussed for linear conductivity and elastic moduli of materials described in Chapters 4–8 of Volume I, which we now describe.

Consider first the simplest possible approximation, which we refer to as the EMA1. Suppose that a fiber is attached to an effective medium—a uniform sheet—under tensile strain and is stretched via a number of links; see Figure 8.5. If the fiber is uniformly stretched along with the sheet, the stress $\sigma_f$ along it would be constant. However, this is not possible as $\sigma_f$ must vanish at the fiber's ends. If the strain is small, we can write down the following equation for $\sigma_f$:

$$\langle l_s \rangle \frac{d\sigma_f}{dx} = c \frac{u_f - u_s}{\langle l_s \rangle}, \tag{54}$$

where $u_f$ and $u_s$ are the local displacements of the fiber and the sheet, respectively, and $c$ is a parameter that depends on $w/l_f$ and $p/p_c$. Since $\sigma = Y_f \epsilon$, where $\epsilon$ is
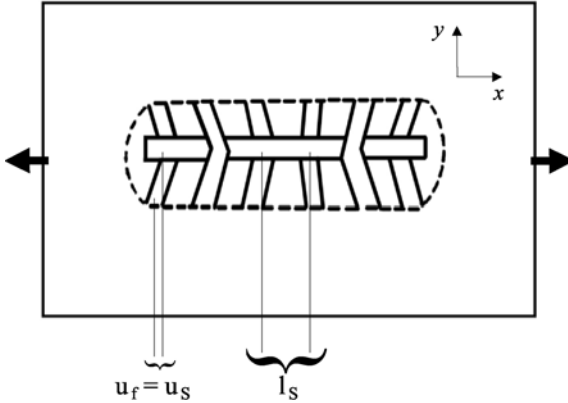
FIGURE 8.5. Schematics of EMA computation of stress along a fiber (after Åström *et al.*, 1994).

the strain, it can easily be verified that the following equation due to Cox (1952),

$$\sigma_f(x, k) = Y_f \epsilon_x \left\{ 1 - \frac{\cosh[k(\frac{1}{2} - x/l_f)]}{\cosh(\frac{1}{2}k)} \right\}, \tag{55}$$

which he derived by a mean-field approximation, satisfies Eq. (54) and the boundary conditions that the stress vanishes at $x = 0$ and $x = l_f$; here, $k = \sqrt{c} l_f / \langle l_s \rangle$. The strain $\epsilon_x$ that appears in Eq. (55) is in fact the strain $\epsilon_f$ in the fiber, but use of $\epsilon_x$ indicates that the strain lies in the $x$-direction. Note that the average shear-lag stress is simply $\langle l_s \rangle d\sigma_f / dx \sim k \langle l_s \rangle / l_f$.

So far, we have assumed that the single fiber embedded in the effective medium is aligned with the direction of the external strain along which the sheet is stretched (see Figure 8.5). In general, however, the fibers are distributed randomly, and therefore one must obtain the orientation dependence of $\sigma_f$. This is, however, straightforward since, in the absence of transverse Poisson contraction, a rotation by an infinitesimal field $\sigma_f$ yields

$$\sigma_f(x, k, \theta) = \sigma_f(x, k) \cos^2 \theta, \tag{56}$$

with $\theta$ being the angle of the fiber with respect to the direction of the external strain. Figure 8.6 compares the predictions of Eqs. (55) and (56) with the simulation results (Åström *et al.*, 1994) in which $k$ has been treated as an adjustable parameter. It is clear that the predictions agree well with the simulation results. These simulations also indicate that $k = p(1 + aw/l_f)/p_c$, where $a$ is a constant.

However, the foregoing treatment is not without problems, especially if it is further developed in order to predict the elastic stiffness of the network, because it actually makes the segment stresses correlated along the fibers with reduced stress close to the fiber ends. A refined treatment of the problem, which we refer to it as the EMA2, can be developed (Åström *et al.*, 2000) if one combines the
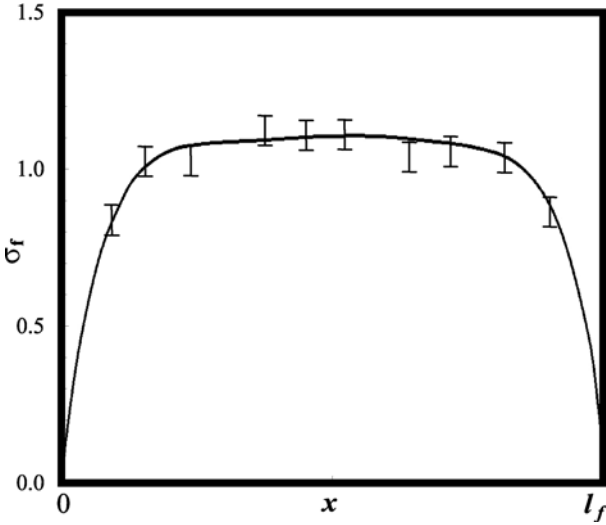
FIGURE 8.6. Average distribution of axial stress along fiber for $p/p_c = 4$ and $l_f/w = 18.8$ (after Åström *et al.*, 1994).

probability distribution for the segments' length with the assumption that the fiber segments deform only in the energetically most-favorable mode, with the modes being bending, stretching, and shearing. Since the center of the fibers are distributed at random in the simulation cell, the probability distribution for their segment length is known, and is given by

$$P(l_s) = \frac{2p}{\pi l_f} \exp\left(-\frac{2p}{\pi l_f} l_s\right), \tag{57}$$

and therefore the average segment length is, $\langle l_s \rangle = \pi l_f/(2p)$.

If the deformation of the fiber network is quasi-static, then the fiber segments will be deformed such that there is force equilibrium at all fiber-fiber bonds, which also define the global minimum of the total elastic energy of the system. This implies that the fiber segments will, in general, be deformed in a way that offers the least elastic resistance. We may define the segments either by bending/shearing or by stretching. According to the stiffness matrix (53), the bending stiffness modulus is $Y_f w^4/l_s^3$, the shear stiffness modulus is $Y_f w^2/[2(1 + \nu_p)l_s]$, while the elongation stiffness modulus is $Y_f w^2/l_s$. Åström *et al.* (2000) assumed that a segment deforms only by bending if the bending modulus is smaller than both the shear and elongation modulus, i.e., if the segment length is such that, $l_s > l_c \equiv w\sqrt{2(1 + \nu_p)}$. On the other hand, if $l_s < l_c$, then the segments are assumed to deform by shearing and stretching. The final ingredient of the model is the assumption that elongation of a segment is proportional to $\cos^2\theta$ [similar to Eq. (56)], while bending and shear are proportional to $\sin 2\theta$. We note that the strain field in the effective-medium treatment does not include any rotation.

We can now compute the total elastic energy $\mathcal{H}$ of the system which is given by

$$\mathcal{H} = p\epsilon_x^2 \left(\frac{L_x L_y}{l_f}\right)\left(\frac{1}{2}Y_f w^2\right) \int_{-\pi/2}^{\pi/2} \frac{\cos^4\theta}{\pi} d\theta \int_0^{l_c} \frac{2p}{\pi l_f} \exp[-2pl/(\pi l_f)]dl + p\epsilon_x^2 \left(\frac{L_x L_y}{l_f}\right)$$

$$\times \left(\frac{1}{2}Gw^2\right) \int_{-\pi/2}^{\pi/2} \frac{\sin^2(2\theta)}{4\pi} d\theta \int_0^{l_c} \frac{2p}{\pi l_f} \exp[-2pl/(\pi l_f)]dl + p\epsilon_x^2 \left(\frac{L_x L_y}{l_f}\right)\left(\frac{1}{2}Y_f w^4\right)$$

$$\times \int_{-\pi/2}^{\pi/2} \frac{\sin^2(2\theta)}{4\pi} d\theta \int_{l_c}^{\infty} \frac{2p}{\pi l_f l^2} \exp[-2pl/(\pi l_f)]dl, \qquad (58)$$

where $G = Y_f/[2(1 + \nu_p)]$. On the other hand, the elastic energy is related to the effective stiffness $C_e$ of the network by, $\mathcal{H} = (1/2)C_e\epsilon_x^2 L_x L_y$, which means that the expression for $C_e$ is given by

$$C_e = \frac{pw^2 Y_f}{8l_f}\left\{\left(\frac{2pw}{\pi l_f}\right)^2 \left[\frac{e^{-z}}{z} - \mathcal{E}_1(z)\right] + \left[3 + \frac{1}{2(1+\nu_p)}\right](1 - e^{-z})\right\}, \qquad (59)$$

where $z \equiv 2pl_c/(\pi l_f)$, and

$$\mathcal{E}_n(z) = \int_1^{\infty} \frac{e^{-zx}}{x^n}dx. \qquad (60)$$

The first test of Eq. (59) is its ability for reproducing the known results in certain limits. Hence, consider first the limit $w \to 0$. If we rescale the network stiffness, $C_e \to C_e/w^2$ when $w \to 0$, the fiber network becomes a central-force network, i.e., a network of simple Hookean springs. Equation (59) then yields $C_e \propto w \to 0$, which is expected since the average coordination number of the network is less than 4, and therefore, as explained in Section 8.7.3 of Volume I, the network cannot be rigid. On the other hand, in the limit $p \to \infty$, which is equivalent to $w/\langle l_s\rangle \to \infty$, Eq. (59) predicts that $C_e \propto Y_f w^2 p/l_f$, implying that in the limit of high $p$ the network stiffness is simply proportional to $Y_f$ multiplied by the density of the fibers in the network, i.e., the network behaves as an elastic continuum, which is the expected behavior.

However, there remains one problem to be addressed. In writing down the expression for the total elastic energy $\mathcal{H}$, Eq. (58), it was assumed that all segments are deformed. However, below the percolation threshold of the network, $C_e = 0$, and no segment is deformed. At, and just above, $p_c$, there are also many segments that carry no load, while for $p \gg p_c$ such segments appear only at the end of the fibers with a density that can be shown to be about $0.55p_c$, independent of $p$ (Åström et al., 1994). Thus, for Eq. (59) to reproduce the correct percolation behavior, one must make a transformation from $p$ to the density $p_l$ of the loaded segments, and Åström et al. (2000) suggested that $p/p_c = p_l/p_c + 0.55 + 0.45/(p_l/p_c + 1)$, which is simply a crossover from $p = p_c$ when $p_l = 0$ to $p_l \to p - 0.55p_c$ in the limits $p \to \infty$ and $p_l \to \infty$. Therefore, one should replace the first $p$ on the
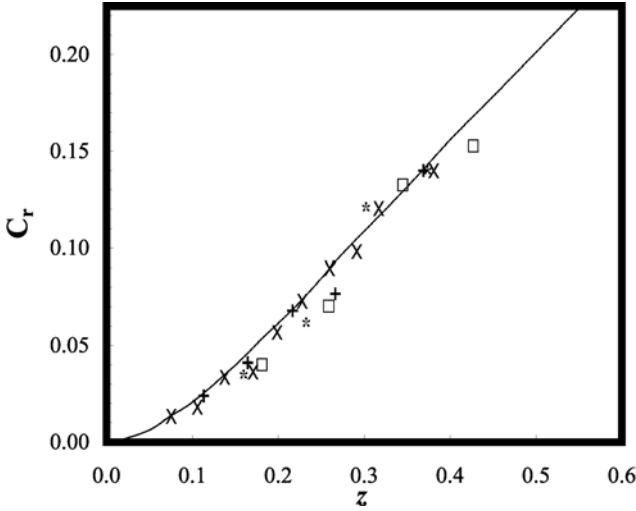
FIGURE 8.7. Comparison of the predictions of Eq. (62) (curve) with the results of numerical simulations for $w = 0.05$ (+), $w = 0.06$ (×), $w = 0.07$ (∗), and $w = 0.08$ (□) (after Åström *et al.*, 2000).

right-hand side of Eq. (58) by $p_l$ given by

$$p_l = \frac{p_c}{2} \left\{ \frac{p}{p_c} - 1.55 + \left[ \left( 1.55 - \frac{p}{p_c} \right)^2 - 4 \left( 1 - \frac{p}{p_c} \right) \right]^{1/2} \right\}. \quad (61)$$

Finally, if we define $z_l = 2 p_l l_c / (\pi l_f)$, and a reduced stiffness $C_r = 16 \sqrt{2(1 + \nu_p)}\, C_e / (\pi w Y_f)$, we obtain

$$C_r = z_l \left\{ \frac{z^2}{2(1 + \nu_p)} \left[ \frac{e^{-z}}{z} - \mathcal{E}_1(z) \right] + \left[ 3 + \frac{1}{2(1 + \nu_p)} \right] (1 - e^{-z}) \right\}. \quad (62)$$

Figure 8.7 compares the predictions of Eq. (62) with the simulation results for various values of $w$, and it is clear that the agreement between the two sets is quite good.

The shape of the stress-strain diagram for the fractured fiber network, shown in Figure 8.4, can also be understood by appealing to the EMA. Here, we discuss how this is accomplished by using the EMA1. The shape of a stress-strain diagram of a fracturing material depends critically on the failure criterion. Suppose, for example, that breaking occurs by axial tension. Then, Eq. (56) predicts that the critical angle $\theta_f$ for failure is given by

$$\theta_f = \arccos \left( \sqrt{\frac{\epsilon_f}{\epsilon_x}} \right), \quad (63)$$

where $\epsilon_f$ is the axial strain for failure. Equation (63) is obtained by writing $\sigma_f(x, k, \theta_f) = Y_f \epsilon_f$ and $\sigma_f(x, k) = Y_f \epsilon_x$ and solving the resulting equation for

$\theta_f$. It can then be shown, using the EMA1 treatment described above, that one obtains the following expression for the stress $\sigma$ as a function of the strain $\epsilon_x$ (Åström *et al.*, 1994):

$$\sigma = \epsilon_x w^2 \left(\frac{2}{\pi}p\right) Y_f \int_{\theta_f(\epsilon_x)}^{\pi/2} [\cos^4\theta + (G/Y_f)\sin^2(2\theta)]d\theta, \qquad (64)$$

which is obtained from the total elastic energy of the system which, within the EMA1, is given by

$$\mathcal{H} = \frac{1}{2}\epsilon_x^2 w^2 Y_f \left(\frac{2}{\pi}p\right) L_x L_y \int_0^{\pi/2} [\cos^4\theta + (G/Y_f)\sin^2(2\theta)]d\theta. \qquad (65)$$

Equation (65) is of course a simplified version of Eq. (58).

As discussed above, as more fibers fail, the fracture zone becomes a narrow, quasi-1D zone. Thus, in order to create such a zone, one assigns an infinitesimally lower failure threshold to a band across the network, and then applies Eq. (64) in this fracture zone which is given a unit width. No fiber fails in the rest of the network, i.e., Eq. (64) is applied with $\theta_f = 0$. The result is shown in Figure 8.8 in which the dashed curve is the equilibrium curve. These predictions are in qualitative agreement with the simulation results shown in Figure 8.4, except that there is a discontinuity in the predicted stress-strain diagram after the elastic regime (the regime of a linear relation between $\sigma$ and $\epsilon_x$) ends, whereas the simulations do not indicate such a sharp and discontinuous change. The disagreement between the simulation results and the EMA1 predictions becomes progressively stronger as more fibers fail. The same qualitative trends would have been obtained, had we used the EMA2 to derive the stress-strain diagram for the fracturing fiber network. Therefore, although the EMA provides some qualitative insight into the
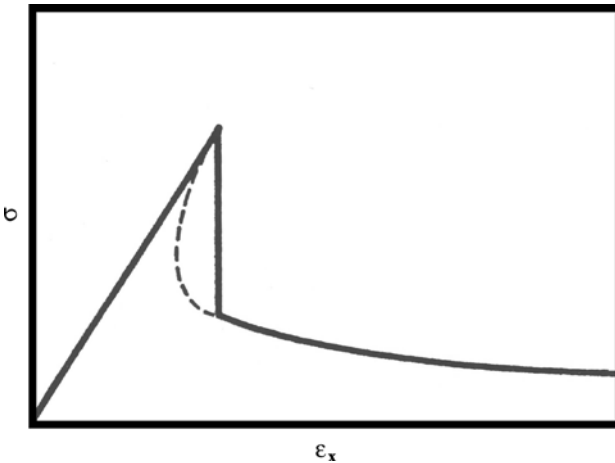


FIGURE 8.8. The stress-strain diagram as predicted by the effective-medium approximation (solid curve).

early stages of a fracture process, it cannot be expected to be accurate as failure of the fibers progresses.

## 8.2    Quasi-static Fracture of Heterogeneous Materials

Lattice models can represent the behavior of fracture of materials if the phenomenological coefficients and properties of the materials, such as their elastic constants and failure threshold (see below), are properly defined and set. As discussed above and also in Chapter 7, use of a finite-element (FE) method for discretizing the continuum equations and studying fracture has been popular among engineers. The discretized equations, and the associated mesh that one obtains in such approaches, resemble a lattice model. However, only very weak spatial disorder can be incorporated into such a model, since strong disorder necessitates use of a very fine structured FE mesh which makes the computations prohibitive. An alternative approach to the FE method is based on identifying the key microstructural features associated with the disorder and relevant to the failure process. One then subsumes all of the details of the mechanical behavior of that material region, including the failure of a region of the material by the nucleation of a stable crack of the same size, into a *local* constitutive law. Disorder is included by allowing the phenomenological coefficients of the constitutive law to vary, from bond to bond, according to some probability distribution. A network of such bonds is then used to numerically calculate local stresses on, and interactions between, the bonds (and sites) under the application of a macroscopic boundary condition. By allowing for failure of such bonds under their local stress or strain (or a combination of both), cracks are formed which may interact with each other, generate new cracks via load transfer, and propagate to macroscopic sizes, leading to material failure. Thus, one is able to account for the nucleation of cracks on the key length scales and also the effect of disorder on such phenomena.

This approach was first used, in a rather primitive form, over 30 years ago. Early efforts for developing discrete models of fracture of materials (Mikitishin *et al.*, 1969; Dobrodumov and El'yashevich, 1973) used lattices in which the bonds were linear springs that could only be stretched (no bending or rotation was allowed). However, because of the computational limitations of their times, and the over-simplified nature of the models, they did not attract wide attention. To our knowledge, modern lattice models of quasi-static mechanical fracture of heterogeneous materials, of the type that are described in this chapter, were first proposed by Sahimi and Goddard (1986).

Generally speaking, three variations of such lattice models have been developed for studying mechanical breakdown in disordered materials. In the first approach, which is completely deterministic, one uses a heterogeneous lattice each bond of which describes the system on a certain length scale, with failure characteristics described by a few key parameters. One then deforms the lattice gradually by applying a boundary condition to the system that resembles what is used in an experiment on fracture of a material, as a result of which the individual bonds

break irreversibly in a certain manner. These models are either quasi-static so that the process time enters the computations only as the number of Monte Carlo steps (or as the number of bonds that are broken), or explicit time-dependence of the fracture process is somehow built into them. This class of models is usually appropriate for materials in which the disorder is *quenched* (fixed in time).

The second and third approaches are probabilistic. One of them (Louis and Guinea, 1987; Hinrichsen *et al.*, 1989; Meakin *et al.*, 1989) draws on an analogy between mechanical breakdown and the dielectric breakdown model of Niemeyer *et al.* (1984) described in Section 5.4.1. As in Niemeyer *et al.*'s model, these models give rise to complex fractal crack patterns, and may be appropriate for systems in which disorder is annealed; comparison between the predictions of such models and fracture of materials with annealed disorder confirms this (see below). The second class of probabilistic models was intended mainly for fracture of polymeric materials. In these models, an elastic element breaks with a temperature-dependent probability, hence taking into account the effect of the activation and elastic energies stored in the element. As we will see later in this chapter, many of the probabilistic models have, in some sense, some type of dynamics built into them.

### 8.2.1   Lattice Models

Consider a 2D network, such as a $L \times L$ triangular or square lattice, or a 3D network such as a $L \times L \times L$ simple-cubic or BCC lattice. Every bond of the lattice represents a Hookean spring or beam. In the former case, every site $i$ of the lattice is characterized by a displacement vector $\mathbf{u}_i$, while in the latter case, in addition to its displacement $\mathbf{u}_i$, site $i$ is also characterized by a rotation vector. Hence, the nearest-neighbor sites are connected by springs or beams. The initial (equilibrium) length of all the springs or beams is the same and, unlike the FE method in which the mesh is made finer where the stress is larger, the initial topology of the network is the same everywhere. The exception to this rule is when one uses the lattice models for studying mechanical and fracture properties of fibrous materials. Such models were already described above and also in Chapter 8 of Volume I, and therefore are not discussed any further.

We consider here the case of a brittle material for which a linear relation between the stress and strain in the spring or beam is valid up to a threshold (defined below). A force law of this type is not, of course, completely realistic, but has long been thought of as a sensible approximation for brittle ceramics (see, for example, the discussion by Lawn, 1993). The displacements $\mathbf{u}_i$ (and the rotations) are computed by minimizing the total elastic energy $\mathcal{H}$ of the system, the exact form of which depends on the type of model that one wishes to study, and the degree of microscopic detail that one incorporates into the model. For example, lattices in which only the central- or stretching (Hookean) forces are operative (Sahimi and Goddard, 1986; Beale and Srolovitz, 1988; Fernandez *et al.*, 1988; Srolovitz and Beale, 1988; Hansen *et al.*, 1989; Arbabi and Sahimi, 1990b; Sahimi and Arbabi, 1993), those in which the bond-bending or angle-changing forces, in addition to central forces (see below) also act on the bonds of the lattice (Sahimi

and Goddard, 1986; Arbabi and Sahimi, 1990b; Sahimi and Arbabi, 1992, 1993, 1996; Sahimi *et al.*, 1993), as well as the Born model described in Chapter 8 of Volume I (see also below) in which the elastic energy of the system consists of the contributions by the central forces and a scalar-like term (Hassold and Srolovitz, 1989; Yan *et al.*, 1989; Caldarelli *et al.*, 1994) have all been utilized.

Let us now describe these lattice models. In general, the elastic energy of the bond-bending (BB) model is given by (Kantor and Webman, 1984)

$$\mathcal{H} = \frac{1}{2}\alpha \sum_{\langle ij \rangle}[(\mathbf{u}_i - \mathbf{u}_j) \cdot \mathbf{R}_{ij}]^2 e_{ij} + \frac{1}{2}\gamma \sum_{\langle jik \rangle}(\delta\boldsymbol{\theta}_{jik})^2 e_{ij}e_{ik}, \qquad (66)$$

where $\alpha$ and $\gamma$ are the central and BB force constants, respectively, $\langle jik \rangle$ indicates that the sum is over all triplets in which the bonds $j$-$i$ and $i$-$k$ form an angle with its vertex at $i$, and $e_{ij} = 1$ if $i$ and $j$ are connected, and $e_{ij} = 0$ otherwise. The first term on the right side of Eq. (66) represents the contribution of the stretching forces, while the second term is due to BB forces. The precise form of $\delta\boldsymbol{\theta}_{jik}$ depends on the microscopic details of the model. In the most general form, if bending of all pairs of bonds that have one site in common, including the collinear bonds, is allowed, then (Arbabi and Sahimi, 1990a)

$$\delta\boldsymbol{\theta}_{jik} = \begin{cases} (\mathbf{u}_{ij} \times \mathbf{R}_{ij} - \mathbf{u}_{ik} \times \mathbf{R}_{ik}) \cdot (\mathbf{R}_{ij} \times \mathbf{R}_{ik})/|\mathbf{R}_{ij} \times \mathbf{R}_{ik}|, & \mathbf{R}_{ij} \text{ not parallel to } \mathbf{R}_{ik}, \\ |(\mathbf{u}_{ij} + \mathbf{u}_{ik}) \times \mathbf{R}_{ij}|, & \mathbf{R}_{ij} \text{ parallel to } \mathbf{R}_{ik}, \end{cases} \qquad (67)$$

where, $\mathbf{u}_{ij} = \mathbf{u}_i - \mathbf{u}_j$. For *all* 2D systems, Eq. (67) is simplified to

$$\delta\boldsymbol{\theta}_{jik} = (\mathbf{u}_i - \mathbf{u}_j) \times \mathbf{R}_{ij} - (\mathbf{u}_i - \mathbf{u}_k) \times \mathbf{R}_{ik}.$$

The BB model has a well-defined continuum counterpart. For most materials to which the BB model is applicable, one has $\gamma/\alpha \leq 0.3$ (Martins and Zunger, 1984).

In the Born model the associated elastic energy is given by

$$\mathcal{H} = \frac{1}{2}\alpha_1 \sum_{ij} \mu[(\mathbf{u}_i - \mathbf{u}_j) \cdot \mathbf{R}_{ij}]^2 e_{ij} + \frac{1}{2}\alpha_2 \sum_{ij}(\mathbf{u}_i - \mathbf{u}_j)^2 e_{ij}, \qquad (68)$$

where $\mathbf{R}_{ij}$ is the unit vector along the line from $i$ to $j$, and $\alpha_1$ and $\alpha_2$ represent, more or less, two adjustable parameters. The first term of Eq. (68) is the energy of a network of central-force springs, i.e., Hookean springs that transmit force only in the $\mathbf{R}_{ij}$ direction, but do not transmit shear forces, whereas the second term is a contribution analogous to scalar transport (for example, the power dissipated in conduction), since $(\mathbf{u}_i - \mathbf{u}_j)^2$ represents the *magnitude* of the displacement difference $\mathbf{u}_i - \mathbf{u}_j$. The Born model can be derived from linear continuum mechanics by discretizing the linear equation that governs the elastic equilibrium of a solid, i.e., $\nabla \cdot \boldsymbol{\sigma} = \mathbf{0}$ (where $\boldsymbol{\sigma}$ is the stress tensor), and using the usual relation, $\boldsymbol{\sigma} = \lambda(\nabla \cdot \mathbf{u})\mathbf{U} + \mu[\nabla\mathbf{u} + (\nabla\mathbf{u})^\mathrm{T}]$, where $\lambda$ and $\mu$ are the usual Lamé constants, and $\mathbf{U}$ is the identity tensor (see Section 8.4 of Volume I for details). If this is done, then one obtains, $\alpha_1 = 2(1 - \nu_p)/(1 + \nu_p)$, and

$\alpha_2 = 2(1 - 3\nu_p)/[4(1 - \nu_p)]$, where $\nu_p$ is the Poisson's ratio. However, in this form, the elastic energy given by Eq. (68) will not be rotationally invariant, thus violating a fundamental physical requirement for an elastic energy representation of a solid material. Therefore, Eq. (68), in which $\alpha_1$ and $\alpha_2$ are treated as adjustable parameters, is a semi-empirical representation of materials.

The Born model may be considered as an analogue of a 3D solid in plane-stress with holes normal to the $x$-$y$ plane, or as a 2D solid with the Poisson's ratio defined as the negative of ratio of the strain in the $y$-direction to that in the $x$-direction, when a stress is applied in the $x$-direction but none is applied in the $y$-direction. Results for a 3D solid in plane-strain can be generated from those of this model using the transformation $\nu_p' = \nu_p/(1 + \nu_p)$, where $\nu_p'$ is the Poisson's ratio for the plain strain.

Let us mention another interesting way of generating a BB model. In their studies of brittle fracture, Chung *et al.* (2001) generated a spring network by molecular dynamics simulation, starting with a random distribution of spheres that interact with each other through certain potentials. The system would then be allowed to reach equilibrium, after which the centers of the spheres that were not separated by a distance larger than a certain limit were connected by springs. Both the central and BB forces were included in the network so obtained.

The spring lattices are suitable models for simulating a fracture process in materials that are under shear or tension. However, one should use the beam model (see Chapter 8 of Volume I for more details) (Herrmann *et al.*, 1989a; de Arcangelis *et al.*, 1989; Tzschichholz, 1992,1995; Tzschichholz *et al.*, 1994; Tzschichholz and Herrmann, 1996) when external compressional forces are imposed on the system, since a spring cannot break under compression. In the beam model, in addition to the central and BB or angle-changing forces, torsional forces also contribute to the elastic energy $\mathcal{H}$ of the lattice. We believe, however, that, except when external compressional forces are imposed on the system, the BB model is a completely realistic representation of the elastic energy of disordered materials. Recall that, as discussed in Chapters 8 and 9 of Volume I, the BB model is capable of describing the elastic properties of polymers, glasses, ceramics and powders, and hence use of more complex models for the elastic energy of the material is not necessary. In addition to the above models, a model based on discretization of the following equation (sometimes referred to as the Lamé equation)

$$(\lambda + \mu)\nabla(\nabla \mathbf{u}) + \mu \nabla^2 \mathbf{u} = \mathbf{0}, \tag{69}$$

where $\lambda$ and $\mu$ are the usual Lamé constants, has also been used (Herrmann *et al.*, 1989b).

Sahimi and Goddard (1986) suggested that three general classes of disorder may be incorporated into such model, which are as follows.

(1) Deletion or suppression of a fraction of the bonds either at random or in a prescribed fashion, so that the material's heterogeneity is of percolation-type. The suppressed or deleted bonds may, for example, represent the microporosity or some type of defect in the system *before* the fracture process began.

(2) Random or correlated distribution of the elastic constants $e_{ij}$ of the bonds. The idea is that in real heterogeneous materials the shapes and sizes of the elastic zones through which stress transport takes place may be statistically distributed, resulting in a different $e_{ij}$ for each zone, or bond in the lattice model, that follows some type of a statistical distribution. Such a model may be appropriate for a composite material that, for example, consists of several constituents, each of which has its own elastic properties.

(3) Random or correlated distribution of the critical thresholds at which the linear constitutive relation that describes the stress-strain relation in the beam or spring breaks down. For example, in *shear* or *tension* each bond may be characterized by a critical length $l_c$, such that if it is stretched beyond $l_c$, it breaks irreversibly. Such a threshold can be estimated experimentally by evaluating macro tensile strength of the material. Alternatively, each bond can be characterized by a critical force (stress) $F_c$ ($\sigma_c$), such that if it suffers a force (stress) larger than $F_c$ ($\sigma_c$), it breaks irreversibly. Under compression, a beam breaks if it is bent too much. The idea for using this type of disorder is that a solid material made up intrinsically of the same material (the same elastic constant $e_{ij}$ everywhere) may contain regions having different resistances to breakage under an imposed external stress or potential due to, for example, the presence of defects during its manufacturing or formation process. Depending on the intended application, we may use any combination of the three types of disorder. For example, one may model the disordered material with fractal lattices with bonds that have statistically-distributed properties (such as their elastic constant). Because of their fractality, such models have low connectivities and large porosities, and may be relevant to transgranular stress corrosion cracking of ductile metal alloys, such as stainless steel and brass (Sieradzki and Newman, 1985). They may also be relevant to stress and crack propagation in weakly-connected granular media, such as sedimentary rocks. We do not, however, consider them here as they have not been studied extensively.

Another important source of disorder in stressed materials is the so-called *residual stress* variations, which are caused by, among other things, thermal expansion mismatch. The appropriate elastic lattice models with bond mismatch were described in Section 9.7 of Volume I. We will not discuss the effect of this type of disorder on fracture, although they can be analyzed by modification of the models that are described here (see, for example, Curtin and Scher, 1990a,b; Sridhar *et al.*, 1994).

After selecting the lattice and the form of the elastic energy of the system (i.e., the types of forces that are operative in the lattice), we specify the type of the heterogeneity that the material contains. If the disorder is of percolation-type (type-1 heterogeneity described above), then its inclusion in the lattice is straightforward and needs no discussion. For types-2 and 3 heterogeneities described above, their statistical distribution must be specified. A statistical distribution that has been used widely is

$$f(x) = (1 - \zeta)x^{-\zeta}, \tag{70}$$

where $x$ is any property of the lattice that is statistically distributed and represents its heterogeneity, and $0 \leq \zeta < 1$. The advantage of the distribution (70) is that, varying $\zeta$ allows one to generate distributions that are very narrow ($\zeta \to 0$) or very broad ($\zeta \to 1$), and therefore one can study the extent to which such extreme distributions affect failure phenomena. Note that $\zeta = 0$ represents a uniform distribution, while Roux *et al.* (1988) showed that, in the limit $\zeta \to 1$, fracture becomes equivalent to a type of percolation. A great advantage of the lattice models is that, *any* type of statistical distribution $f(x)$ of the heterogeneities can be used. For example, de Arcangelis *et al.* (1989) used, in addition to (67), a Weibull distribution

$$f(x) = m \lambda^{-m} x^{m-1} \exp\left[-(x/\lambda)^m\right], \tag{71}$$

where $2 \leq m \leq 10$ supposedly describes many real materials.

After specifying the lattice type, the form of the elastic energy $\mathcal{H}$, and the type of disorder, the boundary conditions must be specified. One can, for example, use shear, uniaxial tension or compression, uniform dilation (i.e., pulling a lattice equally in all directions), or surface cracking which is used for simulating fracture of a thin film of a material attached to a substrate (for example, thin polymeric coatings, or paints, or even mud). In this case, each site of the lattice is connected by a spring to the substrate which has a lattice constant larger than the original lattice. In this way all the bonds are equally overstretched without having applied any force on a boundary of the lattice, implying that no external boundary is in fact needed, and one can use periodic boundary conditions in all directions.

The simulations can now begin. One must compute the distribution of the nodal displacements (and rotations, if such motions are allowed), from which the forces (and stresses) exerted on all the bonds are computed. The procedure for doing so consists of minimizing the total elastic energy of the system with respect to the displacements of the internal nodes of the lattice (and their rotations, if such motion is allowed). Because of the assumption of brittleness, these equations are linear and therefore, subject to the boundary conditions imposed on the system, can be solved by one of several methods that are available for solving such equations. If very high precision is needed, then the conjugate-gradient method (see Chapter 9 for a description of this method) is the best technique to use.

After computing the initial distribution of the stresses (and strains) in the lattice, a criterion for nucleation of the microcracks must be specified. The criterion, however, depends on the type of material that is being studied. For example, if each elastic bond is a rubber band, then it will tear apart when stretched beyond a certain limit. Thus, for example, we assign a threshold $l_c$ for the length of the bonds, which is selected from the probability density functions described above. Then, in terms of $l_c$, the breaking criterion is that a bond breaks if its length in the deformed lattice exceeds its $l_c$. Alternatively, among all the bonds that have exceeded their $l_c$, the one with the *largest* deviation from its $l_c$ breaks first. The idea is that in a deformed material, the *weakest* point of the system fails first.

However, if the elastic bond represents, for example, a glass rod, then it will break if it is bent too much. One must of course use a lattice of beams for modeling such a material. Therefore, a good strategy would be devising a breaking criterion

that is a combination of both stretching and bending. For example, in the beam model one can use the criterion (de Arcangelis *et al.*, 1989) that the beam with largest value of the following quantity

$$p_b = \left(\frac{F}{F_c}\right)^2 + \max\left\{\frac{|\mathcal{M}_i|}{\mathcal{M}_c}, \frac{|\mathcal{M}_j|}{\mathcal{M}_c}\right\}, \qquad (72)$$

breaks, where $F$ is the longitudinal force acting along the beam and $F_c$ its critical value, and $\mathcal{M}_i$ and $\mathcal{M}_j$ are the moments applied on the two adjacent sites $i$ and $j$ of the beam, with $\mathcal{M}_c$ being the critical threshold of the moment. Both $F_c$ and $\mathcal{M}_c$ are distributed according to some probability density functions and can, in fact, represent the heterogeneity of the material (type-3 disorder discussed above). The first term of (72) describes breaking due to stretching, while the second term is representative of breaking due to bending. One can distribute, for example, $\mathcal{M}_c$ over a broader range than $F_c$ (and vice versa) as a measure of susceptibility of the material to breaking by bending as opposed to stretching.

Other failure criteria of this type can be, and have been, used. For example, a combination of $F_c$ and $l_c$ can also be used for setting up a breaking criterion (Arbabi and Sahimi, 1990b): one breaks that bond for which the ratio $\mathcal{R} = l_m l / l_c$ is *maximum*, where $l$ is the *current* length of the bond in the deformed lattice and $l_m$ is the maximum length that any bond in the lattice has, or break the bond for which $U = F l_c / F_m$ is *minimum*, where $F$ is the total force that the bond suffers, and $F_m$ is the maximum force on any bond in the lattice. Sridhar *et al.* (1994) used the criterion that a bond breaks if its strain energy exceeds a critical value. The flexibility that the lattice models afford one in using almost any type of failure criterion is one important advantage of such models.

Another advantage of such lattice models is that, depending on the intended application, any failure criterion can be used without imposing any undue difficulty on the computations. For example, often portions of a material are damaged but do not break during the fracture process. The damaged area can be modeled by lowering the breaking threshold of the bonds that are in the vicinity of the growing crack. It is clear that the bonds that are damaged in this way are more likely to break in the next step of the simulations than the undamaged, unbroken bonds. One can also model short-lived damage by considering at step $n$ of the simulations the quantity

$$p_b' = p_b(n) + \beta_0 p_b(n-1), \qquad (73)$$

where $p_b(n)$ is the usual quantity that one uses for breaking criterion at step $n$ *without* considering the damage [for example, Eq. (72)], and $0 < \beta_0 < 1$. Then, the breaking criterion must be based on $p_b'$, and $p_b(n-1)$ indicates the effect of "memory" of the damage that occurred during the previous step of the simulations.

In some fracture processes, such as stress corrosion, the damage accumulates during the entire process, and therefore the breaking criterion must somehow reflect this. A simple algorithm for taking this effect into account is as follows (Herrmann *et al.*, 1989b). One assigns a counter $c(n)$ to each bond of the lattice that is susceptible to damage, e.g., those in the "vicinity" of the growing crack,

where the boundaries of this region must also be specified. For example, one can consider all the unbroken bonds that are up to a certain distance from the cracked area. At the beginning of the simulations, all the counters are set to zero, $c(0) = 0$, for all the bonds that may be damaged. Then, at each step $n$ of the simulations, one calculates $p_b(n)$, the quantity based on which the breaking criterion is applied, and computes

$$\alpha(n) = \frac{1 - c(n-1)}{p_b(n)}. \tag{74}$$

The bond with the smallest value of $\alpha$, namely $\alpha_m$, is then broken. Then, all the counters are updated by

$$c(n) = \alpha_m p_b(n) + \beta c(n-1), \tag{75}$$

so that $c(n)$ "accumulates" the damages from all the previous steps of the process. Clearly, if $\beta = 1$, then the damage is irreversible. The limit $\beta \to 0$ corresponds to criterion (73) with $\beta_0 = 1$.

    The failure process is then initiated by selecting the bond (or bonds) that must break. Two different "dynamics" of fracture propagation can be studied. In model 1 only *one* bond is broken at each stage of the simulation, which is equivalent to assuming that the rate at which the elastic forces relax throughout the network is much faster than the breaking of one bond. In model 2 (and depending on the failure criterion), *all* bonds that meet the failure criterion are broken. Most of the studies so far have used model 1, and the properties of model 2 have not been studied extensively. Breaking a bond is equivalent to removing it from the lattice (by, for example, setting its elastic constant to zero), after which one recalculates the stress and strain distributions for the new configuration of the lattice, select the next bond(s) to break, and so on. If the external stress or strain imposed on the lattice is not large enough to break any bond, it is gradually increased. The simulation continues until a sample-spanning crack is formed. As mentioned above, instead of removing a failed bond from the lattice, one may reduce its elastic constant. In this case, one observes the interesting phenomenon of *crack arrest* (Li and Duxbury, 1988).

### 8.2.1.1    Shape of the Macroscopic Fracture

The number of the cracks, their size distribution, and the shape of the macroscopic fracture, if one is formed, all depend on a variety of factors, including the broadness of the distribution of the material's heterogeneity, the dimensionality of the system, the boundary conditions applied to the system, and the interplay between the quenched disorder, which tries to *delocalize* the propagating microfractures, and stress enhancement at the tip of the microfractures which attempts to *localize* the fracture. For example, if an external strain is applied slowly, a great number of microcracks form before a macroscopic fracture (network) is formed. The number of the nucleated cracks depends on the broadness of the heterogeneity distribution in the system. If this distribution is very broad, then the system is a mixture of

very strong and very weak regions. For the growing crack to take advantage of the weak regions, it must find its path which could be quite tortuous, so much so that it may result in fragmentation of the material (Åström and Timonen, 1997a). We will discuss the phenomenon of fragmentation later in this chapter. On the other hand, when weak disorder is present in a solid, a catastrophic crack is formed quickly that spans the system, and therefore the mechanical failure of the system is very fast.

The number and shape of the cracks also depend on the dimensionality of the system. In 3D stress enhancement at the tip of the growing cracks is much weaker than in 2D, so that even mild disorder can give rise to very complex fracture pattern in 3D. Moreover, the type of the boundary conditions applied to the system is also important. If a stress (instead of a strain) is applied to the material, then the system would fail very quickly soon after the first microcrack is formed, even if the stress is applied slowly to the system. Therefore, similar to real fracture tests, there are significant differences between a stress-controlled and a strain-controlled fracture test in such models. Figure 8.9 shows three different stages of fracture of a 2D triangular lattice. The bonds of the lattice suffer only stretching (central) forces, disorder was generated in the lattice by removing 10% of the bonds at random before deformation of the lattice began, and the deformation was caused by applying a strain to one face of the lattice, while the opposite face was fixed. The top panel shows the initial configuration of the lattice before it was deformed. The middle panel represents the system when it loses its rigidity. Recall from Chapter 8 of Volume I that if only central forces act on the bonds, the lattice loses its rigidity at a fraction of uncut springs which is significantly larger than the connectivity threshold. For example, the connectivity threshold of the triangular lattice is at $p_c = 0.347$, while a triangular lattice of Hookean springs loses its rigidity at $p \simeq 0.641$. However, the rotational freedom of the bonds (which do not contribute to the total elastic energy of the system) is not lost yet. At this point, the lattice contains a strip of bonds with *zero* shear modulus, which can be thought of as a type of shear band. Note that the shear band can move at most one bond before the shear modulus is re-established. Final failure of the lattice, which occurs due to formation of a macroscopic crack that splits the system into two pieces, occurs very close to the shear band. In contrast with Figure 8.9, Figure 8.10 presents an example of the microcracks that are produced in the same lattice in which both the central and bond-bending (angle-changing) forces contribute to the elastic energy of the system, and the distribution of the heterogeneities is broad, following Eq. (70).

Other interesting crack shapes can also be produced by such models. For example, Herrmann *et al.* (1989b) considered a square lattice of beams with periodic conditions in the horizontal direction, while the top and bottom of the lattice was sheared. They defined a probability $p_b$ slightly more general than Eq. (69), given by

$$p_b = [F^2 + q \cdot \max(|\mathcal{M}_i|, |\mathcal{M}_j|)]^\eta, \tag{76}$$

where $F$ is the traction (and/or compression) force applied on a beam, and $q$ and $\eta$ are two free parameters, with $q$ being a measure of the affinity of a beam to

FIGURE 8.9. At the top is the ini-
tial configuration of the system in
which 10% of the bonds have been
removed at random. In the mid-
dle is the system after the rigidity
failure (when the fraction of bro-
ken bonds is less than the rigidity
threshold of the system), while
the bottom system is at the com-
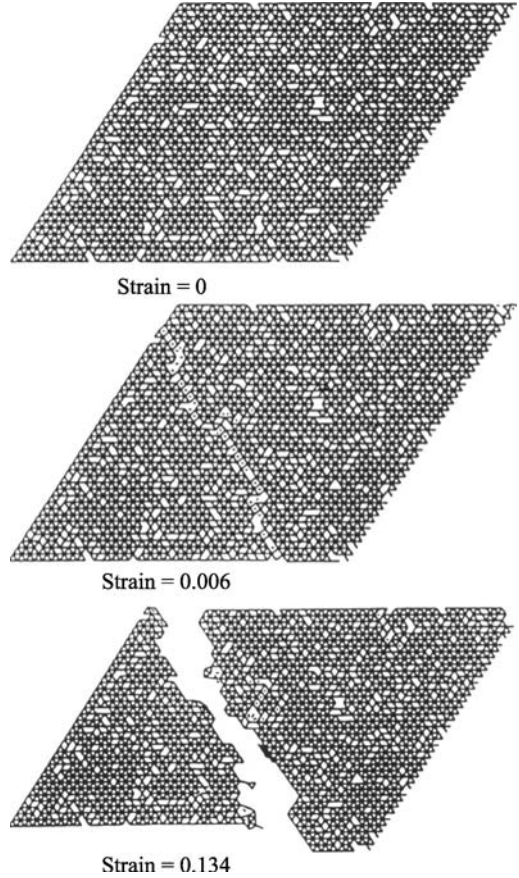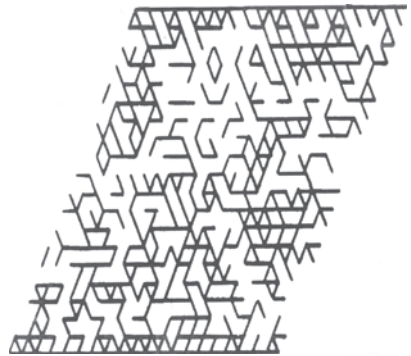plete fracture point (after Beale
and Srolovitz, 1988).



Strain = 0

Strain = 0.006

Strain = 0.134

FIGURE 8.10. Fracture pattern in a triangu-
lar network with stretching and bond-bending
forces and a broad distribution of the critical
lengths $l_c$.



breaking by bending, and $\eta$ having no apparent physical meaning. A beam was re-
moved at the center of the lattice in order for microcracking to be initiated. Figure
8.11 presents three fracture patterns. The one at the top was produced by using a
breaking criterion based on Eqs. (74) and (75), where $p_b$ was assumed to be given
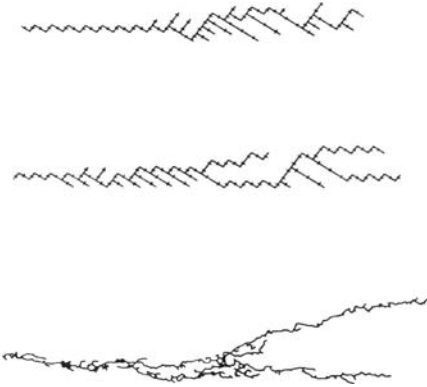
FIGURE 8.11. The top and the middle figures show the fracture patterns (rotated by 45°) in a 60 × 60 network. The pattern at the bottom is the experimental fracture pattern (see the text) (after Herrmann *et al.*, 1989b).

by Eq. (76), and $q = 0$, $\beta = 1$ and $\eta = 1$, while the middle one was generated based on Eq. (73) with $\beta_0 = 1$ and $\eta = 0.2$. Evidently, decreasing $\eta$ increases the tendency of the growing fracture to branch out. Also shown at the bottom of Figure 8.11 is the fracture pattern in an alloy, $Ti_{11.5}Mo_6Zr_{4.5}Sn$, aged 100 hours at 750 K and cracked under increasing stress intensity. There are certain similarity between the experimental pattern and what is shown in the middle of the figure.

Let us point out that, under certain conditions, the set of the broken bonds (the microcracks) can form a fractal network. The fractality of this set, and that of the sample-spanning fracture, will be discussed later in this chapter.

### 8.2.1.2   Dependence of the Elastic Moduli on the Extent of Cracking

One of the most interesting bits of information, which is also experimentally accessible, is the behavior of the elastic moduli of the system as the breaking process proceeds. Shown in Figure 8.12 are the Young's moduli $Y$ of three distinct models undergoing fracture (Arbabi and Sahimi, 1990b) and their dependence on the fraction $p$ of the unbroken bonds. Three distinct sets of results are shown. One is for a BCC lattice with central-force springs (recall from Chapter 8 of Volume I that in 3D the rigidity percolation thresholds of the simple-cubic lattice of Hookean springs with no BB forces is 1, and therefore it cannot be used in such simulations), while the second set is for a simple-cubic lattice in which both the central and bond-bending forces contribute to its deformation. The third set is the Young's modulus of a percolating simple-cubic network with central and bond-bending forces in which a fraction $1 - p$ of the bonds has been removed *at random*. Clearly, a bond breaking process in which the bonds are broken according to the force that they suffer, or according to the deviation of their length from a threshold value, weakens the system much faster than a random (percolation) bond-breaking process and, as a result, the system fails much sooner than a percolation system. Also shown in this figure are some experimental data on the Young's modulus of ceramic and glasses during their microcracking. The data were in terms of the microporosity of the materials, which was taken to be proportional to $(1 - p)$. As can be seen, the predictions of the fracture model with the central and bond-bending forces are
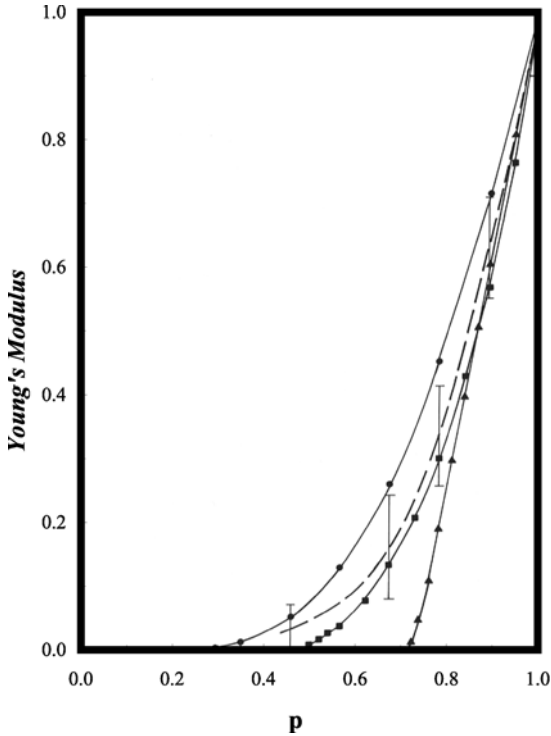
FIGURE 8.12. Young's modulus of the central-force (triangles) and the bond-bending ($\Box$) models (with $\gamma/\alpha = 0.04$) versus the fraction $p$ of unbroken bonds. Also shown are the modulus of a simple-cubic lattice in bond percolation (circles), as well as the experimental data (dashed curve) (after Arbabi and Sahimi, 1990b).

within the range of the experimental data. Moreover, except for $p \simeq 1$, the results produced by the central-force model do not agree with the experimental data because such systems fail at high values of $p$. For $p \geq 0.5$, the results predicted by the random percolation model do not agree with the data as well as the lattice model of fracture, presumably because the percolation threshold of the lattice is low ($p_c^B \simeq 0.25$) and, as a result, the predicted modulus is somewhat large.

An interesting study by Curtin (1997) demonstrated the effect of disorder on quasi-static fracture. He investigated the fracture toughness of heterogeneous materials using a simple-cubic lattice of springs with distributed toughness. He found that the overall toughness of the lattice or, equivalently, the stress $\sigma_f$ to initiate the first microcrack, is a random variable that depends on the width of the toughness distribution of the individual springs. This by itself is not surprising (see below). What was interesting was the finding that, for narrow distributions, the toughness was found to be controlled by the nucleation of the kinks at the weakest springs, whereas for broad distributions the toughness was controlled by the highly rigid regions of the system that pin the growth of the fracture front. However, the difference between the toughnesses of materials with narrow and broad distributions

was found to be small, hence suggesting that simple disorder *alone* (such as a narrow and uniform distribution of the thresholds) cannot solely be responsible for the variety of fracture behavior seen in experiments, and more complex factors must play an important role. We will come back to this issue later in this chapter.

### 8.2.1.3    Fracture Strength of Materials with Strong Disorder

The failure strength $\sigma_f$ of a solid is usually determined in a tensile test. There are several definitions of the failure stress which depend on the nature of the tensile test. In a *stress-controlled* test the sample fails at the highest stress in the stress-strain diagram. In the context of network models considered here, this usually occurs at the point where the first bond breaks. In a *strain-controlled* test, on the other hand, the strain is incremented and stress is the dependent variable. As the stress is finite for all strains, the failure stress in this case corresponds to the point where the stress first drops to zero. We define stress, or fracture, strength $\sigma_f$ of a system as the lowest externally applied stress at which the system breaks down. One can hypothesize that the eventual failure of the system is governed by the most critical flaw in the system, i.e., the weakest part of the system. Hence, calculation of the full distribution function of fracture strength $\sigma_f$ reduces to the computation of the distribution function of the most critical flaw in the system. It can be shown that this is an excellent approximation for the failure stress of the system in a *stress-controlled* tensile test. However, the fracture strength is not a *self-averaging* property. That is, the distribution of the fracture strength of even large samples of nominally the same material is *not* a delta function, implying that the *most probable* fracture strength and its *average* are *not* equal. Therefore, the full distribution must be computed or measured.

We first consider fracture strength of a strongly disordered material. The material's heterogeneity is generated by a percolation algorithm, i.e., before deformation and cracking of the material begins, only a fraction $p$ of the system is solid material; the rest of it is pre-fracture vacancy or voids generated by some type of defect. This means, in the context of the lattice models discussed in this chapter, that a fraction $p$ of the bonds (springs or beams) are present and the rest are cut *before* deformation of the lattice begins.

Monte Carlo simulations of quasi-static fracture of triangular and simple-cubic lattices of springs with central and bond-bending forces were carried out by Sahimi and Arbabi (1993) to check whether the fracture strength of materials with strong (percolation-type) disorder exhibits universality near the percolation threshold. We assume that

$$\sigma_f \sim (p - p_c)^{T_f}, \tag{77}$$

where $T_f$ is a new critical exponent that is *not* equal to $f$, the critical exponent that characterizes the elastic moduli of linear elastic materials near the percolation threshold $p_c$. Equation (77) is now rewritten as, $\sigma_f \sim \xi_p^{-\hat{T}_f}$, where $\xi_p$ is the correlation length of percolation and $\hat{T}_f = T_f/\nu$, with $\nu$ being the critical exponent of $\xi_p$. Thus, one may use the standard finite-size scaling method to estimate $\hat{T}_f$.

Briefly, according to the finite-size scaling theory, for length scales $L \ll \xi_p$, we must replace the correlation length $\xi_p$ by $L$, the linear size of the network, and hence in this regime $\sigma_f \sim L^{-\hat{T}_f}$. Thus, one carries out a series of simulations with lattices of various sizes, at the percolation threshold $p_c$ of the infinite lattice, and estimates $\sigma_f$ for each $L$. Since only small and moderate $L$ can be used, the results are fitted to $\sigma_f = L^{-\hat{T}_f}[a_0 + a_1 h_1(L) + a_2 h_2(L)]$, where $a_0$, $a_1$ and $a_2$ are three constants, and $h_1(L)$ and $h_2(L)$ are correction-to-scaling functions that take into account the effect of finite sizes of the network. The most accurate estimate of $\hat{T}_f$ is obtained with (Sahimi and Arbabi, 1991), $h_(L) = (\ln L)^{-1}$ and $h_2(L) = L^{-1}$.

Monte Carlo simulations in 2D (Sahimi and Arbabi, 1993) yielded, $T_f \simeq 2.42 \pm 0.14$, in good agreement with the measurement of Benguigui et al. (1987), $T_f \simeq 2.5 \pm 0.4$, described in Section 8.2.2, while in 3D (Sahimi and Arbabi, 1993), $T_f \simeq 2.64 \pm 0.30$. In addition to the Monte Carlo results of Sahimi and Arbabi (1993), less extensive molecular dynamics computations of Ray and Chakrabarti (1985b), and lattice simulations of Beale and Srolovitz (1988) also seemed to support the validity of Eq. (77).

If, instead of a lattice, one considers a continuum in which the distribution of the local transport properties may have certain singular properties, then one has a new universality class for the transport exponents. The same is true about the exponent $T_f$. For example, if one punches holes in the material at random, then the material, before its deformation begins, resembles the Swiss-cheese model, i.e., one in which spherical inclusions (or circular inclusions in 2D) are randomly distributed in a uniform matrix. The elasticity exponent $f$ of the Swiss-cheese model, that characterizes the power-law behavior of its elastic moduli near the percolation threshold $p_c$, is larger than that of the lattice model. Therefore, we may also expect higher values of $T_f$ in order to describe the power-law dependence of the fracture strength of the Swiss-cheese model near the percolation threshold, and this is indeed true. In fact, using a scaling analysis, Chakrabarti et al. (1988) proposed that for the Swiss-cheese model,

$$T_f \simeq \frac{1}{2}(f + d\nu + x) = \frac{1}{2}(f_{sc} + d\nu), \tag{78}$$

where $f_{sc} = f + x$ is the elasticity exponent of the Swiss-cheese model, and $x = 3/2$ and $5/2$ for $d = 2$ and $3$, respectively. Thus, for a 2D material, one obtains $T_f \simeq 4.06$, in good agreement with the experimental measurement of Benguigui et al. (1987) who also measured the fracture strength of a 2D Swiss-cheese model (by punching holes into their material) and reported that, $T_f \simeq 4.0 \pm 0.1$.

Rigorous upper and lower bounds for $T_f$ have also been derived (Ray and Chakrabarti, 1985a,1988). In a lattice near the percolation threshold, representing a composite material with strong disorder, the nearest-neighbor uncut bonds or sites form a type of "super-lattice"—a very large cluster made of tortuous links (long strands made of many bonds)—that cross each other at nodes that are separated by an average distance $\xi_p$, the correlation length of percolation. Therefore, in a $d$-dimensional system, the external stress $\sigma$ is shared by $\xi_p^{d-1}$ number of parallel

links, implying that the stress per link is $\sigma_l \simeq \sigma \xi_p^{d-1}$. Then, the total strain is, $\epsilon = \sigma/Y \simeq \sigma \xi_p^{f/\nu}$, where we have used the power-law dependence of the Young's modulus on $p$ near $p_c$. The total strain is shared by $\xi_p^{-1}$ number of links, which means that the strain per link is given by, $\epsilon_l \simeq \sigma \xi_p^{1+f/\nu}$. Therefore, the elastic energy per link is given by

$$\mathcal{H}_l \simeq \sigma_l \epsilon_l \simeq \sigma^2 \xi_p^{d+f/\nu}. \tag{79}$$

We now recall that only the bonds in the backbone of the lattice, i.e., the multiply-connected bonds of the sample-spanning cluster, contribute to stress and/or strain transfer in the system, the total number $M$ of which is, $M \propto \xi_p^{D_{bb}}$, where $D_{bb}$ is the fractal dimension of the backbone. If $\mathcal{H}_l$ is shared equally by all such bonds, then the energy per bond $\mathcal{H}_l/M$ is of the order of $\sigma^2 \xi_p^{d-D_{bb}+f/\nu}$. However, because of the assumption of equal sharing by all the $M$ bonds of the links, this value of $\mathcal{H}_l/M$ underestimates the strain energy per bond, since many of the bonds either do not support any stress at all or support very small stresses, because the backbone of a percolation lattice is multiply connected. Thus, if we assume a fixed energy threshold for breaking each bond, lattice fracture will occur for $\sigma_f = \sigma$ for which $\mathcal{H}_l/M \sim \sigma^2 \xi_p^{d-D_{bb}+f/\nu}$ exceeds the threshold. Therefore, if $\sigma_f$ follows Eq. (77), then, because we underestimate $\mathcal{H}_l/M$, we must have

$$T_f \geq \frac{1}{2}[f + (d - D_{bb})\nu]. \tag{80}$$

On the other hand, recall that the number of singly-connected (red) bonds, i.e., those that, if cut, would split the backbone into two pieces, is $M_r \sim \xi_p^{1/\nu}$. If all the strain energy is supported by such bonds, then one obtains an overestimate of the elastic energy per bond, $\mathcal{H}_l/M_r \sim \sigma^2 \xi_p^{d+f/\nu-1/\nu}$. Using the same argument that we utilized for the lower bound, we obtain

$$\frac{1}{2}[f + (d - D_{bb})\nu] \leq T_f \leq \frac{1}{2}(f + d\nu - 1). \tag{81}$$

Using the numerical values of the various exponents, $f \simeq 3.96$ and $3.75$, and $D_{bb} \simeq 1.675$ and $1.87$, for 2D and 3D systems, respectively, we obtain, $2.22 \leq T_f \leq 2.81$ for $d = 2$, and $2.37 \leq T_f \leq 2.7$ in $d = 3$. These bounds are perfectly consistent with the experimental and numerical estimates of $T_f$ given above.

Somewhat sharper bounds were proposed by Bergman (1986) who suggested that

$$f - \nu D_{min} \leq T_f \leq f - 1, \tag{82}$$

where $D_{min}$ is the fractal dimension of the shortest paths, or the chemical paths, on the backbone of percolation clusters. These bounds, together with the estimates, $D_{min} \simeq 1.13$ and $1.34$ for 2D and 3D systems, respectively, imply that

$$2.45 \leq T_f \leq 2.96, \quad d = 2, \tag{83}$$

$$2.58 \leq T_f \leq 2.76, \quad d = 3, \tag{84}$$

which agree nicely with the simulation results. Moreover, in both 2D and 3D the estimated $T_f$ is close to the lower bound (82), and therefore the relation $T_f = f - v D_{min}$ cannot be ruled out.

### 8.2.1.4  Distribution of Fracture Strength

Traditionally, the Weibull distribution (WD) has been used in fitting the fracture strength data for various materials. For a sample of a material of linear size $L$, the cumulative WD is given by

$$F_L(\sigma_f) = 1 - \exp(-cL^d \sigma_f^m), \tag{85}$$

where $c$ and $m$ are constant, and $d$ is the dimensionality of the system [see Eq. (18) which is the same as (85) with $cL^d = \lambda^{-m}$]. However, as discussed in Sections 5.25 and 5.4.2.1, Duxbury and Leath (1987) formulated a new Gumble distribution (GD) which is given by

$$F_L(\sigma_f) = 1 - \exp\left[-cL^d \exp\left(\frac{-k}{\sigma_f^\delta}\right)\right], \tag{86}$$

where $k$ and $\delta$ are also constant. Although Eq. (86) is supposed to be valid for materials that are far from their percolation threshold, i.e., those that have few pre-fracture vacancies or voids, as our discussion below indicates, it is not clear that this is actually the case. Both Eqs. (85) and (86) can be derived by appealing to percolation statistics, or any other theory that can provide the statistics of clusters of vacancies or voids in a material. Similar ideas were utilized to derive the failure distribution for electrical and dielectric breakdown of heterogeneous materials (see Chapter 5). We now describe the derivation of these two distributions for fracture strength of materials.

Suppose that a solid material of linear size $L$ contains $n$ cracks, each with failure probability $q_i(\sigma)$, $i = 1, 2 \cdots, n$, under an applied stress $\sigma$. To make the analysis manageable, we further assume that the stress-released regions of each of the cracks are separate and do not overlap, and that $F(\sigma)$ is the cumulative failure probability of the entire sample under stress $\sigma$. Then (Ray and Chakrabarti, 1985a)

$$1 - F_L(\sigma) = \prod_{i=1}^{n}[1 - q_i(\sigma)] \simeq \exp\left[-\sum_i q_i(\sigma)\right] = \exp[-L^d \rho(\sigma)], \tag{87}$$

where $\rho(\sigma)$ is the density of cracks *weaker* than the stress $\sigma$ or, equivalently, the density of the cracks that will start propagating at and above the stress $\sigma$. Equation (87) is due to the fact that the sample material survives if each of the cracks within it survives. We now show that the WD and GD correspond to two limiting cases of Eq. (87).

Consider first the derivation of the WD. This distribution arises if $\rho(\sigma)$ is given by a power law in $\sigma$, which will be the case if the density $\rho(l)$ of the linear cracks is of power-law type. This type of distribution does not arise in materials with a random distribution of vacancies or voids, unless the material is precisely at its

percolation threshold $p_c$, or is in a state above $p_c$ but is viewed at length scales $L \ll \xi_p$, where $\xi_p$ is the correlation length of percolation. Moreover, such power-law distributions do arise in real materials in which there are significant long-range correlations, without necessarily being near $p_c$. Suppose then that $g(l) \sim l^{-\tau}$, where $\tau$ is a critical exponent that characterizes the power-law distribution of the clusters of vacancies or voids (in the language of percolation theory, these are clusters of broken or unoccupied bonds). We further assume that the breaking stress (fracture strength) $\sigma_f$, beyond which a crack of length $l$ is created, is related to $l$ by

$$\sigma_f \sim \frac{1}{l^{1/\delta}}, \tag{88}$$

where the value of $\delta$ is determined by the nature of the surface of the crack in general, and its surface roughness in particular. Note that Eq. (88) represents a generalization of the Griffith law (see Chapters 6 and 7) which predicts a $1/l^{1/2}$ singularity for $\sigma_f$ (i.e., $\delta = 2$). This generalization is justified on the ground that the Griffith law was derived under the assumption that the cracks are smooth. However, as discussed in Sections 7.8.7 and 7.8.14, when the cracks are not smooth and have a rough and self-affine surface (which experiments show that this is often the case), then one must use the generalized Griffith law in which case the parameter $\delta$ can be treated as an adjustable parameter. If so, given Eq. (88), we can write $\rho(l) \sim \sigma_f^{\tau\delta}$ which, when substituted in Eq. (84), yields the Weibull distribution with the Weibull exponent $m = \tau\delta$. If, for example, we define $\sigma_q$ as the stress at which a certain percentage $q$ of the system fails (i.e., set $F = q$), we obtain

$$\sigma_q \sim \frac{1}{L^{d/m}}, \tag{89}$$

so that the stress depends strongly on the system's size.

Consider now the opposite limit in which the material is far from its percolation threshold. Then, we know from percolation theory that $\rho(l)$, i.e., the probability of having a cluster of vacancies (cracks) of linear size $l$ is given by

$$\rho(l) \sim \exp(-l/\xi_p) \sim \exp(-1/\xi_p\sigma_f^{\delta}), \tag{90}$$

where the generalized Griffith law, Eq. (88), was again invoked. Equation (90), when substituted in Eq. (87), yields an equation for $F_L(\sigma_f)$ which is equivalent to Eq. (86), which predicts that $\sigma_q$, the stress beyond which a certain fraction $q$ of the system fails, is given by

$$\sigma_q \sim \frac{1}{(\ln L)^{1/\delta}}, \tag{91}$$

which indicates a very weak logarithmic size dependence.

One can use the lattice models of fracture to directly test the accuracy of Eqs. (85) and (86) and the conditions under which they may provide accurate fits of experimental data. However, given that the results of such simulations may not be accurate enough to distinguish between the two distributions, each of which contains two adjustable parameters, a more sensitive test of the validity of these two

distributions can be made if we rewrite Eqs. (85) and (86) in alternative forms. If we define a quantity $A$ by

$$A = -\ln\left\{-\frac{\ln[1 - F_L(\sigma_f)]}{L^d}\right\}, \tag{92}$$

then the WD can be rewritten as

$$A_w = a_1 \ln\left(\frac{1}{\sigma_f}\right) + b_1, \tag{93}$$

while the GD is rearranged as

$$A_G = a_2\left(\frac{1}{\sigma_f^\delta}\right) + b_2. \tag{94}$$

These two equations predict linear variations of $A_w$ with $\ln(1/\sigma_f)$ and of $A_G$ with $1/\sigma_f^\delta$. The exact value of $\delta$ has not been determined, but in general $1 \le \delta \le 2$.

The conditions under which Eq. (85) or (86) may provide accurate representation of fracture strength data for a given material are not completely clear yet. It appears (Sahimi and Arbabi, 1993b) that in highly heterogeneous solids neither equation is very accurate, although the WD appears to perform better. On the other hand, in weakly-disordered materials, or those far from the percolation threshold (i.e., a material with few vacancies or voids), the GD may be a better representation of the distribution of fracture strengths. For example, Figure 8.13 presents the fit of the fracture simulation results to a GD with $\delta = 1$ for a triangular network with the central and bond-bending forces, in which before deformation and fracture of the lattice started, 10% of the bonds had been removed at random. Simulations with the central force (Beale and Srolovitz, 1988) and the Born models (Hassold and
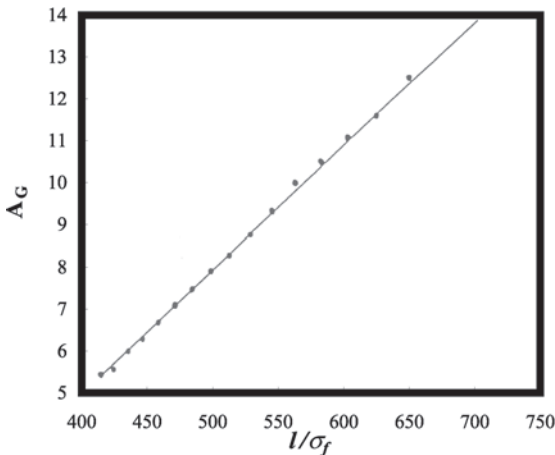


FIGURE 8.13. Fit of the simulation results for fracture of a triangular network with stretching and bond-bending forces (with $\gamma/\alpha = 0.1$) to the Gumbel distribution. Before fracturing, 10% of the bonds were removed at random (after Sahimi and Arbabi, 1993).

Srolovitz, 1989) also seem to indicate that the GD is accurate if the system is far from $p_c$, although Hassold and Srolovitz (1989) reported equally accurate fit of the data with the WD, and Hansen *et al.* (1989) also reported some deviations from the GD in their central-force model. Curtin and Scher (1992) discussed the conditions under which a WD may be appropriate for representing the distribution of fracture strength. We note, however, that as the percolation threshold is approached, neither distribution seems to be very accurate.

Fracture behavior of a material at its percolation threshold $p_c$, or equivalently above its $p_c$ but at length scales $L \ll \xi_p$, depends on the broadness of the distribution of its heterogeneities, and deserves to be discussed (Sornette, 1988). In a percolation system far from $p_c$, there are many multiply-connected paths, called macro-links, which support stress transport. In such a system, the distribution of fracture strength is the result of one or both of the following factors:

(1)  Fluctuations of the individual characteristics of the microscopic regions (for example, their breaking threshold $l_c$ and/or their elastic constants) of the material.

(2)  Fluctuations of the macro-link sizes $\mathcal{L}$ around the percolation correlation length $\xi_p$. If the characteristics of the microscopic regions are all the same (i.e., if the material is not made of different constituents with different properties), the first factor cannot contribute to the distribution of fracture strength. As $p_c$ is approached, two changes take place: First, one has fewer macro-links and, secondly, compared with those of the long macro-links, the contributions of the shorter macro-links to stress transport become negligible. Thus, macro-link to macro-link fluctuations also decrease. At $p_c$, there is only *one* huge macro-link, and therefore all the fluctuations disappear completely and the distribution of fracture strength must be a delta function. However, if, for example, the elastic constants of the microscopic regions of the material are statistically-distributed quantities (which is often the case in real heterogeneous solids), then region-to-region fluctuations exist and the distribution of fracture strength is a meaningful quantity to define, calculate, or measure. This is particularly important for disordered materials modeled by continuum percolation which usually possess a broad distribution of the elastic constants of the channels through which stress transport takes place. The distribution of fracture strength in such materials is a Weibull-like distribution, rather than the GD. This is supported by the lattice simulations of Sahimi and Arbabi (1993).

### 8.2.1.5  Size-Dependence of Fracture Properties

An important problem in fracture mechanics of materials is the dependence of the various properties on the sample size. Moreover, equations such as (77), (88) and (89) already state that certain features of fracture may have scale-invariant properties. Let us now summarize the equations that express the sample size dependence of certain fracture properties, and discuss whether they are universal. We describe here such properties for the discrete and deterministic models of fracture. Later in this chapter, we discuss the same properties for other models of fracture.

Similar to electrical and dielectric breakdown problems discussed in Chapter 5, an important property is the average failure stress as a function of the sample size $L$ at a fixed $p$, where $(1 - p)$ is the fraction of the vacancies or voids in the material (or $p$ is the fraction of the intact bonds in the lattice models) before fracture has begun. To obtain this quantity, we first note that, in Eq. (86), the constants $c$ and $k$ must depend on $p$. Thus, setting $F_L(\sigma_f)$ to any constant value between 0 and 1 and solving for $\sigma_f$, we obtain from Eq. (83) the following expression

$$\sigma_f^\delta = \frac{c}{A(p) + B(p) \ln L}, \tag{95}$$

where $A(p)$ and $B(p)$ are simple functions of $p$, and $c$ is a constant. In particular, $B(p)$ is small for $p$ close to unity, and $B(p) \to \infty$ as $p \to p_c$. Simulations with the central-force model (Beale and Srolovitz, 1988), the Born model (Hassold and Srolovitz, 1989), and the central-force and bond-bending model (Sahimi and Arbabi, 1993) all support the accuracy of Eq. (95).

In general, one is also interested in the scaling behavior of the external stress $\sigma$ or force $F$ for breaking the lattice and its dependence on its linear size, since in practice not only can this quantity be measured easily, but it also provides us with important insight into a material's structure (that is, whether the material is strong and difficult to break versus being weak and easily breakable, both of which depend on its morphology and size). Since this force is, for example, proportional to $\epsilon Y$, where $\epsilon$ is a displacement or strain, a plot of the stress $\sigma$, or the force $F$, versus $\epsilon$ can be compared with the traditional stress-strain diagrams that are measured for composite solids. However, changing the parameters of the lattice model will result in a wide variety of such diagrams. Therefore, instead of presenting the results for each model and lattice size separately (which would be impossible), we may try to collapse the data for all values of $L$, the linear size of the lattice, onto a single curve. If the data collapse is possible, then such diagrams possess universal properties which could be exploited for practical purposes.

Figure 8.14 presents (Sahimi and Arbabi, 1993) the results for the triangular lattice with sizes $L = 50$ and 70, in which both the central and bond-bending forces act on the springs. The thresholds were distributed according to Eq. (70) with $\zeta = 0$ (a uniform distribution). However, as can be seen, the data collapse is not complete and three distinct regimes can be discerned. The first regime represents the initial stages of crack growth and is far from the maximum of the curve. In this regime microcracking propagates at a relatively slow rate and is more or less similar to a percolation process, as the springs break essentially at random. As microcracking proceeds, one eventually arrives in the second regime in the vicinity of the maximum in which microcracking is intense and the lattice is relatively close to its macroscopic failure point. Beyond the maximum, the system is in the so-called *post-failure* regime, and is highly sensitive (unstable) to small variations in the applied stress or strain. The general features and shapes of these curves are in good qualitative agreement with direct experimental measurements and observations for brittle fracture in various types of disordered solids.
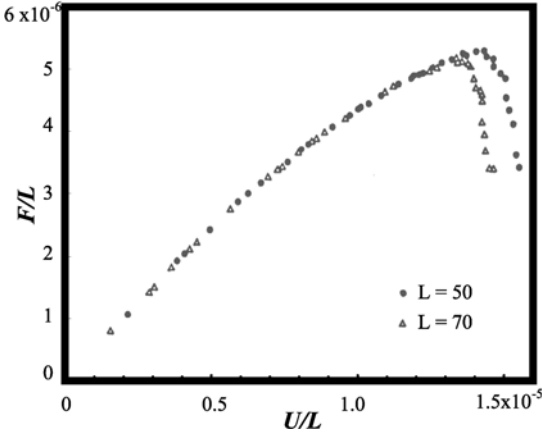
FIGURE 8.14. Collapse of stress-strain data for fracture of a triangular network with stretching and bond-bending forces (after Sahimi and Arbabi, 1993).

To obtain quantitative information on the scaling of $F$ with $L$, we may assume a scaling form. de Arcangelis *et al.* (1989) assumed that

$$F \sim L^{\Omega_1} \phi(\epsilon/L^{\Omega_2}), \qquad (96)$$

where $\Omega_1$ and $\Omega_2$ are two presumably universal exponents, and $\phi(x)$ is the scaling function. An equation similar to (96) was also used for data collapse in the electrical breakdown problem [see Eq. (5.31)]. Based on their simulations of the beam and the central-force models, and using Eqs. (70) and (71) as the distribution of the thresholds or heterogeneities, de Arcangelis *et al.* (1989) claimed that $\Omega_1 = \Omega_2$. Moreover, in 2D they found that, $\Omega_1 = \Omega_2 \simeq 0.75$. On the other hand, Arbabi and Sahimi (1990b) and Sahimi and Arbabi (1993) utilized the following equation

$$F \sim \frac{L^{\Omega_1}}{(\ln L)^\psi} \phi(\epsilon/L^{\Omega_1}), \qquad (97)$$

and argued that the logarithmic corrections, although seemingly weak, are necessary because their existence is predicted by approximate analytical theories [see Eq. (86) with $\psi = 1/\delta$]. If both $\Omega_1$ and $\psi$ are varied in order to obtain the most complete data collapse, one obtains, in 2D, $\Omega_1 \simeq 1 \pm 0.1$ and $\psi \simeq 0.1$. The value of $\psi$ is small and may thus be subject to relatively large uncertainties. Since simulation of very large networks is currently not feasible, the most accurate way of deciding whether Eq. (96) or (97) provides a more accurate fit of the data is by fitting the data to both equations and calculating the squared residual errors (i.e., the difference between the data and their predictions by the fitted equation) that each fit produces. Values of $\Omega_1$ and $\psi$ are insensitive to the parameter $\zeta$ of the distribution of the critical threshold, Eq. (70), unless $\zeta \to 1$. The results of de Arcangelis *et al.* (1989) were also insensitive to the parameters of the distributions that they used. The estimated values of the exponents for 3D systems are (Arbabi

and Sahimi, 1990b; Sahimi and Arbabi, 1993), $\Omega_1 \simeq 2 \pm 0.1$, and $\psi \simeq 0.2$, which, together with the results for 2D lattices, suggest that for a $d$-dimensional system

$$F \sim \frac{L^{d-1}}{(\ln L)^{\psi}} \phi(\epsilon/L^{d-1}), \tag{98}$$

where $0 \leq \psi \leq 0.2$. Our unpublished simulation results for quasi-static fracture of a BCC lattice with central forces and lattice sizes of up to $L = 32$ indicated that the results do follow Eq. (98). Note that Eq. (98) has a simple interpretation: $L^{d-1}$ is the surface area on which the external force is applied, and $(\ln L)^{\psi}$ is the manifestation of the sample-size effect on the fracture process.

One can also study the variations of $F$ with $N_b$, the number of bonds that break during fracture. We assume that

$$F \sim L^{\Omega_3} \phi(N_b/L^{D_f}). \tag{99}$$

An equation similar to (99) was used in the electrical breakdown problem; see Chapter 5. Equation (99) implies that if the force $F$ is plotted versus $N_b/L^{D_f}$, then the results for various lattice sizes and parameters of the model should collapse onto each other. Figure 8.15 presents such a data collapse (Sahimi and Arbabi, 1993). de Arcangelis *et al.* (1989) found for their 2D models that, $\Omega_3 \simeq 0.75$ and $D_f \simeq 1.7$, consistent with their results obtained with Eq. (96). On the other hand, Arbabi and Sahimi (1990b) and Sahimi and Arbabi (1993) found that $\Omega_3 \simeq 1 \pm 0.05$, and $D_f \simeq 1.7 \pm 0.1$ in 2D, and $\Omega_3 \simeq 2 \pm 0.1$, and $D_f \simeq 2.3 \pm 0.2$ in 3D, consistent with their results and Eqs. (97) and (98). Note that $D_f$ represents the fractal dimension of the set of all the broken bonds.

As discussed by de Arcangelis *et al.* (1989), there are other interesting scaling features of these models. For example, one can study how the stress at the maximum of the diagram of the type shown in Figure 8.14 scales with the linear size $L$ of
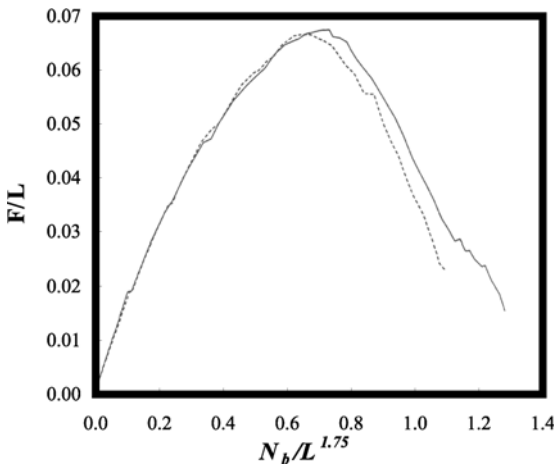


FIGURE 8.15. Collapse of the data according to Eq. (99) (after Sahimi and Arbabi, 1993).

the sample. Similarly, one can look at the scaling of the number of the broken bonds $N_b$ that corresponds to this stage of the fracture at the maximum. Finally, the number of broken bonds at the end of the fracture process also scales with $L$ as $L^{D_f}$, with $D_f \simeq 1.7$ in 2D.

### 8.2.2    Comparison with the Experimental Data

Experimental measurements of the fracture strength of composite materials with percolation disorder have been carried out. Benguigui *et al.* (1987) measured the strain and stress of a perforated metal (aluminum or copper) sheets and of a 2D diluted elastic network near $p_c$. The samples were prepared by two different methods. In one, random holes, with a diameter slightly larger than the length of a bond of a square lattice drawn on the material, were punched in the material. The material so prepared was then used in a strain-controlled experiment in which it was elongated by applying a strain to the material and increasing it monotonically and continuously until it failed macroscopically. The stress was measured during the entire experiment. In the second experiment the material was prepared by cutting the interhole bonds at random such that the hole size was smaller than the lattice unit cell. A stress was then applied to the material and increased monotonically, during which the elongation was measured.

Benguigui *et al.* (1987) found that in strain-controlled experiments the stress necessary to break the first bond was always larger than those for breaking the subsequent bonds, and thus fracture occurred by a cascade effect. As discussed earlier, in strain-controlled experiments, the fracture strength is defined as the stress necessary to break the first bond, whereas in stress-controlled experiments it is the stress by which the sample breaks. Figure 8.16 presents their experimental data. Moreover, Benguigui *et al.* (1987) found that, as $p_c$ was approached, the fracture stress $\sigma_f$ vanishes according to power law (77) with $T_f \simeq 2.5 \pm 0.4$, and therefore, as expected, $T_f$ is not identical with the elasticity exponent $f$ of elastic percolation networks $f(d = 2) \simeq 3.96$. This estimate of $T_f$ is completely consistent with the Monte Carlo estimate of Sahimi and Arbabi (1993), $T_f \simeq 2.42 \pm 0.14$, mentioned above, and also with the theoretical bounds (81) and (82).

Benguigui *et al.* (1987) also found that the stress for elongation of the sample appeared to diverge at $p_c$ according to a power law with an exponent $1.4 \pm 0.2$. Similar, but less precise, experiments were carried out by Sieradzki and Li (1986) who measured the fracture stress of a system composed of a 2mm-thick plate of aluminum with holes punched at positions corresponding to a triangular network of 21 rows and 20 columns. The fracture stress was determined by obtaining the full load-displacement curve for the sample to failure.

The validity of Eqs. (85) and (86) for representing the distribution of fracture strength of composite materials was also tested against experimental data. van den Born *et al.* (1991), who measured the mechanical strengths of highly porous ceramics, reported that the size dependence of their data was well described by both distributions, but for the failure pressure dependence, the GD with $\delta = 1$ was
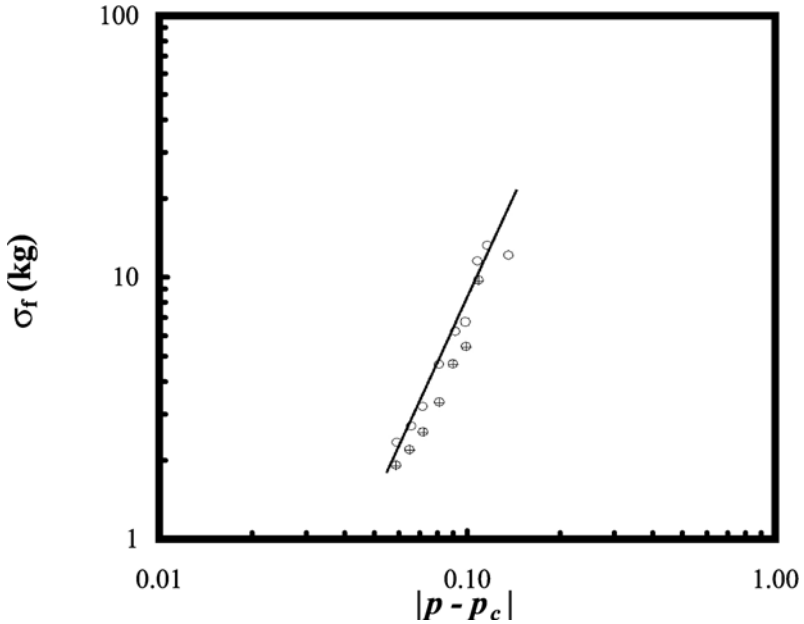
FIGURE 8.16. Logarithmic plot of the experimental data for fracture load (open circles) and yield load (filled circles) versus $p - p_c$ (after Benguigui *et al.*, 1987).

found to be more accurate. Figure 8.17 presents their data and their fit to the GD. Evidently, the pre-fracture porosity of the material used in these experiments was very low, so that the porous ceramic was far above its percolation threshold.

In another set of interesting experiments, Li and Sieradzki (1992) studied mechanical breakdown of random porous Au, a new material specifically designed for their experiments. They used digital image analysis to characterize the microstructures of their samples which varied by more than two orders of magnitude in length scale. The porous Au underwent a microstructurally controlled ductile-brittle transition. Such a transition had already been predicted by Sahimi and Goddard (1986), based on the broadness of the heterogeneity distribution in their lattice model of fracture, and in the numerical simulations of Kahng *et al.* (1988) using the scalar (fuse) model fracture processes (see Chapter 5). These results provide support for usefulness of the lattice models for describing quasi-static fracture processes in composite solids. Other relevant experimental data are discussed below.

## 8.2.3  *Percolation Versus Quasi-static Brittle Fracture*

In practice, in addition to the macroscopic properties of a fractured material, one is also interested in the distribution of its *local* properties, such as the distribution of the forces that are exerted on various parts of the material. This is an important distribution as it identifies the regions of the material that may break, which can then
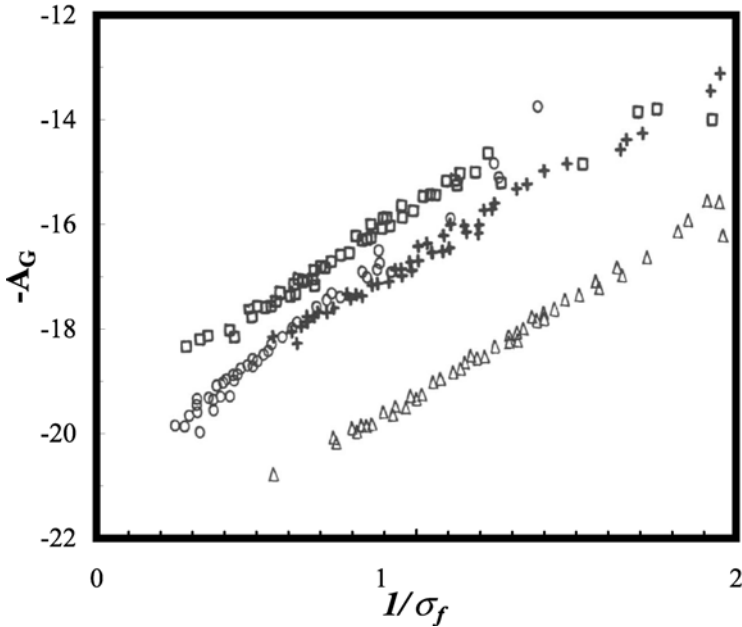
FIGURE 8.17. Fit of experimental data for fracture strength of porous silicon extrudates to the Gumbel distribution (after van den Born *et al.*, 1991).

be used in the design of tougher or better composites. Constructing this distribution can also help one to understand the similarities between fractured and percolation systems. Percolation is usually a static process in which failure of a bond has nothing to do with the stress or strain field in the lattice. On the other hand, the growth of cracks in a disordered solid is in general a dynamic and nonlinear phenomenon which does not occur at random, but depends upon the stress or strain field in the solid. However, under certain experimental conditions the accumulation of damage and the growth of cracks in a solid occur essentially at random as in, for example, a solid material that is under rapid thermal cycling, or a solid in which the heterogeneities are broadly distributed (Sahimi *et al.*, 1993), e.g., natural rock, in which case a percolation process may be able to describe fracture. Moreover, in most cases, percolation phenomena typically represent second-order (continuous) phase transitions, whereas many fracture processes that take place in nature resemble first-order phase transitions, although the precise nature of the geometrical phase transition that occurs during fracture—first order versus second order—remains controversial (see, for example, Zapperi *et al.*, 1997; Andersen *et al.*, 1997; Moreno *et al.*, 2000). Therefore, it is important to understand the extent of the similarities between fracture and percolation processes, because if there are similarities between the two, then percolation phenomena, which are well-understood and also much easier to study, may help us gain a deeper understanding of fracture of disordered solids.

There are at least two ways of comparing a fractured lattice with a percolating one. The first method is based on comparing the force distributions (FDs), i.e., the distributions of the forces that are exerted on the bonds of the lattice, and their moments in the two systems. We already described in Chapter 8 of Volume I (see Sections 8.6.3 and 8.11.4) the FD of elastic percolation networks, first computed by Sahimi and Arbabi (1989). As discussed by Sahimi and Arbabi (1993), the initial stages of brittle fracture and percolation processes in a lattice are more or less similar. That is, during the initial stages of fracture growth, the bonds that break are distributed essentially at random in the lattice, unless the lattice is uniform, or its disorder is very weak. In these initial stages, the stress enhancement at the tip of a given microcrack is not strong enough to ensure that the next bond that breaks would be at the tip of the present microcrack. However, as more microcracks nucleate the effect of stress enhancement becomes stronger, and deviations from random percolation increase. Beyond a certain point in the growth of the cracks, there will be no similarity between the two processes. Hence, one is naturally interested to locate the point at which a fracture process starts to deviate from a percolation phenomenon. The key clue is already provided in the stress-strain diagrams discussed above. Equations (96)–(98) are manifestation of finite-size scaling which represent the fracture data for various lattice sizes up to the maximum of the stress, beyond which it breaks down. This type of finite-size scaling is also valid for percolation networks for *any* $p$ (the fraction of the intact bonds) in the interval $p_c \leq p \leq 1$ (so long as the linear size $L$ of the network, $L < \xi_p$, where $\xi_p$ is the correlation length of percolation), albeit with different exponents and scaling functions. Therefore, in the type of disordered lattices that we are considering here, fracture and percolation are more or less similar up to the maximum in the stress-strain diagram of the fractured system, i.e., in the regime in which finite-size scaling is applicable to the fracture data, but they are not similar beyond this point.

Similar to elastic percolation networks analyzed in Chapter 8 of Volume I, one can also calculate the moments $M(q)$ of the FD in a fractured lattice and study their scaling with the lattice's linear size. Thus, one writes

$$M(q) = \sum_i n_{F_i} F_i^q, \tag{100}$$

where $n_{F_i}$ is the number of bonds that suffer a force with a magnitude $F_i$. Similar to elastic percolation networks, each moment of the FD scales with the linear size $L$ of the fractured lattice as

$$M(q) \sim L^{-\tilde{\tau}(q)}. \tag{101}$$

Herrmann *et al.* (1989a) and de Arcangelis *et al.* (1989) calculated the corresponding exponents $\tilde{\tau}(q)$ for 2D fractured lattices at two different points in the system. One was just before the lattices failed and a sample-spanning fracture was formed. Figure 8.18 presents their results, indicating that each moment of the FD of fractured lattices scales with the linear size $L$ of the system with a different exponent. The second point was at the maximum of the stress, where the system is entering
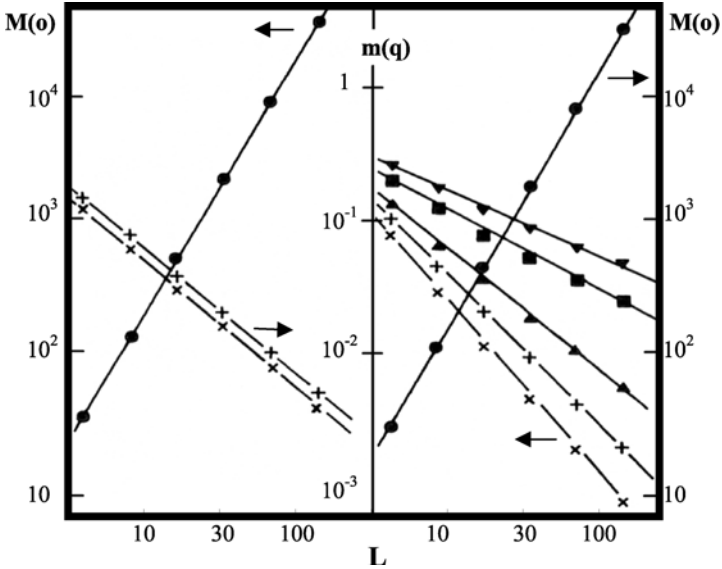
FIGURE 8.18. Rescaled moments $m(q) = [M(q)/M(0)]^{1/2}$ versus the linear size $L$ of the lattice. The left figure gives the moments at the maximum current (where there is constant gap scaling), whereas the right figure gives the moments right before the network fails (after de Arcangelis *et al.*, 1989).

the post-failure regime. In this case, the exponents that characterize the moments of the FD near $p_c$ can be obtained from one of the exponents.

To understand better the difference between percolation and fractured lattices, one important point to remember is that, while elastic properties of percolation networks are controlled by the *low* moments of their FD (for example, the elastic moduli are proportional to the second moment of the FD), fracture properties of the same systems are controlled by the *high* moments of the FD. This is due to the fact that fracture and breakdown occur where the *largest* loads (for example, stress) are concentrated in the system, and the effect of such regions is manifested only by the high moments of the FD.

## 8.2.4   Universal Fixed Points in Quasi-static Brittle Fracture

We now discuss another universal, and potentially very useful, aspect of quasi-static fracture of disordered materials. In a series of simulations, Sahimi and Arbabi (1992) computed three properties of a lattice as it underwent deformation and brittle fracture. First, the thresholds $l_c$ were distributed according to Eq. (70) and the elastic modulus $C_{11}$ of the lattice during fracture was measured. Next, the *same* fully-connected lattice (i.e., with the *same* values of $l_c$) was used to measure the shear modulus $\mu$ of the lattice during fracture, caused by shearing the system. This is equivalent to fracturing two identical samples and measuring their $C_{11}$ and $\mu$
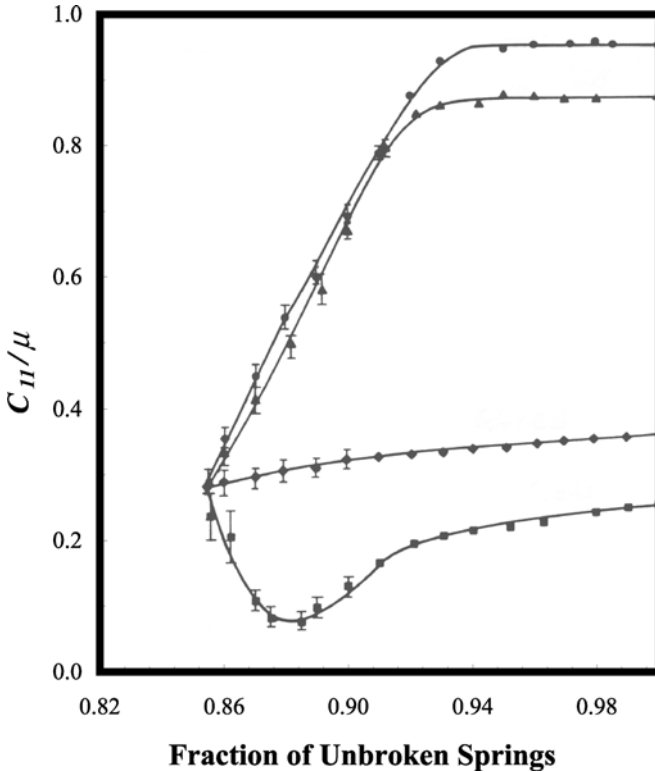
FIGURE 8.19. Dependence of the ratio $r = C_{11}/\mu$ on the fraction of unbroken bonds during fracture of a triangular lattice with stretching and bond-bending forces. The results, from top to bottom, are for $\gamma/\alpha = 0.0, 0.01, 0.3$ and $1$ (after Sahimi and Arbabi 1992).

independently. In Figure 8.19 we present typical results for the ratio $r = C_{11}/\mu$ as a function of the fraction of unbroken springs, for various values of $\gamma/\alpha$, the ratio of the elastic constants of the bond-bending model [see Eqs. (66) and (67)]. The last points of these curves represent $C_{11}/\mu$ right before the system failed macroscopically. We refer to this as the incipient fracture point (IFP). As can be seen, even though the initial states of the systems (i.e., their initial values of $r$ with no spring broken) are different, they all approach the same value of $r$ as the IFP is approached. Note also that, initially $r$ remains essentially constant, implying that the initial value of $r$ is not sensitive to the presence of a few cracks or even a collection of localized cracks. However, as damage accumulates and the cracks grow, a turning point (TP) appears and $r$ changes drastically. Because $\gamma/\alpha = 0$ corresponds to a lattice in which only central forces are present, Figure 8.19 indicates that this behavior is independent of the microscopic force law of the system. The behavior of the system for $\gamma/\alpha = 1$ is particularly interesting. Initially, $r$ remains essentially constant. However, as damage accumulates a TP appears beyond which $r$ decreases and reaches a minimum. But near the IFP, $r$ rises again

and approaches its value at the IFP which appears to be universal. Simulations (Sahimi and Arbabi, 1992) indicated that the value of $r$ at IFP is *universal* and independent of $\gamma/\alpha$ and $\zeta$, the parameter of the threshold distribution, Eq. (70) (unless, of course, $\zeta \to 1$). For 2D *isotropic* lattices (for example, a triangular lattice) simulations indicated that

$$\frac{C_{11}}{\mu} \simeq \frac{5}{4}. \tag{102}$$

The emergence of a universal fixed point at the IFP may mean that in a disordered solid that undergoes quasi-static brittle fracture, the approach of $r$ to its universal value at the IFP may be interpreted as the "signature" of the material's approach to its global failure point. Although Figure 8.18 indicates that for certain values of $\gamma/\alpha$ one may have a non-monotonic variation of $r$ with the accumulated damage (which, from an experimental view, implies that the closeness of $r$ to its universal value cannot be used for detecting the approach of the system to its global failure point), for most real materials one has $\gamma/\alpha \leq 0.3$, and for such values of $\gamma/\alpha$ the approach of $r$ to its value at the IFP is *always* monotonic.

What is the theoretical explanation for the apparent universality of $r$? It is not difficult to show that the dependence of $C_{11}$ and $\mu$ on the fraction of unbroken springs, as the IFP is approached, is similar to each other. As such, $r$ represents an *amplitude ratio*, and it is known (Aharony, 1980) from statistical mechanics that certain amplitude ratios are *universal*. The apparent universality of $r$ may mean that, much like renormalization group theory of critical phenomena, universal fixed points may be used for classifying various fracture processes. To do this, we recall from Section 8.11.8 of Volume I that, it has been suggested (Bergman and Kantor, 1984; Schwartz *et al.*, 1985; Arbabi and Sahimi, 1988) that, in elastic percolation networks $C_{11}/\mu$ takes on a universal value as $p_c$ is approached. For 2D isotropic percolation networks near $p_c$ one has, $C_{11}/\mu \simeq 3$, which is different from the corresponding value for brittle fracture considered here. Sahimi and Goddard (1986) suggested that brittle fracture of very heterogeneous networks is more or less similar to a percolation process, and Roux *et al.* (1988) presented evidence that in the limit of infinite disorder [i.e., the limit $\zeta \to 1$; see Eq. (70)] the quasi-static lattice models of brittle fracture represent a type of percolation process. For example, fracture in natural rock, a highly heterogeneous solid with scale-dependent properties, may be a realization of this (Sahimi *et al.*, 1993). On the other hand, in most solids disorder is finite, and simulations indicated that even for $\zeta = 0.8$ the value of $r$ at IFP is very different from that of elastic percolation networks at $p_c$, indicating that the limit $\zeta = 1$ may be a type of singular point, so that even for $\zeta = 1 - \delta (\delta \ll 1)$, one should still obtain the value of $r$ at the IFP described here and not that of percolation networks at $p_c$. Note that for 2D isotropic systems, the Poisson's ratio is given by $\nu_p = 1 - 2/r$, implying that for isotropic materials at the IFP, $\nu_p$ takes on a universal value.

Therefore, it has been proposed (Sahimi and Arbabi, 1992) that, value of the Poisson's point at the IFP may be used to classify various universality classes

of quasi-static brittle fracture processes in disordered solids. Specifically, it has been proposed that there are *two* distinct universality classes. One is for weakly-disordered materials that are under a uniform external load (stress or strain). In such solids the growth of a crack at a point depends on the environment around that point, and therefore the damage accumulation is *not* at random. The universality class of such solids is described by the fixed point described by Eq. (102). Examples include most engineering solid materials that are typically not too heterogeneous, with a few defects, or small (laboratory-scale) pieces of rock that are microscopically disordered but macroscopically homogeneous with no large scale variations in their elastic properties. The second fixed point is for systems in which damage accumulates essentially at random. Such solids, which include highly heterogeneous media, belong to the universality class of the fixed point of elastic percolation networks at the $p_c$. Examples may include natural rock at large length scales (of the order of a few hundred meters or more) with spatially-varying properties, and solids that undergo rapid thermal cycling.

An important question is whether the universality of the Poisson's ratio at the IFP can be tested experimentally. One system in which these ideas can be tested is natural rock. It has been suggested (Sahimi and Arbabi, 1992,1996; Sahimi *et al.*, 1993; Robertson *et al.*, 1995) that the quasi-static models of brittle fracture that we have described so far may also describe fracture of natural rock, since rock fracture is an extremely slow process. If so, then the fracture pattern and the universal fixed point predicted by the model should be observable in rock. This is in fact the case. Natural rock contains large fractures, in the form of a complex and interconnected network. Despite their obvious significance, characterization of fractured rock, and how the fractures are formed and become connected is not as well-developed as that of unfractured porous media. However, this is changing very fast now and such ideas as scaling, fractals, and percolation concepts are beginning to find their proper place in the field of characterization of fractured rock (for reviews see, for example, Sahimi, 1993b, 1994a, 1995b).

One of the first systematic studies of fractured rock was carried out by Barton and co-workers (Barton and Larsen, 1985, Barton *et al.*, 1987; Barton and Hsieh, 1989; Barton, 1992) as part of the effort by the United States Geological Survey to characterize the geologic and hydrologic framework at Yucca Mountain, Nevada, which is being considered by the United States Department of Energy as a potential underground repository for high-level radioactive wastes. Barton and Larsen (1985) developed the *pavement method* of clearing a subplanar surface and mapping the fracture surface in order to measure connectivity, trace length, density and fractal scaling of the fractures, in addition to their orientation, surface roughness and aperture. Each of these parameters is important in predicting the hydraulic characteristics of the network and in working out the history of its development in relation to the regional tectonics. The most significant observation of the Yucca Mountain study was that the fractured pavements have a fractal geometry, i.e., the fracture pattern is scale-invariant. Thus, it was possible to represent the distribution of the fractures ranging from 20 cm to 20 m by the fractal dimension $D_f$ defined

as

$$D_f = \frac{\log N_\ell}{\log(1/\ell)} \qquad (103)$$

where $N_\ell$ is the number of fractures of length $\ell$. Using Eq. (103) (the standard box-counting method described in Chapter 1), fractal dimensions of the fractured surfaces at Yucca Mountain were found to be in the range 1.5–1.7.

Another study was undertaken for the Geysers geothermal field in northeast California (Sahimi *et al.*, 1993). This field, from which heat and vapor are extracted for use in power plants generating electrical power, covers an area of more than 35,000 acres and is one of the most significant geothermal fields in the world. Using the box-counting method, Sahimi *et al.* (1993) determined the fractal dimension of the fracture surfaces of the Geysers field and found, as did Barton and co-workers, that at small length scales the fracture pattern is fractal with $D_f \simeq 1.5 - 1.7$, whereas at much larger length scales, $D_f \simeq 1.9$. These results were interpreted with the help of the quasi-static models of fracture described above (Sahimi *et al.*, 1993). In particular, note that this range of fractal dimension $D_f$ is essentially the same as what one finds with the quasi-static lattice models of fracture.

There is also convincing experimental evidence indicating that in fractured rock the Poisson's ratio $\nu_p$ may take on a universal value at the IFP, in agreement with the prediction of the quasi-static lattice models of brittle fracture described above. The existence of a universal fixed point in fractured rock can be directly tested by experimental measurements, since for 3D systems

$$\frac{V_p}{V_s} = \left( \frac{C_{11}}{\mu} + \frac{1}{3} \right)^{1/2}, \qquad (104)$$

where $V_p$ and $V_s$ are the velocities of the shear and compressional waves in the medium, respectively, which can be measured by established experimental procedures (Brace and Orange, 1968). Sammonds *et al.* (1989) fractured four sandstone samples at four different confining pressures, and measured $V_p$ and $V_s$. Different confining pressures result in different fracture patterns since they control the closure of pre-existing cracks and nucleation and growth of new microcracks. At the three lowest confining pressures the corresponding fracture patterns were found to be brittle-like, and from their results one finds (Sahimi and Arbabi, 1992) that $V_p/V_s \simeq 1.14 \pm 0.04$ at the IFP for all the three fractured sandstones, implying a universal value for $C_{11}/\mu$. At the highest confining pressure fracture was ductile-like, and although, as expected, the stress-strain diagrams of the sample was not similar to that of brittle fracture, their results indicated that even for this case $V_p/V_s \simeq 1.1$, beyond the point at which stress became independent of strain (which is typical of stress-strain diagrams of ductile fracture), consistent again with the value for brittle fracture at the lower confining pressures. These data provide strong experimental support for the existence of universal fixed points at the IFP. Note that since for 3D isotropic systems, we have

$$\nu_p = \frac{3(C_{11}/\mu - 1)}{2 + 6(C_{11}/\mu - 1)}, \qquad (105)$$

the experimental data of Sammonds *et al.* (1989) imply that for their samples, $\nu_p \simeq 0.1$ at the IFP. On the other hand, if we use $C_{11}/\mu \simeq 4/3$, the corresponding value for 3D isotropic elastic percolation networks at $p_c$, we find $\nu_p \simeq 1/4$. Thus, over the length scale used in these experiments, the sandstones examined by Sammonds *et al.* (1989) must have been relatively homogeneous, since their Poisson's ratio at the IFP ($\nu_p \simeq 0.1$) is different from that of percolation networks at $p_c$ ($\nu_p \simeq 1/4$), which corresponds to highly heterogeneous systems. These results also support the validity of classifying fracture processes according to the value of the Poisson's ratio at the IFP.

## 8.3   Dynamic Brittle Fracture

As discussed in Chapter 7, until a few years ago, there was a classical unsolved problem in dynamic fracture: While classical analysis based on the linear continuum mechanics would predict that brittle fracture in materials should speed up until it reaches the Rayleigh wave speed $c_R$, experiments indicated that the velocity of fracture propagation never reaches this limit. It typically reaches about 40% of the limit, and almost never more than 60% of it. The tip of the crack also heats up by hundreds of degrees (see, for example, Green and Pratt, 1974; Fuller *et al.*, 1975) and, moreover, it emits high-frequency waves (Gross *et al.*, 1993). Despite many attempts and many proposed mechanisms for these experimental observations, the problem remained unsolved for many years. Lattice models of quasi-static fracture described and discussed earlier in this chapter cannot obviously shed any light on the resolution of these issues. However, a few lattice models of dynamic fracture, that have been developed over the past few years to address this problem, have provided definitive answer to this classical problem. These models are either probabilistic, in the sense that a bond of the lattice is broken with a certain probability, or are fully deterministic and therefore, in this sense, are similar to what was described earlier in this chapter.

Why can such lattice models of dynamic fracture be useful? Aside from our experience with lattice models of quasi-steady fracture which provided significant insight into such phenomena, and in contrast to continuum models of cohesive zone, where the correct starting equations are not yet known with certainty, and instabilities that are in qualitative agreement with experiments are difficult to predict, calculations in a lattice or a crystal provide a framework within which the starting point is unambiguous, and instabilities resembling those seen in experiment arise naturally. Here, we describe and discuss some of the most important theoretical results relating to the instability predicted by such models. If the lattice contains no disorder, then it is possible to study the motion of a crack in a macroscopic sample, but not describing the motion of every atom in detail. In this way, questions about the behavior of cohesive zone and the precise nature of crack motion can be resolved without any additional assumptions. A surprising fact is that it is possible to obtain a large variety of analytical results for fracture in arbitrarily large systems. Furthermore, the qualitative lessons following from these calculations also seem to be quite general.

One may wonder about the appropriateness of such models for the dynamic behavior of a crack in a crystal. The critical question, which could have also been asked about the lattice models of quasi-static fracture is, is the lattice essential, or can one make the underlying lattice go away by taking a continuum limit? To our knowledge, all attempts for describing the cohesive zone of brittle materials in a continuum framework have run into severe difficulties (see, for example, the discussion by Langer and Lobkovsky, 1998). These problems do not arise in lattice models described here, nor are they encountered in the atomic-scale MD simulations that will be described in Chapter 9.

The simplicity of the ideal brittle crystal is somewhat misleading in a number of respects. Therefore, before introducing the models, a few natural questions regarding the generality of their predictions are posed and commented on.

(1) Does the simple force law employed between the nodes of the lattice (the atoms in the crystal) neglect any essential aspect of the dynamics? Experience with the MD simulations of dynamic fracture, to be described in Chapter 9, indicates that the same qualitative results are also observed in a brittle ideal crystal (that is, one without disorder). Moreover, we already know that, the same type of force law, when utilized in the quasi-static case, produced very insightful results.

(2) The most interesting calculations that have been carried out so far (see below) are for a strip geometry. Is this geometry too restrictive? The general formulation of fracture mechanics provides an answer to this question, since it tells us that as long as the conditions of small-scale yielding are satisfied, the behavior of a crack is entirely governed by the structure of the stress fields in the near vicinity of the crack tip. These fields are solely determined by the flux of energy to the crack tip. A given energy flux can be provided by an infinite number of different loading configurations, but the resulting dynamics of the crack will all be the same. As a result, the specific geometry used to load the system is irrelevant and no generality is lost by the use of a specific loading configuration. This fact is born out not only by the experimental work described in Chapters 6 and 7, but also by the lattice models of quasi-static fracture described earlier in this chapter.

(3) Are the predictions for a perfect (without disorder) crystal or lattice relevant to amorphous materials? This is still an open question. However, the results of the lattice calculations appear to be remarkably robust. Adding weak quenched disorder to the crystal has little qualitative effect (recall that the same conclusion was reached with the lattice models of quasi-static fracture). The effects of topological disorder are not known. However, it is well-known (see, for example, Holland and Marder, 1998) that when the temperature of a brittle crystal increases above zero, the velocity gap (see below) becomes narrower, and its behavior is reminiscent of what is seen in experiments performed on amorphous materials, which were described in Chapters 6 and 7.

One of the earliest *exact* calculation with lattice models of dynamic fracture was carried out by Kulakhmetova *et al.* (1984) who showed that it is possible to

find exact analytical expressions for Mode I fractures moving in a square lattice. They derived exact relationships between the energy flux to a crack tip and its velocity, observed that phonons must be emitted by moving cracks, and calculated their frequencies and amplitudes. Later calculations by Marder and Gross (1995) extended these results to other lattice geometries, allowed for a general Poisson's ratio, showed that there is a minimum allowed crack velocity (which was found when steady fracture motion was linearly stable), computed the point at which steady motion becomes unstable to a branching instability, and estimated the spacing between branches. These calculations are extremely elaborate, with the analytical expressions being so lengthy that they do not even fit on printed pages (see Gross, 1995)! This is particularly true about the calculations for Mode I equations. For this reason, we summarize some of the results for Mode I, and then proceed to describe in detail how the calculations are carried out in the case of anti-plane shear, Mode III, where the algebra is much less demanding, but most of the ideas are the same. But, let us first describe a relatively simple lattice model of dynamic fracture which could predict some aspects of dynamic fracture.

One of the earliest numerical simulation of dynamic fracture, based on a lattice model, was carried out by Mori *et al.* (1991). A triangular lattice was used in their work, each bond of which was assumed to be a Hooke's spring, if the threshold for its breaking was not exceeded. They assumed that a spring breaks irreversibly if it is stretched beyond a given threshold. Each spring was also characterized by a spring constant $\alpha$. The nodes of the lattice were occupied by particles of mass $m$. Suppose that $\mathbf{e}$ is the elongation vector of the springs that are connected to a particle at position $\mathbf{R}_i$. The equation of motion for the particle at time $t$ is

$$m\frac{d^2\mathbf{R}_i}{dt^2} = -\mathcal{D}\frac{d\mathbf{R}_i}{dt} - k\mathbf{e}, \tag{106}$$

where the first term on the right-hand side is a damping term with $\mathcal{D}$ being the damping constant; this term essentially represents some type of friction. $\mathcal{D}$ is not a parameter that can be measured easily in any experiment, and thus should be treated as a free parameter of the model. Setting $d\mathbf{R}_i/dt = \mathbf{v}_i$ (where $\mathbf{v}_i$ is the velocity of node $i$), we obtain two equations that govern $\mathbf{R}_i$ and $\mathbf{v}_i$ which, when written in a finite-difference form, are given by

$$\mathbf{v}_i(t + \Delta t) = \left(1 - \Delta t\frac{\mathcal{D}}{m}\right)\mathbf{v}_i(t) - \Delta t\, k\mathbf{e}, \tag{107}$$

$$\mathbf{R}_i(t + \Delta t) = \mathbf{R}_i(t) + \frac{\Delta t}{m}\mathbf{v}_i(t), \tag{108}$$

where $\Delta t$ is the time increment.

To begin the simulations, an initial microcrack is inserted in the middle of the lattice. As the initial condition, the lattice is stretched by an amount $\ell$. Equations (107) and (108) are then solved at time $t$ and the springs are examined to identify those that have exceeded their threshold. Such springs are broken, the time $t$ is advanced by $\Delta t$, Eqs. (107) and (108) are solved again, the next springs(s) is (are) broken, and so on. An important parameter of the model is $I = \ell_0/\ell$, where $\ell_0$
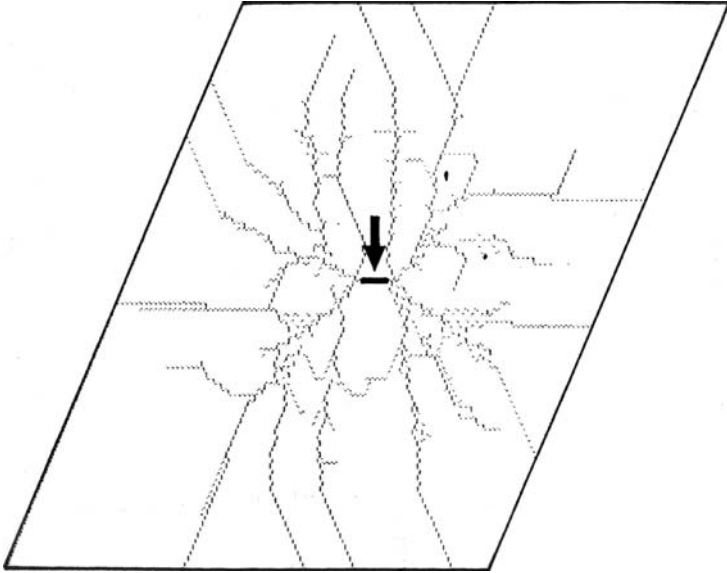
FIGURE 8.20. Fracture pattern in the dynamic model of Mori *et al.* (1991) for $I = 1.0083$. The arrow indicates the location of the initial microcrack.

is the initial length of the springs, and $\ell$ is the initial amount of stretching that the lattice has suffered. Increasing $I$ is, in some sense, equivalent to increasing the temperature of the system. Complex fracture pattern can emerge, depending on the value of $I$. It was found that for $I < 1.0085$ the fracture pattern was tree-like, and in fact no microcrack was formed if no initial crack was inserted in the system. However, for $I > 1.0085$ the microcracks became connected and formed a network and, moreover, even with no initial microcrack in the system, fractures were formed "spontaneously." Figure 8.20 shows the typical fracture pattern for $I = 1.0083$.

In what follows, we describe the exact calculations carried out by Marder and co-worker. Some of our discussions follow closely that of Fineberg and Marder (1999), while the rest are based on the review by Sahimi (1998).

### 8.3.1  Dynamic Fracture in Mode I

The model consists of a triangular lattice of atoms with a lattice constant $a$. Let $\mathbf{u}_i$ be the displacement of a mass point from its equilibrium location, and assume that the energy of the material is a sum of two-body terms, i.e., those that depend on two particles at a time. We then linearize the energy to lowest order in particle displacements. Translational invariance of the lattice implies that the force between particles 0 and 1 depends only upon $\mathbf{U}_1 = \mathbf{u}_1 - \mathbf{u}_0$. However, the force can be a general linear functional of $\mathbf{U}_1$. One way of writing down an expression for such a general linear functional is to decompose the force between particles 0 and 1

into a component along $\hat{\mathbf{d}}_{\|1}$, which is along the line that connects two neighboring atoms, and a component that is along $\hat{\mathbf{d}}_{\perp 1}$ which is perpendicular to it. Hence, the first component corresponds to central forces between atoms, whereas the second component is a non-central force. Recall from Chapter 8 of Volume I (see also Chapter 9) that in real materials non-central forces between atoms are the rule rather than the exception. Suppose now that the restoring force parallel to the direction of equilibrium bonds is proportional to $\mathbf{F}_{\|}$, while that perpendicular to this direction is proportional to $\mathbf{F}_{\perp}$. Therefore, if $\mathbf{U}_1 = (U_{1x}, U_{1y})$, then the force due to the displacement of the particle along $\mathbf{U}_1 = \mathbf{u}_{i-1,j+1} - \mathbf{u}_{i,j}$ is given by

$$F_{\|}\hat{\mathbf{d}}_{\|1}(\mathbf{U}_1 \cdot \hat{\mathbf{d}}_{\|1}) + F_{\perp}\hat{\mathbf{d}}_{\perp 1}(\mathbf{U}_1 \cdot \hat{\mathbf{d}}_{\perp 1}) = F_{\perp}\left(\frac{-1}{2}, \frac{\sqrt{3}}{2}\right)\left(\frac{-1}{2}U_{1x}, \frac{\sqrt{3}}{2}U_{1y}\right)$$

$$+F_{\perp}\left(\frac{\sqrt{3}}{2}, \frac{1}{2}\right)\left(\frac{\sqrt{3}}{2}U_{1x}, \frac{1}{2}U_{1y}\right). \tag{109}$$

If, by a similar method, one adds up contributions from other particles, one obtains for the force due to neighbors

$$\mathbf{F}(m, n) = \sum_{j=1}^{6} \sum_{q=\|,\perp} F_q\hat{\mathbf{d}}_{qj}[\mathbf{U}_j(m, n) \cdot \hat{\mathbf{d}}_{qj}] \tag{110}$$

By varying the constants $F_{\|}$ and $F_{\perp}$, one can obtain any desired values of shear and longitudinal wave speeds, which are given by

$$c_l^2 = \frac{3a^2}{8m}(F_{\perp} + 3F_{\|}), \tag{111}$$

$$c_t^2 = \frac{3a^2}{8m}(3F_{\perp} + F_{\|}), \tag{112}$$

where $m$ is the mass of each particle. In addition to the forces between the neighboring atoms, it is possible to take into account the effect of complex dissipative functions that depend upon particle velocities. Marder and Gross (1995) added a term to the equations so as to reproduce the experimentally-measured frequency dependence of sound attenuation in Plexiglas. There is a slight technical restriction in the calculations of Marder and Gross (1995) in that, forces right on the crack line are required to be central. However, more detailed calculations indicate that, when this technical limitation is removed, the results do not change significantly. Figure 8.21 shows schematics of steady-state propagation of a crack in the lattice loaded in Mode I.

Detailed, lengthy and difficult calculations show that, in Mode I loading, many details of the relation between loading and fracture velocity depend upon ratios of the sound speeds and also the frequency dependence of dissipation, implying that there exists no universal curve that can describe Mode I fracture. When only central forces are present ($F_{\perp} = 0$), it is difficult (but not impossible) to have steady-state fracture propagation, implying that the range of loads for which cracks can propa-
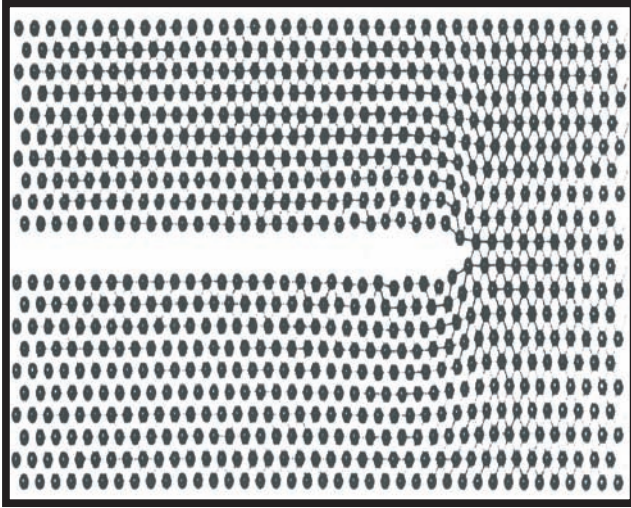
FIGURE 8.21. Schematics of steady-state motion of fracture moving in an ideal brittle crystalline material under Mode I loading (after Fineberg and Marder, 1999).

gate in a stable fashion is small, and depends upon the amount of dissipation. Thus, the presence of non-central forces is essential to having stable fracture growth. This also has important implications for MD simulations of dynamic fracture propagation in materials. For example, the proper form of the interatomic potential for representing the non-central forces, such as the bond-bending ones, is an important issue in the MD simulations which will be discussed in detail in Chapter 9.

Given the complexities of exact calculations for Mode I fracture, it is wiser to consider a simpler geometry for which the analytical methods can be discussed in complete detail, and all the essential ideas that are needed for understanding Mode I can be explored with much less elaborate calculations. This is the subject of the next section.

## 8.3.2 *Dynamic Fracture in Mode III*

The main results described in this section are for the relation between loading and fracture velocity, the velocity gap, and the calculation of the point at which steady crack motion becomes unstable. Calculations of Marder and Gross (1995) indicate that steady-state fractures, when they exist, are always linearly stable, and therefore their stability analysis will not be discussed.

We consider a fracture moving in a strip composed of $2(N + 1)$ rows of mass points. Figure 8.22 presents the schematics of the system. The bonds are linear springs until they break at a separation of $2u_f$. The location of each mass point is described by $u(m, n)$, interpreted as the height of mass point $(m, n)$ into or out of the page, where $m$ takes on integer values, while $n$ takes values of the form $\frac{1}{2}, \frac{3}{2}, \cdots, N + \frac{1}{2}$. The force between adjacent mass points is determined by the
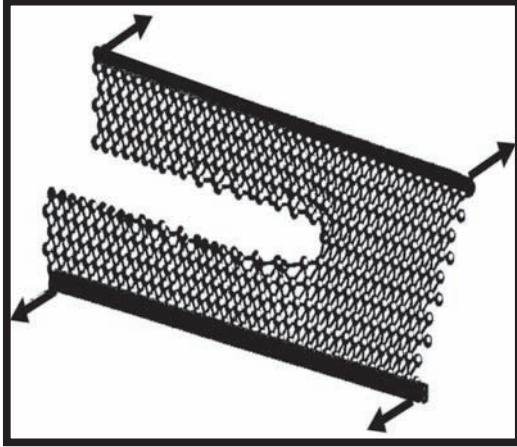
FIGURE 8.22. Dynamic fracture of triangular crystal in anti-plane shear (after Fineberg and Marder, 1999).

difference in the height between them. Hence, in terms of $u(m, n)$, the equation of motion of the system is given by

$$\frac{d^2 u}{dt^2} = -b \frac{du}{dt} + \frac{1}{2} \sum_{m', n'} \mathcal{F}[u(m', n') - u(m, n)], \qquad (113)$$

where $b$ is the coefficient of a small dissipative term. However, other functional forms for the dissipation term can be used. The sum in Eq. (110) is over all nearest neighbors $(m', n')$ of $(m, n)$. Here,

$$\mathcal{F}(u) = u \Theta (2u_f - | u |), \qquad (114)$$

describes ideal brittle springs, with $\Theta$ being the step function. The boundary condition is given by

$$u[m, \pm(N + \frac{1}{2})] = \pm u_N. \qquad (115)$$

It is important to determine the value of $u_N$ for which there is just enough stored energy per unit length to the right of the fracture to break the pair of bonds connected to each lattice site on the crack line. For $m \gg 0$ one has

$$u(m, n) = \frac{n u_N}{N + \frac{1}{2}}, \qquad (116)$$

and the height difference between mass points with adjacent values of $n$ is

$$U_{\text{right}} = \frac{n_N}{N + \frac{1}{2}}. \qquad (117)$$

Therefore, the energy stored per unit length in the $(2N + 1)$ rows of bonds is

$$\frac{1}{2} \times (2 \text{ upper bonds/site}) \times (\text{rows}) \times (\text{spring constant}) \times U_{\text{right}}^2$$

$$= \frac{1}{2} \times 2(2N + l) \times \frac{1}{2} \times \left(\frac{u_N}{N + \frac{1}{2}}\right)^2 = 2N_0 u_N^2, \tag{118}$$

where $N_0 = 1/(2N + 1)$. The energy required to break two bonds each time the crack propagates by a unit length is given by

$$\frac{1}{2} \times (2 \text{ bonds/site}) \times (\text{spring constant}) \times (\text{separation at fracture})^2$$

$$= \frac{1}{2} \times 2 \times \frac{1}{2} \times (2u_f)^2 = 2u_f^2. \tag{119}$$

Therefore, by equating Eqs. (118) and (119) one obtains the proper dimensionless measure of the external driving force,

$$\Delta = \frac{u_N \sqrt{N_0}}{u_f}, \tag{120}$$

which reaches 1 as soon as there is enough energy to the right of the crack to break the bonds along the crack line. Note that $\Delta$ is linearly related to the displacement $u_N$ imposed at the edges of the strip.

Slepyan (1981) and Kulakhmetova *et al.* (1984) developed the techniques for solving problems of this type. However, there are some differences between details of their solution and what is presented here, because Eq. (113) describes motion in a strip, which is what is considered here, rather than an infinite plate that they studied. Both geometries have certain advantages. The strip is preferable to the infinite plate if one wishes to compare the predictions with the results of numerical simulations. On the other hand, using the infinite plate results in certain natural limits. Moreover, the lattice considered here is a triangular rather than a square lattice used by Slepyan and co-workers.

We assume that a fracture moves in steady state, so that one by one, the bonds connecting $u(m, \frac{1}{2})$ with $u(m + 1, -\frac{1}{2})$ or $u(m, -\frac{1}{2})$ break because, as a result of the driving force described by Eq. (115), the distance between these atoms exceeds the limit set by Eq. (115). Assuming that the times at which the bonds break are known, the original nonlinear problem is immediately transformed into a linear problem. However, once the solution of the linear problem has been derived, one must verify that,

(1) bonds break at the time that they are supposed to. If we impose this condition, we obtain a relation between the crack velocity $v$ and dimensionless loading $\Delta$.
(2) Conversely, no bonds break when they are not supposed to. This condition results in a velocity gap on the low end of the velocity range, and leads to crack tip instabilities above a critical energy flux (see below).

We note here that steady states in a perfect lattice or crystal are more complex than those in a continuum, since in the latter case a steady state acts as $u(x, vt)$,

whereas the closest to such a state that one can come in a triangular lattice is by having the symmetry

$$u(m, n, t) = u(m + 1, n, t + 1/v), \tag{121}$$

and also

$$u(m, n, t) = -u\left[m, -n, t - \left(\frac{1}{2} - g_n\right)/v\right], \tag{122}$$

which implies in particular that

$$u\left(m, \frac{1}{2}, t\right) = -u\left[m, -\frac{1}{2}, t - 1/(2v)\right]. \tag{123}$$

Here

$$g_n = \begin{cases} 0 & \text{if } n = \frac{1}{2}, \frac{5}{2} \cdots \\ 1 & \text{if } n = \frac{3}{2}, \frac{7}{2} \cdots \\ \text{mod}(n - \frac{1}{2}, 2) & \text{in general.} \end{cases} \tag{124}$$

By assuming that a fracture is in steady state, one can eliminate the variable $m$ entirely from the equation of motion. We define

$$u_n(t) = u(0, n, t), \tag{125}$$

and write the equations of motion in steady state as

$$\frac{d^2 u_n}{dt^2} = \frac{1}{2} \begin{bmatrix} u_{n+1}[t - (g_{n+1} - 1)/v] & +u_{n+1}(t - g_{n+1}/v) \\ +u_n(t + 1/v) & -6u_n(t) + u_n(t - 1/v) \\ +u_{n-1}[t - (g_{n-1} - 1)/v] & +u_{n-1}(t - g_{n-1}/v) \end{bmatrix} - b\frac{du_n}{dt} \tag{126}$$

if $n > 1/2$. For $n = 1/2$ we have

$$\frac{d^2 u_{1/2}}{dt^2} =$$

$$\frac{1}{2} \begin{bmatrix} u_{3/2}(t) & +u_{3/2}(t - 1/v) \\ +u_{1/2}(t + 1/v) - 4u_{1/2}(t) & +u_{1/2}(t - 1/v) \\ +[u_{-1/2}(t) - u_{1/2}(t)]\Theta(-t) & +[u_{-1/2}(t - 1/v) - u_{1/2}(t)]\Theta[1/(2v) - t] \end{bmatrix} - b\frac{du_{1/2}}{dt} \tag{127}$$

We have assumed that $t = 0$ is the time at which the bond between $u(0, \frac{1}{2}, t)$ and $u(0, -\frac{1}{2}, t)$ breaks, and therefore, by symmetry, the time at which the bond between $u(0, \frac{1}{2}, t)$ and $u(1, -\frac{1}{2}, t)$ breaks is $1/(2v)$.

Above the fracture line, Eqs. (126) and (127) are completely linear, and thus it is not difficult to derive the solution that yields the motion of every atom with $n > 1/2$ in terms of the behavior of an atom with $n = 1/2$. If we Fourier transform

Eq. (126), we obtain

$$
-\omega^2 u_n(\omega) =
\left\{
\begin{array}{l}
\frac{1}{2} u_{n+1}(\omega)[e^{i\omega(g_{n+1}-1)/v} + e^{i\omega g_{n+1}/v}] \\[4pt]
+\frac{1}{2} u_n(\omega)(e^{i\omega/v} - 6 + e^{-i\omega/v}) \\[4pt]
+\frac{1}{2} u_{n-1}(\omega)[e^{i\omega(g_{n-1}-1)/v} + e^{i\omega g_{n-1}/v}]
\end{array}
\right\} + ib\omega
\tag{128}
$$

where $\omega$ is the Fourier transform variable conjugate to the time $t$. If

$$
u_n(\omega) = u_{1/2}(\omega) e^{k(n-1/2) - i\omega g_n/(2v)}
\tag{129}
$$

then by substituting Eq. (129) into Eq. (128), and using the fact that $g_n + g_{n+1} = 1$, we obtain

$$
-\omega^2 u_{1/2}(\omega) =
$$

$$
\left\{
\begin{array}{l}
\frac{1}{2} u_{1/2}(\omega) e^{k}[e^{i\omega(g_{n+1}+g_{n-2})/(2v)} + e^{i\omega(g_{n+1}+g_n)/(2v)}] \\[4pt]
+\frac{1}{2} u_{1/2}(\omega)(e^{i\omega/v} - 6 + e^{-i\omega/v}) \\[4pt]
+\frac{1}{2} u_{1/2}(\omega) e^{-k}[e^{i\omega(g_{n-1}+g_{n-2})/(2v)} + e^{i\omega(g_{n-1}+g_n)/(2v)}]
\end{array}
\right\} + ib\omega \, u_{1/2}(\omega)
\tag{130}
$$

so that

$$
\omega^2 + ib\omega + 2\cosh(k)\cos[\omega/(2v)] + \cos(\omega/v) - 3 = 0.
\tag{131}
$$

Let

$$
z = \frac{3 - \cos(\omega/v) - \omega^2 - ib\omega}{2\cos(\omega/2v)},
\tag{132}
$$

which is equivalent to

$$
y = z + \sqrt{z^2 - 1},
\tag{133}
$$

with $y = e^k$.

One now constructs a solution which meets all the boundary conditions by writing

$$
u_n(\omega) = u_{1/2}(\omega) e^{-i\omega g_n/2v}
\left\{
\frac{y^{[N+1/2-n]} - y^{-[N+1/2-n]}}{y^N - y^{-N}}
\right\}
+ \frac{u_N(n - \frac{1}{2})}{N} \frac{2\alpha}{\alpha^2 + \omega^2}
\tag{134}
$$

This solution equals $u_{1/2}$ for $n = 1/2$, and equals $2\alpha u_N/(\alpha^2 + \omega^2)$ for $n = N + \frac{1}{2}$. Introducing $\alpha$ is necessitated by the fact that for $n = N + \frac{1}{2}$, $u(m, n, t) = u_N$. Instead of working with the Fourier transform of this boundary condition, which is a delta function and difficult to work with, it is better to use the following boundary condition

$$
u_{N+1/2}(t) = u_N e^{-\alpha|t|},
\tag{135}
$$

and let $\alpha \to 0$ at the end of the calculation. In the analysis that follows, frequent use is made of the fact that $\alpha$ is small.

From a physical point of view, the most interesting variable is not $u_{1/2}$, but the distance between the bonds which will actually break. Hence, we define

$$U(t) = \frac{u_{1/2}(t) - u_{-1/2}(t)}{2} = \frac{u_{1/2}(t) + u_{1/2}(t + 1/2v)}{2}, \tag{136}$$

to rewrite Eq. (127) as

$$\frac{d^2 u_{1/2}}{dt^2} = \frac{1}{2} \begin{bmatrix} u_{3/2}(t) & +u_{3/2}(t - 1/v) \\ +u_{1/2}(t + 1/v) - 4 & +u_{1/2}(t) + u_{1/2}(t - 1/v) \\ -2U(t)\Theta(-t) & -2U(t - 1/2v)\Theta[1/(2v) - t] \end{bmatrix} - b\frac{du_{1/2}}{dt}. \tag{137}$$

If we Fourier transform this expression, use Eq. (134), and define

$$U^{\pm}(\omega) = \int_{-\infty}^{\infty} U(t)\Theta(\pm t)\exp(i\omega t)d\omega, \tag{138}$$

we obtain

$$u_{1/2}(\omega)F(\omega) - (1 + e^{i\omega/2v})U^-(\omega) = -\frac{u_N}{N}\frac{2\alpha}{\omega^2 + \alpha^2}, \tag{139}$$

where

$$F(\omega) = \left\{ \frac{y^{[N-1]} - y^{-[N-1]}}{y^N - y^{-N}} - 2z \right\} \cos(\omega/2v) + 1. \tag{140}$$

We now use Eq. (136) in the form

$$U(\omega) = \frac{1}{2}(1 + e^{-i\omega/2v})u_{1/2}(\omega) \tag{141}$$

to obtain

$$U(\omega)F(\omega) - 2\cos^2(\omega/4v)U^-(\omega) = -\frac{u_N}{N}\frac{2\alpha}{\omega^2 + \alpha^2}. \tag{142}$$

If we write

$$U(\omega) = U^+(\omega) + U^-(\omega) \tag{143}$$

we finally obtain

$$U^+(\omega)Q(\omega) + U^-(\omega) = u_N N_0 \left( \frac{1}{\alpha + i\omega} + \frac{1}{\alpha - i\omega} \right), \tag{144}$$

where

$$Q(\omega) = \frac{F(\omega)}{F(\omega) - 2\cos^2(\omega/4v)}. \tag{145}$$

To derive the right-hand side of Eq. (144), we used the facts that $F(0) = -1/N$, and that $\alpha$ is very small, so that the right-hand side of this equation is a delta function.

We now utilize the Wiener–Hopf technique (see, for example, Noble, 1958) to write

$$Q(\omega) = \frac{Q^-(\omega)}{Q^+(\omega)}, \tag{146}$$

where $Q^-$ and $Q^+$ are free of poles and zeroes in the lower and upper complex $\omega$ planes, respectively. This decomposition can be carried out with the explicit formula

$$Q^\pm(\omega) = \exp\left[\lim_{\varepsilon \to 0} \frac{1}{2\pi} \int \frac{\ln Q(\omega')}{i\omega \mp \varepsilon - i\omega'}\, d\omega'\right]. \tag{147}$$

We now separate Eq. (144) into two parts, one of which has poles only in the lower half plane, while the other part has poles only in the upper half plane:

$$\frac{U^+(\omega)}{Q^+(\omega)} - \frac{u_N N_0}{Q^-(0)} \frac{1}{-i\omega + \alpha} = \frac{u_N N_0}{Q^-(0)} \frac{1}{i\omega + \alpha} - \frac{U^-(\omega)}{Q^-(\omega)}. \tag{148}$$

Because the right- and left-hand sides of Eq. (148) have poles in opposite sections of the complex plane, they must separately equal a constant. However, the constant must be zero, because otherwise $U^-$ and $U^+$ will behave as a delta function near $t = 0$, and therefore

$$U^-(\omega) = \frac{u_N N_0 Q^-(\omega)}{Q^-(0)(\alpha + i\omega)}, \tag{149}$$

and

$$U^+(\omega) = \frac{u_N N_0 Q^+(\omega)}{Q^-(0)(\alpha - i\omega)}, \tag{150}$$

which provide an explicit solution for $U(\omega)$. Numerical evaluation of Eq. (147), and $U(t)$ using Eqs. (149) and (150) is fairly straightforward (using, for example, fast Fourier transforms). However, in carrying out the numerical transforms, one must carefully analyze the behavior of the functions for large values of $\omega$. If these functions decay as $1/(i\omega)$ [the inverse Fourier transform of which is $\Theta(t)$], we should subtract the $1/(i\omega)$ off before the numerical transform is performed, after which we should add this term back with the appropriate step function. Conversely, if the functions to be Fourier transformed have $\Theta(t)$ discontinuity, it is best to subtract off the appropriate multiple of $e^t \Theta(t)$ before Fourier transforming, and then add on the appropriate multiple of $1/(1 - i\omega)$. A solution of Eqs. (149) and (150) constructed in this manner is given in Figure 8.23.

One can now derive a relation between the dimensionless displacement $\Delta$ and the crack velocity $v$. Recall that making the transition from the original nonlinear problem posed by Eq. (113) to the linear problem expressed by Eq. (126) relies on assuming that bonds along the crack line break at time intervals of $1/2v$. Because of the symmetries expressed by Eqs. (121)–(123), it is sufficient to guarantee that

$$u(t = 0) = u_f. \tag{151}$$

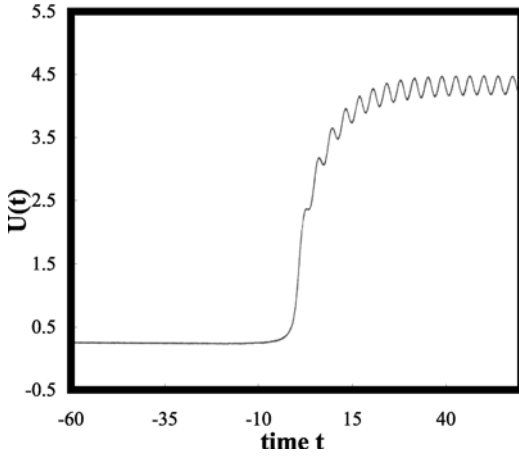FIGURE 8.23. Predictions of Eqs. (149) and (150) for $v = 0.5$, $N = 9$ and $b = 0.01$ (after Fineberg and Marder, 1999).

All the displacements are proportional to the boundary displacement $u_N$, and hence Eq. (151) fixes a unique value of $u_N$ and its dimensionless counterpart, $\Delta$. This means that, once one assumes that the fracture moves in steady state at a velocity $v$, there is a unique $\Delta$ to make it possible.

To derive Eq. (151), we must require that

$$\lim_{t \to 0^-} \frac{1}{2\pi} \int U^-(\omega) \exp(-i\omega t) d\omega = u_f. \tag{152}$$

This integral can be evaluated by inspection. We know that for $t > 0$,

$$\int U^-(\omega) \exp(-i\omega t) \, d\omega = 0, \tag{153}$$

and that any function that for large $\omega$ behaves as $1/(i\omega)$ has a step function discontinuity at the origin. Therefore, Eqs. (149) and (152) become

$$u_f = u_N N_0 \frac{Q^-(\infty)}{Q^-(0)}. \tag{154}$$

Since Eq. (145) yields $Q(\infty) = 1$, one obtains from Eq. (147) that,

$$Q^-(\infty) = Q^+(\infty) = 1. \tag{155}$$

As a result, Eq. (154) and the definition of $\Delta$ [Eq. (120)] yield

$$\Delta = \frac{Q^-(0)}{\sqrt{N_0}}. \tag{156}$$

To make this result more explicit, we use Eq. (147) and the fact that

$$
Q^-(0) = \exp\left\{\frac{1}{2\pi}\int\frac{1}{2}\left[\frac{\ln Q(\omega')}{\varepsilon - i\omega'} + \frac{\ln Q(-\omega')}{\varepsilon + i\omega'}\right]d\omega'\right\}
$$
$$
= \exp\left\{\frac{1}{2\pi}\int\left[\frac{1}{-2i\omega'}\ln\frac{Q(\omega')}{\bar{Q}(\omega')} + \frac{\varepsilon}{\varepsilon^2 + \omega'^2}\ln Q(0)\right]d\omega'\right\}.
$$

(157)

Therefore,

$$
Q^-(0) = \sqrt{N_0}\exp\left[-\frac{1}{2\pi}\int\frac{1}{2i\omega'}\ln\frac{Q(\omega')}{\bar{Q}(\omega')}d\omega'\right].
$$

(158)

Inserting Eq. (158) into Eq. (156) yields

$$
\Delta = \exp\left[-\frac{1}{2\pi}\int\frac{1}{2i\omega'}\ln\frac{Q(\omega')}{\bar{Q}(\omega')}d\omega'\right].
$$

(159)

In order to obtain an expression that is correct not only for Mode III model considered here, but also for more general cases, we rewrite Eq. (159) as

$$
\Delta = C\exp\left[-\frac{1}{2\pi}\int\frac{1}{2i\omega'}\{\ln Q(\omega') - \overline{\ln Q(\omega')}\}d\omega'\right],
$$

(160)

where $C$ is a constant of order unity that is determined by the geometry of the lattice. For example, for the triangular lattice loaded in Mode III, $C = 1$, whereas $C = 2/\sqrt{3}$ for the same lattice loaded in Mode I (Marder and Gross, 1995). The advantage of Eq. (160) is that it is suitable for numerical evaluation, since there is no uncertainty regarding the phase of the logarithm.

When $b$, the coefficient of the dissipation term, becomes sufficiently small, $Q$ is real for real $\omega$ except in the small neighborhood of the isolated roots and poles that are near the real $\omega$-axis. Suppose that $r_i^+$ are the roots of $Q$ with negative imaginary part (since they belong to $Q^+$), $r_i^-$ are the roots of $Q$ with positive imaginary part, and $p_i^\pm$ are the poles of $Q$. Equation (160) can be written as

$$
\Delta = C\sqrt{\prod_i r_i^- p_i^+ / r_i^+ p_i^-}, \quad \text{for } b \to 0.
$$

(161)

In deriving Eq. (161) use was made of the fact that, away from a root or pole of $Q$, the integrand of Eq. (160) vanishes. Together with Eqs. (149) and Eqs. (150), Eqs. (160) and (161) constitute the formal solution of the model. Since $Q$ is a function of the steady-state velocity $v$, Eq. (159) relates the external driving force imposed on the crystal, $\Delta$, to the velocity $v$ of the crack. The results of a typical calculation for $N = 9$ are presented in Figure 8.24, which show clearly the velocity gap.

Having determined the formal solution of the model, we can now investigate a few important issues.

### 8.3.2.1  Phonon Emission

At $\Delta = 1$ just enough energy is stored to the right of the fracture tip to break all bonds along the crack line. However, since all steady states occur for $\Delta > 1$, not
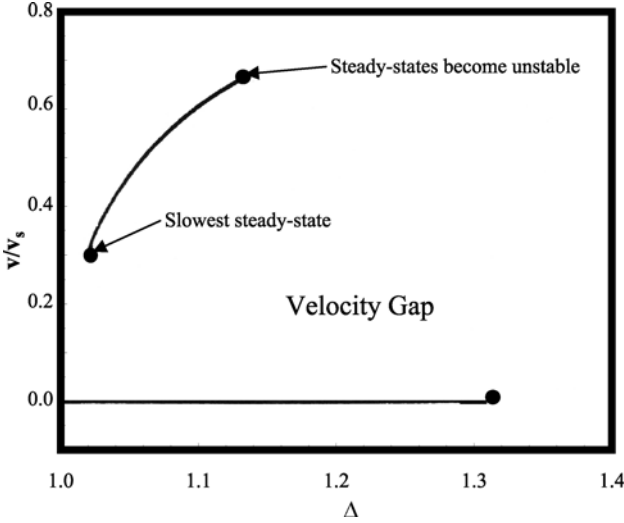
FIGURE 8.24. Fracture velocity $v$, normalized by sound velocity $v_s = \sqrt{3}/2$, versus the driving force $\Delta$, using $N = 9$. The end points of the lower curve indicate the linearly-stable lattice-trapped states (after Fineberg and Marder, 1999).

all the energy stored to the right of the fracture tip will be used for breaking bonds. What happens to the remaining energy depends on the amount of dissipation $b$, and the distance from the fracture tip at which one inspects the system. In the limit $b \to 0$, travelling waves leave the fracture tip and carry energy off in its wake, with the amount of energy that they carry off being independent of $b$. For all non-zero values of $b$, however, the travelling waves will eventually decay, and the extra energy will have been absorbed by dissipation, but the value of $b$ determines whether one views the process as microscopic or macroscopic.

The frequencies of the radiation emitted by the crack have a simple physical interpretation. Consider the motion of a particle through a lattice, in which the phonons are described by the dispersion relation, $\omega_\alpha(\mathbf{k})$. If the particle moves with a constant velocity $\mathbf{v}$ and interacts with the various ions according to some function $\mathcal{F}$, then, to linear order, the motions of the ions can be described by a matrix $\mathbf{M}$ which describes their interactions with each other as

$$m\frac{d^2 u_i^{(l)}}{dt^2} = -\sum_{j,n} M_{ij}(\mathbf{R}^{(l)} - \mathbf{R}^{(n)})u_j^{(n)} + \sum_n \mathcal{F}_i(\mathbf{R}^{(n)} - \mathbf{v}t). \qquad (162)$$

Multiplying both sides of Eq. (162) by $e^{i\mathbf{k}\cdot\mathbf{R}^{(l)}}$, summing over $l$, and letting $\mathbf{K}$ and $\Omega$ be, respectively, the reciprocal lattice vectors and the volume of a unit cell, yield

$$m\frac{d^2 u_i(\mathbf{k})}{dt^2} = \sum_j M_{ij}(\mathbf{k})u_j(\mathbf{k}) + \frac{1}{\Omega}\sum_{\mathbf{K}} \exp[i\mathbf{v}\cdot(\mathbf{k}+\mathbf{K})t]\mathcal{F}_i(\mathbf{k}+\mathbf{K}). \qquad (163)$$

By inspection we can show that the lattice frequencies excited in this way are those which in the extended zone scheme satisfy the following equation (Ashcroft and Mermin, 1976)

$$\omega(\mathbf{k}) = \mathbf{v} \cdot \mathbf{k}. \tag{164}$$

If we pretend that the crack is a particle, we can use Eq. (164) to predict the phonons that the crack emits.

There are actually two phonon dispersion relations to consider, one for propagating radiation far behind the crack tip, and the other for propagating radiation far ahead of the tip. Far behind the crack tip, all the bonds are broken, and therefore to find the travelling one we must set $U^- = 0$ in Eq. (142), since all the bonds behind the tip are broken, and also $u_N = 0$, because phonons can propagate without any driving term, which then lead us to $F(\omega) = 0$. Similarly, because no bond far ahead of the crack tip is broken, $U^- = U$, and the condition for phonons is $F(\omega) - 2\cos^2(\omega/4v) = 0$. Therefore, according to Eq. (145), the roots and poles of $Q(\omega)$ are the phonon frequencies behind and ahead of the fracture, which are also the quantities that appear in Eq. (161).

### 8.3.2.2    Forbidden Fracture Velocities

After verifying that bonds along the crack line break when they should, we must also verify that they have not been stretched enough to break earlier than they should have. That is, not only must the bond between $u_{0+}$ and $u_{0-}$ reach length $2u_f$ at $t = 0$, but also must be the first time at which that bond stretches to a length greater than $2u_f$. For $0 < v < c$ with $c \simeq 0.3$ (the precise value of $c$ varies with $b$ and $N$), this condition is violated. The states have the unphysical character that masses rise above height $u_f$ for $t < 0$, with the bond connecting them to the lower line of masses remaining, however, intact, and then they descend whereupon the bond breaks.

Since the solution of Eqs. (126) and (127) is unique, but does not in this case solve Eq. (113), no solutions of Eq. (113) exist at all at these velocities. Once the crack velocity has dropped below a lower critical value, all the steady states will have this character. This argument indicates that no steady state in the sense of Eq. (121) can exist. It is also possible to search for analytical solutions that are periodic and travel *two* lattice spacings before repeating, but to our knowledge no solutions of this type have yet been found. One can numerically verify that if a fracture is allowed to propagate with $\Delta$ just above the critical threshold, and $\Delta$ is then very slowly lowered through the threshold, the fracture stops propagating. It does not slow down noticeably; rather, suddenly, the moving fracture emits a burst of radiation that carries off its kinetic energy, and stops in the space of an atom, which is why there exists a velocity gap.

### 8.3.2.3    Nonlinear Instabilities

Recall that we assumed in the calculations for predicting the steady states that the only bonds which break are those which lie on the crack path. This assumption

can be tested by using the numerical solutions of Eqs. (149) and (150), since the solution fails above a critical value $\Delta_c$ of $\Delta$. The sound speed $v_s(= c_R)$ equals, in dimensionless units, $\sqrt{3}/2$, and thus we rescale velocities by this value. For example, for $N = 9$, at a velocity, $v_c/v_s = 0.666\cdots$, $\Delta_c = 1.158\cdots$, and the bond between $u(0, \frac{1}{2})$ and $u(1, \frac{1}{2})$ reaches a distance of $2u_f$ shortly after the bond between $u(0, \frac{1}{2})$ and $u(0, -\frac{1}{2})$ breaks. The steady-state solutions that are obtained when the lattice is strained with larger values of $\Delta$ are inconsistent; only dynamical solutions more complex than steady states, involving the breaking of bonds off the fracture path, are possible. To investigate these states, one must numerically solve Eq. (113). Such numerical simulations have been carried out (Marder and Liu, 1993). The results show that just above the threshold at which horizontal bonds begin to break, the distance between these extra broken bonds diverges. The reason is that breaking a horizontal bond takes energy from the fracture and slows it down below the critical value. The fracture then tries once more to reach the steady state, and only in the last stages of the approach does another horizontal bond break, hence beginning the process again.

A rough estimate of the distance between broken horizontal bonds can be obtained as follows. Let $\ell_h(t)$ be the length of an endangered horizontal bond. One must view the problem in a reference frame moving with the fracture tip, and therefore at every time interval $1/v$ one shifts attention to a bond one lattice spacing to the right. When $\Delta$ is only slightly greater than $\Delta_c$, the length of such a bond, viewed in a moving frame, should behave, before it breaks, as

$$\ell_h \sim 2u_f + \frac{\partial \ell_h}{\partial \Delta}(\Delta - \Delta_c) - \delta\ell \exp(-bt). \tag{165}$$

Here, $\partial \ell_h/\partial \Delta$ denotes the rate at which the steady-state length of $\ell_h$ would change with $\Delta$, if this bond were not allowed to break, and $\delta\ell$ describes how much smaller than its steady-state value the bond ends up after the breaking event occurs. Marder and Gross (1995) showed that deviations from steady states die away at long times as $\exp(-bt)$. From Eq. (165) we can estimate the time between breaking events by setting $\ell_h = 2u_f$ and solving the equation for $t$. The result is that the frequency $f$ with which horizontal bonds break should scale above the critical strain $\Delta_c$ as

$$f \sim -\frac{b}{\ln[(\delta\ell)^{-1}(\Delta - \Delta_c)\partial \ell_h/\partial \Delta]}. \tag{166}$$

However, the accuracy of Eq. (166) has proven to be difficult to check.

### 8.3.2.4  The Connection to the Yoffe's Criterion

The basic reason for the branching instability seen above is the lattice analogue of the Yoffe's instability (see Section 7.6.3), but on small scales. The Mode III calculation finds that the critical velocity for the instability to frustrated branching events is indeed close to the value of $0.6c_R$ predicted by Yoffe in the continuum. However, the critical velocity seen experimentally in amorphous materials is about $\frac{1}{3}c_R$, not about $\frac{2}{3}c_R$. This discrepancy can be due to some combination of three factors.

(1) The force law between the atoms is actually much more complex than ideal breaking bonds. Gao (1993) has pointed out that the Rayleigh wave speed $c_R$ in the vicinity of a crack tip may be significantly lower than its value far away from the tip, because material is being stretched beyond the range of validity of linear elasticity.
(2) The experiments are at room temperature, while the calculations are at zero temperature.
(3) The experiments are in amorphous materials, while the calculations are for a crystal or perfect lattice.

Before closing this section, let us summarize the main results obtained from the ideal brittle crystal.

(i) For a range of loads above the Griffith point, a fracture can be trapped by the crystal, i.e., it neither moves nor heals, although it is energetically possible for the fracture to move (Thomson *et al.*, 1971; Thomson, 1986). However, molecular dynamics simulations, to be described in Chapter 9, indicate that lattice trapping occurs at very low temperatures and in fact disappears at room temperature.

(ii) Steady-state fracture motion exists, and is a stable attractor for a range of energy flux.

(iii) Steadily-moving fractures emit phonons with frequencies that can be computed from a simple conservation law.

(iv) The relation between the fracture energy and velocity can be computed.

(v) The slowest steady state runs at around $0.20c_R$. There is no slower-moving steady-state fracture, and therefore there is a velocity gap.

(vi) At an upper critical energy flux, steady-state fractures become unstable and generate frustrated branching events in a fashion reminiscent of experiments in amorphous materials described in Chapter 7.

### 8.3.3   The Effect of Quenched Disorder

The robustness of the branching phenomena described in the perfect crystal model, which contained no types of heterogeneities, has been further illustrated by numerical studies of a few lattice models in which the effect of disorder was explicitly taken into account. We briefly describe some of these models and their most important predictions. However, before doing so, let us point out that there is a qualitative difference between quenched disorder and thermal fluctuations which is quite important to brittle fracture. If a fracture encounters a tough spot in a material, it can stop propagating forever. However, although thermal fluctuations might halt a fracture temporarily, they are just as likely to give energy to a static fracture and induce it to start moving again. Lattice-trapped fractures are not completely static in the presence of thermal fluctuations; they creep ahead with some probability (Marder, 1996). When the rate of creep increases to speeds that are of the order of the sound speed, the distinction between creeping and running fractures vanishes, and therefore the velocity gap disappears.

Furukawa (1993) modified the model proposed by Mori *et al.* (1991) (see Section 8.3) by including a shear friction in the model. Thus, in this model the equation of motion is given by

$$m\frac{d\mathbf{v}_i}{dt} = -\mathcal{D}_1\mathbf{v}_i - \mathcal{D}_2\sum_j(\mathbf{v}_i - \mathbf{v}_j) + \sum_j\mathbf{F}(\mathbf{R}_j - \mathbf{R}_i), \qquad (167)$$

where the second term on the right-hand side represents shear friction or dissipation, while the last term is the force term which was taken to be $\mathbf{F}(\mathbf{R}) = \mathbf{R}f(R)/R$, where $R$ is the magnitude of $\mathbf{R}$. The function $f(R)$ was selected to be $f(R) = R - 1$ for $R \leq R_0$, and $f(R) = (R_0 - 1)\exp[-\kappa(R - R_0)]$ for $R \geq R_0$, where $R_0$ and $\kappa$ are two constants. A bond breaks if $R > R_0$. Both square and triangular lattices were used. Fracture formation was initiated by inserting an initial microcrack at the center of the lattice. In one case the lattice spacing in the direction of the macroscopic deformation was 1, and the remaining lattice distances were $R_e$ with $1 < R_e < R_0$. An interesting prediction of the model was that, when $\kappa = \infty$, the fracture velocity $v = \ell/t$, where $\ell$ is the distance between the central microcrack and the most distant broken bond at time $t$, follows a power law:

$$v \sim \frac{(R_e - 1)^{x_1}}{(R_0 - R_e)^{x_2}}, \qquad (168)$$

where the velocity has been scaled by the 1D sound velocity. In most cases $x_1 = x_2 = 1$, except in the square lattice without the dissipation term ($\mathcal{D}_2 = 0$), in which case $x_1 = 1$ and $x_2 = 1/3$. In this case fracture propagation was subsonic, whereas it was supersonic on the triangular lattice if $(R_0 - R_e)/(R_e - 1)$ was small (see Sections 7.7.1, 7.8.15, and 9.8.3.3 for discussions of supersonic fracture propagation). In the second case that was studied, all the lattice spacings were equal. A variety of interesting fracture patterns were obtained, some of which are shown in Figure 8.25. Also obtained were oscillatory fracture patterns, some of which are also shown in Figure 8.25. As discussed above and in Chapters 6 and 7, such oscillatory fracture patterns contribute to the dynamic instability observed during fracture propagation.

Another interesting deterministic lattice model of dynamic fracture was proposed by Rautiainen *et al.* (1995). In their model each bond $ij$ of a square lattice was an elastic element with an interaction energy that was described by

$$\mathcal{H}_{ij} = \frac{\alpha}{2}[(\mathbf{u}_i - \mathbf{u}_j) \cdot \mathbf{R}_\parallel]^2 + \frac{\gamma}{2}[(\mathbf{u}_i - \mathbf{u}_j) \cdot \mathbf{R}_\perp]^2, \qquad (169)$$

where $\mathbf{u}_i$ is the displacement of node $i$, and $\mathbf{R}_\parallel$ and $\mathbf{R}_\perp$ are the unit vectors parallel and perpendicular to the vector connecting $i$ and $j$ in the undeformed lattice, respectively. To include disorder in the model, a fraction $q$ of the bonds was removed at random before fracture simulations began. To introduce dynamics into the model, a dissipation mechanism was included in the model by incorporating in it Maxwellian viscoelasticity which allows the description of relaxation and dissipation of elastic energy as a dynamical decay of the local forces. The constitutive
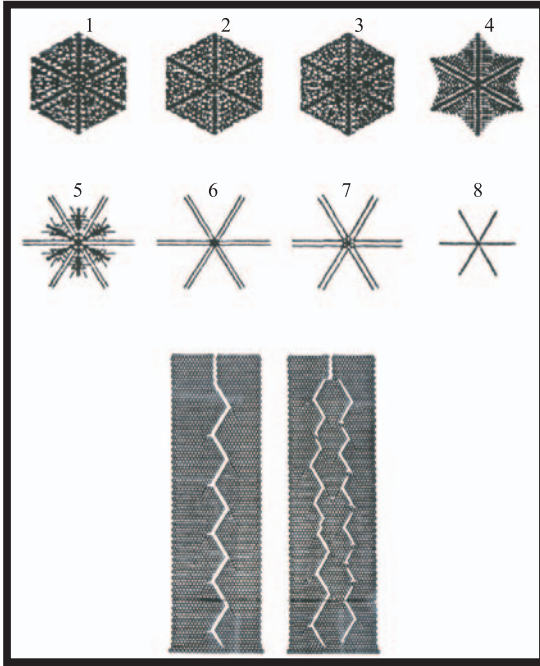
FIGURE 8.25. Regular and oscillatory fracture patterns obtained by the dynamic model of Furukawa (1993). The regular patterns were obtained under isotropic tension.

equation for the forces acting at each bond at time $t$ was taken to be

$$\frac{dF_{ij}}{dt} = \frac{dF_H}{dt} - \frac{1}{t_r}F_{ij}, \tag{170}$$

where $F_{ij}$ is the force between $i$ and $j$ arising out of their interaction, $F_H$ is the elastic force derived from Eq. (169), and $t_r$ is a phenomenological parameter which can be considered as a relaxation time scale. In effect, each bond is replaced by a Maxwellian viscoelastic element—a spring and a dashpot in series, and was considered as broken if its length exceeded a critical value.

The model predicts brittle fracture in the limit of slow straining. At finite strain rates, damage development becomes ductile with increasing dissipation. For small $t_r$ and $q$, the number of broken bonds increases rapidly at the initial stages of the fracture history. However, after some time damage accumulation stops and the system resists rupture. This is due to local viscoelastic dissipation which arrests crack growth. For large $t_r$ the number of broken bonds increases slowly, and the manner by which the bonds break is correlated. Thus, ductility increases with decreasing $t_r$. Hassold and Srolovitz (1989) had already shown, using a Born and the quasi-static lattice model described in Section 8.2, that crack arrest can be controlled by varying $\alpha/\gamma$. Therefore, the role of $\alpha/\gamma$ in the quasi-static model is played by the relaxation parameter $t_r$. The crack velocity was found to be approximately,
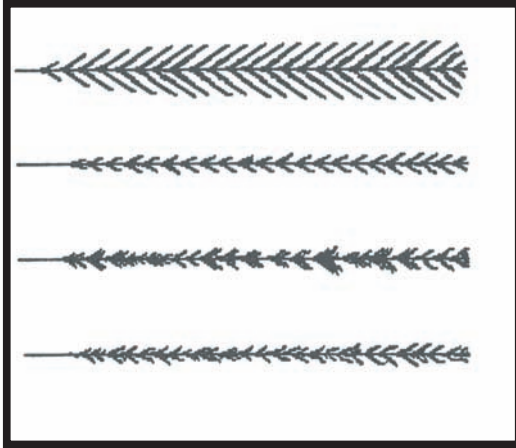
FIGURE 8.26. Fracture patterns in the dynamic model of Heino and Kaski (1996). The top pattern corresponds to $t_r = 0$, while the rest are for $t_r \neq 0$.

$v \simeq \sqrt{\gamma}$. A dynamic Born model was developed by Marfin *et al.* (2000) that also produced many of these predictions.

In another lattice model of dynamic fracture, Heino and Kaski (1996) used the same Hamiltonian as in Eq. (169), but the lattice was more complex. Each bond of the lattice consisted of two perpendicular springs in fixed directions representing tensile and bending behavior. A dashpot was connected in series to each spring. In a tensile experiment, the sites move which results in elongation or shrinking of the springs and dashpots. Disorder was introduced into the model by assuming that each bond has its own $\alpha$ and distributing it according to a uniform distribution. The ratio $\alpha/\gamma$ for each bond was held constant though, which corresponds to varying the Young's modulus of the bonds. A bond was considered as broken when its strain exceeded a threshold. Interesting fracture patterns were generated by the model by varying its parameters, some of which are shown in Figure 8.26. As can be seen, there is a dominating fracture that eventually spans the system and causes it to fail. However, many *daughter* cracks also appear essentially symmetrically on both sides of the main fracture. They appear periodically and advance a short distance before dissipation damps their growth. If the network is more disordered, the periodicity disappears and the fracture pattern becomes irregular.

However, the most interesting result emerging from this model was the behavior of the crack tip velocity as a function of the crack width. Initially, the velocity increased rapidly, corresponding to the emergence of a straight fracture. However, after some time the oscillatory daughter fractures appeared, and thus the velocity also started to oscillate with the crack width, and hence with the time. Increasing $\alpha$ increased the crack speed and its oscillation frequency, but decreased the length of the daughter fractures, although the angle that they made with the main fracture was unaffected by $\alpha$. Heino and Kaski (1997) combined a finite-element method (for discretizing the continuum elasticity equations and obtaining a finite-element

mesh) with the bond-breaking process used in lattice models. The model provided the same type of predictions as those of their earlier model.

Åström and Timonen (1996) used a square lattice of beams (see Section 8.2) to study dynamic fracture. Their model was the dynamic analogue of the lattice models of quasi-static brittle fracture described in 8.2. The network was strained by an amount $\epsilon$ in the $y$ direction. The sites at the top and bottom edges of the network were constrained to remain at their original positions, while the sites on the left and right edges were free to move without constraints. The dynamics of the model was calculated using a discrete form of the Newton's equation of motion including a linear viscous dissipation term (see also Section 8.4 below). A beam was considered as broken if the strain on it exceeded a pre-assigned threshold value. Åström and Timonen (1996) showed that, by tuning the strain and the ratio of axial to bending stiffness of the beams, a fracture can propagate either straight, or branch, or bifurcate. The fact that such features can be obtained by both the Born model of Heino and Kaski (1996), described above, and by the beam model of Åström and Timonen, indicates the universality of such features. Moreover, for the branching fracture, Åström and Timonen (1996) found that their trajectories follow a power law,

$$ y \sim x^{0.7}, \tag{171} $$

where $x$ and $y$ are, respectively, the directions parallel and perpendicular to the direction of the main crack, with the origin being the point at which the micro-branch begins. This result is in complete agreement with the power-law observed in experiments in both polymethylmethacrylate (PMMA) and glass (see Section 7.10.4).

### 8.3.4   Comparison with the Experimental Data

The predictions of lattice models of dynamic fracture are in agreement with the results of a set of spectacular experiments by Fineberg *et al.* (1991, 1992) and Gross *et al.* (1993). Their work and many earlier experiments have reported many interesting features of dynamic crack propagation in materials, which were described in detail in Chapters 6 and 7, and are summarized here. Recall, from Chapters 6 and 7, the main features of a fracture surface in a brittle material, namely, the mirror, mist, and hackle sequence: An initially smooth and mirror-like fracture surface becomes misty, and then evolves into a rough hackled region. It has also been reported (see, for example, Döll, 1975; Kusy and Turner, 1977) that in some brittle materials, such as PMMA, the fracture pattern exhibits characteristic wavelength, that surface roughness increases with crack speed (see, for example, Langford *et al.*, 1989, and references therein), and that periodic stress waves are emitted from the tip of the rapidly moving cracks in a wide variety of materials (see, for example, Rosakis and Zehnder, 1985; Dally *et al.*, 1985, and references therein). As described in detail in Chapter 7, Fineberg *et al.* (1991, 1992) carried out precise experiments to study fracture propagation in brittle plastic PMMA and showed that, there is a critical velocity $v_c$ beyond which the velocity of fracture tip begins to oscillate, the dynamics of the fracture changes abruptly, and a periodic fracture

pattern is formed. For $v > v_c$ the amplitude of the oscillations depends linearly on the mean velocity of the propagating fracture. Thus, the motion of fractures is governed by a dynamical instability, and explains why the velocity of their tip does not attain the limiting Rayleigh velocity $c_R$ predicted by the linear elasticity theory.

In another set of beautiful experiments, Gross *et al.* (1993) used two materials, the PMMA and soda-lime glass, to show that all features of dynamics of crack propagation in the two materials, such as acoustic emission, crack velocity, and surface structure, exhibit quantitative similarity with each other. Thus, there exist universal characteristics of fracture energy in most materials that are the result of energy dissipation in a dynamical instability. Perhaps the most spectacular experiments were carried out by Sharon *et al.* (1995) and Sharon and Fineberg (1996) using the brittle plastic PMMA. They identified the origin of the dynamical instability during fracture propagation as being the nucleation and growth of the daughter cracks which limit the speed of the propagating crack tip. The daughter fracture carries away a fraction of the energy concentrated at the tip of the moving crack, thus lowering the velocity of the tip. After some time, the daughter crack stops growing, and thus the crack tip velocity increases, until a new daughter fracture starts to grow, and so on. They also observed that the branching angle for a longer daughter fracture was smaller than that of the shorter daughter fractures. These features are all produced by the dynamic lattice models of Heino and Kaski (1996,1997) described above. The computations of Marder and Liu (1993) for a perfect crystal, that were described and discussed above, also agree with these data.

As already mentioned, oscillatory fracture patterns were also observed in the experiments of Yuse and Sano (1993). They imposed a temperature gradient along a thin glass plate, from a hot region to a cold one. A microcrack was introduced in the glass, and the glass was pushed. As the plate started to move the crack jumped ahead of the thermal gradient and stayed there. It was observed that if the plate moves slowly, the growing crack remains straight and stable. However, increasing the velocity to a critical value $v_c$ gives rise to a transition whereby the fracture path begins to oscillate and an instability appears. At still higher velocities crack branching appears; see Figure 6.6. Ronsin *et al.* (1995) also provided experimental data for brittle fracture propagation in thin glass strips, using a thermally induced stress field. In their experiments the temperature field was controlled by the width $w$ of the plate, and induced thermal expansion in the sample. It was observed that for widths below a critical value $w_c$ no fracture was formed. For $w_c < w < w_o$, where $w_o$ is a second critical width for the onset of oscillatory cracks, straight fractures were formed and propagated with a constant speed. For $w > w_o$ oscillatory fractures were generated that became more irregular as $w$ was increased beyond $w_o$.

## 8.4   Fracture of a Brittle Material by an Impact

Fracture of a brittle material by an impact, such as fracture of glass by an impact, is an important phenomenon that most of us have seen with our own eyes. Simple and interesting experiments performed by Shinkai (1994) demonstrated this
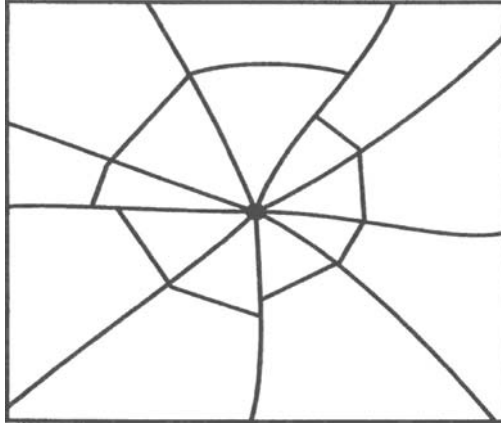
FIGURE 8.27. Schematics of fracture pattern on impact on a thin glass plate. Hertzian fracture is at the center (after Åström and Timonen, 1997b).

phenomenon very nicely. In his experiments thin square-shaped glass plates were supported at the edges, and a small and heavy ball was dropped on them from a point above the plates' centers. The resulting fracture patterns consisted of three types of cracks. If the impact velocity was high, only a circular hole at the point of impact was formed, which is usually referred to as a *it Hertzian fracture*. However, at lower impact velocities, radial and tangential cracks also appeared. An example is shown in Figure 8.27. The radial cracks were fairly straight and directed outwards from the point of impact, while the tangential ones formed a more or less circularly-symmetric crack. At still lower impact velocities, only radial cracks were formed, while if the impact velocity was too low, no crack was formed at all.

A simple and interesting model of dynamic fracture, in which the material cracks as a result of an impact, was proposed by Åström and Timonen (1997b). They used a triangular lattice of beams. The dynamics of the system were calculated using a discrete form of Newton's equations of motion given by

$$\left(\frac{1}{\Delta t^2}\mathbf{M} + \frac{1}{2\Delta t}\mathbf{D}\right)\mathbf{u}(t + \Delta t) = \left(\frac{2}{\Delta t^2}\mathbf{M} - \mathbf{C}\right)\mathbf{u}(t) - \left(\frac{1}{\Delta t^2}\mathbf{M} - \frac{1}{2\Delta t}\mathbf{D}\right)\mathbf{u}(t - \Delta t),$$
(172)

where $\mathbf{M}$ is a diagonal mass matrix, $\mathbf{C}$ is the stiffness matrix, $\mathbf{D}$ is a diagonal damping matrix, $\mathbf{u}$ is the vector that contains the displacements of the sites from their equilibrium positions, and $\Delta t$ is the length of the discrete time step. The stiffness matrix that was used was that of slender beams (that is, a beam in which bending is much larger than shear deformation). The boundary conditions imposed on the lattice in the $xy$ plane were such that the sites at the boundaries lattice are constrained to remain at their original positions, while the sites in a circular area in the middle of the lattice were forced to move a distance $-vt$ in the $z$ direction. For the lattice in the $xz$ plane only the sites at its left and right edges were constrained to
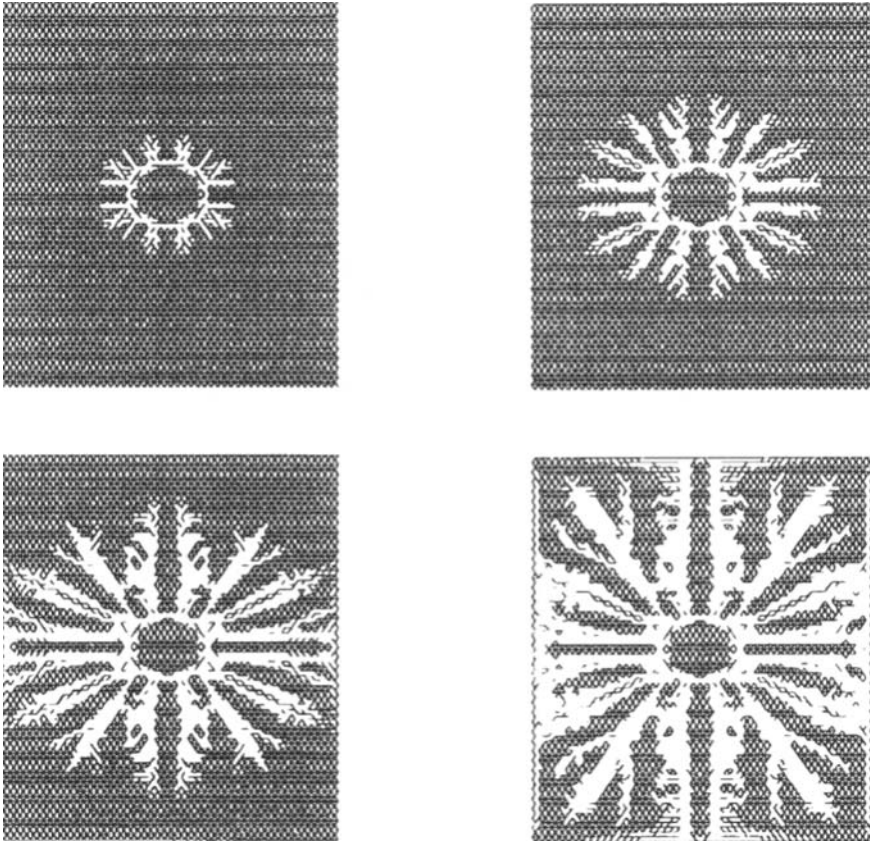
FIGURE 8.28. Radial fractures in a lattice with an externally applied in-plane strain. The results, from top left corner and clockwise, are for times $t = 200, 350, 700$ and $500$ (after Åström and Timonen, 1997b).

remain at their original positions, while a number of sites in the middle of the upper boundary moved in the negative $z$-direction. A beam was considered as broken if its elongation exceeded its pre-assigned threshold. However, beam fracture was not instantaneous. It was assumed that, at the threshold, the Young's modulus of a beam begins to decrease linearly in time until it reaches zero. The rate at which the modulus decreased was expressed by a parameter $r = c/\Delta t$, where $c$ is a constant.

Figure 8.28 presents the crack patterns in a lattice with an externally applied in-plane strain: $r = 0.1/\Delta t$, $v = 1/600\Delta t$, and the in-plane strain, $\epsilon_{xy} = 0.167$. The numbers refer to the time steps at which the snapshots were taken. All the qualitative features of the experimental crack patterns obtained by Shinkai (1994) were reproduced by this rather simple model, and therefore it may be used for studying other aspects of fracture of brittle materials by an impact.

## 8.5    Dynamic Fracture of Materials with Annealed Disorder

All the lattice models that we have so far described and discussed are deterministic in the sense that, a bond of the lattice breaks if its elongation, or the force that it suffers, exceeds some threshold value. Therefore, at any given time, the bond to be broken is identified deterministically, if the stress or strain field in the lattice is given. As such, these models are appropriate for materials in which the disorder is quenched. We now briefly describe the probabilistic lattice models of dynamic fracture which, as discussed in Section 8.2, may be appropriate for fracture of materials in which the disorder is annealed. As in the case of dielectric and electrical breakdown (see Chapter 5), the probability $p_b$ that a bond of a lattice, representing a spring or a beam, breaks is related to the force $F$ that it suffers. Similar to the dielectric breakdown model of Niemeyer *et al.* (1984) described in Section 5.4.1, $p_b$ is assumed to have the following form

$$p_b \propto F^{\eta}, \tag{173}$$

where, as in Niemeyer *et al.*'s model, $\eta$ is a parameter of the model. An exponential form, $p_b \propto \exp(cF^{\eta})$, has also been used (where $c$ is a constant) which would then be compatible with the theory of chemical processes according to which the reaction rate $\mathcal{R}$ of a chemical process, such as the rate of breaking the interatomic bonds, is given by

$$\mathcal{R} = k(T) \exp(-\mathcal{H}_a/k_B T), \tag{174}$$

where $k(T)$ is a temperature-dependent pre-exponential factor, $\mathcal{H}_a$ is the activation energy, and $k_B$ is the Boltzmann's constant.

   The value of $\eta$ in Eq. (173) is selected based on the physics of the problem. For example, one may argue that if an external force is exerted on the system, the activation energy is reduced by an amount which is proportional to $F^2$, and thus one may assume that $\eta = 2$. On the other hand, one may also argue that a bond breaks when its length has reached a critical threshold, and that for a harmonic potential the required energy to reach this length is proportional to $F$, and thus $\eta = 1$.

   The process time is explicitly incorporated into this model by the following algorithm. Each time a bond breaks, the process time is increased by an amount $\Delta t$ given by

$$\Delta t = \frac{1}{N p_b^{(m)}}, \tag{175}$$

where $N$ is the total number of unbroken bonds in the network, and $p_b^{(m)}$ is the maximum probability of breaking for *any* bond in the network at time $t$. If the bond breaking rates are actually known (i.e., if $p_0$ is known), then Eq. (175) provides an absolute time scale for the system; otherwise the true process time is directly proportional to the time scale calculated by (175).

FIGURE 8.29. Fracture pattern in a viscoelastic material (after Van Damme *et al.*, 1987a).

Such a model may be able to explain some aspects of experimental observations of Van Damme and co-workers (Van Damme *et al.*, 1986, 1987a,b). They displaced clay dispersions by air or water in a Hele–Shaw cell, an essentially 2D system consisting of two parallel glass sheets with a small gap between them. When the clay concentration was low, the displacement pattern was similar to diffusion-limited aggregates (DLA) described in Chapters 1 and 5. However, at high clay concentrations a transition was seen from a DLA-like displacement to a fracture pattern; see Figure 8.29. This is a good example of a system in which the disorder is annealed, since as the injected water or air displaces the clay, the material must continuously adjust itself to accommodate the fact that the displacing fluid is pushing its way into the clay. Thus, disorder in this system changes with the time.

Curtin and Scher (1991,1992,1997) and Curtin *et al.* (1997) studied lattice models of dynamic fracture in which each node of the lattice could be a site for nucleation of defects or fractures. It was assumed that the probability $p_i$ per unit time of nucleating a crack at site $i$ is a monotonically-increasing function $p_i[\sigma_i(t)]$ of the local stress $\sigma_i$ at $i$, and was taken to be $p_i(t) = A\sigma_i(\theta)^\eta$, where $\eta$ is a parameter of the model. Many interesting predictions emerge from this model either analytically or by numerical simulations, including, (1) failure is more abrupt as $\eta$ increases; (2) failure times scale inversely with the logarithm of the system size raised to some power, and (3) the distribution of failure times is Gumbel-like (double exponential; see Section 8.2.1.4) and becomes broader as $\eta$ increases, implying that failure becomes *less* predictable as it becomes more abrupt.

Another model that had some dynamics built into it was developed by Louis and Guinea (1987), and further developed by Fernandez *et al.* (1988) and Meakin *et al.* (1989). In this model a triangular lattice of Hooke's springs is used. It is assumed that only those bonds that are on the surface of the growing fracture can break.

An initial microfracture is inserted into the network at its center. Three different ways of breaking the bonds were considered. In model I only those bonds that join pairs of sites that are both on the fracture perimeter were broken. In model II any bond associated with a damaged node—one that had five or fewer unbroken bonds—was broken, while in model III any of the bonds associated with any of the sites at the fracture perimeter was broken. The probability of breaking any bond $i$ is given by

$$p_b = \frac{(\ell_i - \ell_0)^\eta}{\sum_{j=1}^{N}(\ell_j - \ell_0)^\eta},\qquad(176)$$

where $\ell_i$ is the length of bond $i$, $\ell_0$ is its equilibrium length, and $N$ is the number of the unbroken bonds on the perimeter of the fracture. The network was initially diluted by a small amount to prevent unwanted nonlinearities, and a constant force was applied on the boundaries of the network. The model produced fractal fracture patterns with a fractal dimension $D_f$ that depended on the parameter $\eta$ and the boundary conditions at the perimeter of the growing fracture. The fracture patterns were quite similar to diffusion-limited aggregates (see Chapters 1 and 5). Other types of boundary conditions were also used in this model. For example, shear strain and uniaxial tension were both used (Hinrichsen *et al.*, 1989) in which case the fracture pattern had an $X$-like shape. If only bonds in tension were allowed to break, then only one arm of the $X$-shape grew.

## 8.6    Fracture of Polymeric Materials

We already described in Section 6.16.1 the general fracture properties of polymers. Mechanical stability of polymeric materials and their resistance to fracture are clearly of great practical importance. Commercial products made of plastics and other types of polymeric materials must preserve their form under the allowed external loads, and therefore one must ensure that they do not develop fracture and break. For this reason, fracture of polymeric materials has been an active research field for a long time.

Deformation of a polymeric solid material may include, in addition to the reversible part, an irreversible plastic flow or yielding which sets in when the stress becomes large enough to surmount the yield point. Compared to metals and ceramics, the yield point of polymeric materials is low, so that moderate forces or stresses are often sufficient to trigger the yield process. Thermal effects play an important role, and the environment can also affect the properties of polymeric materials if fluids or gases can penetrate into these materials and degrade them.

Extensive experimental studies have shown that there are two mechanisms of yielding in polymeric materials which have very different appearances and can be easily determined. These two mechanisms are as follows.

(1)  The first mechanism is called *shear yielding*. Consider, for example, polyethylene. If we stretch a sample of this material with a constant rate, then its

load-extension (stress-strain) diagram will have the following characteristics. The stress increases at first but then, when the yield point is reached, it passes through a maximum and a neck develops somewhere in the polymer which can extend up to the full extent of the sample. This experiment takes place under an essentially constant tensile force, and eventually elongates the material by several times its original length. If the stretching continues, the force will increase up to the point of break. Shear yielding is typical for partially polycrystalline polymers, and has also been seen in many amorphous polymers, such as polycarbonate. Shear bands are also formed, and are oriented along the directions of maximum shear stress.

(2) Polystyrenes exhibit a quite different behavior. If we again carry out an experiment similar to what was described above and form the stress-strain diagram, we will find that it has the following characteristics. The force increases at first but then, after a slight bending, the material breaks before a maximum is reached. Inspection will then indicate the formation of many void containing microdeformation zones. As described in Chapter 6, these localized zones of plastic flow are called crazes, and *crazing* is usually used to refer to this second mechanism of yield in polymeric materials. However, we must keep in mind that shear yielding and crazing are *not* mutually exclusive, and more often than not they both are operative during deformation of polymeric materials. The stress condition and the temperature of the system are the controlling factors in deciding which mechanism is dominant in a polymeric material. Clearly, the amount of flow before fracture determines the ductility of a polymeric material, so that brittle polymers break without exhibiting much preceding flow.

Many of the precise experiments on dynamics fracture of brittle materials that were described in Chapters 6 and 7 were actually carried out with polymeric materials, such as PMMA, and therefore all the continuum and discrete models of brittle fractures that were described in Chapter 7 and this chapter are applicable to such materials. Similarly, theories of ductile behavior of materials can also be used for investigating this mode of deformation in polymeric materials. In addition, we summarized in Chapter 6 the most important fracture properties of polymeric materials (see Section 6.16.1). Given that the field of polymer fracture is well-described and documented (see, for example, Kinloch and Young, 1983; Kausch, 1987), there is no need here for a long discussion of fracture behavior of polymeric materials. Instead, we restrict ourselves to a brief discussion of recent discrete models of dynamic fracture in such materials.

There is strong evidence that applying a stress $\sigma$ to a polymer reduces the activation energy by an amount that is proportional to $\sigma$. This means that the bond breaking rate, or the probability that a bond breaks, can be written as

$$p_b = p_0 \exp[-(\mathcal{H}_a - \Omega_a \sigma)/k_B T], \tag{177}$$

where $\Omega_a$ is the activation volume (in 3D) or surface (in 2D) of the system. In a network model $\Omega_a \sigma$ is replaced with $L_a F$, where $L_a$ is the activation length of the

bonds. If all the bonds are equivalent, then in the limit $T \to 0$ the bond with the largest strain will always break first, and thus in this limit the probabilistic models reduce to the deterministic ones described earlier in this chapter. In the opposite limit, $T \to \infty$, the bond breaking process becomes completely random and thus represents a percolation process. The process time is incorporated into this model by the algorithm described above using Eq. (175).

Aside from the early work of Dobrodumov and El'yashevich (1973), the first model of this type was developed by Termonia and Meakin (1986). In their 2D model the probability $p_i$ that a bond $i$ breaks is given by

$$p_i \propto \exp[\alpha_i e_i^2 / (2k_B T)], \tag{178}$$

where $\alpha_i$ is the elastic constant of bond $i$, and $e_i$ is its elongation. Fractal fracture patterns were generated by this model with a fractal dimension $D_f \simeq 1.3$ in 2D. Termonia and Smith (1986) and Termonia *et al.* (1985, 1986) developed probabilistic models of mechanical and fracture properties of polymer fibers, which possess a very complex morphology. In their models there is a distinction between the primary bonds—those that are parallel to the fiber axis—and the secondary bonds—those that are perpendicular to the fiber axis. The primary bonds are strong covalent bonds, while the secondary bonds are the much weaker van der Waals and hydrogen bonds. Defects are also included in these models by removing a fraction of the bonds before deformation of the system is started. The bonds break with a probability given by Eq. (177), but the activation energies and volumes for the primary bonds were about 2 orders of magnitude larger than those for the secondary bonds, while $p_0$ was assumed to be the same for both types of bonds. The secondary bonds were allowed to reform between adjacent sites if their coordinates in the direction of the primary bonds became equal, whereas primary bonds were not allowed to do so. Figure 8.30 shows the fracture patterns generated by this model (under isothermal condition), which are in agreement with typical experimental observations.

Termonia and Smith (1987,1988) developed models of polymer deformation and failure in which the effect of chain slippage and the release of entanglements
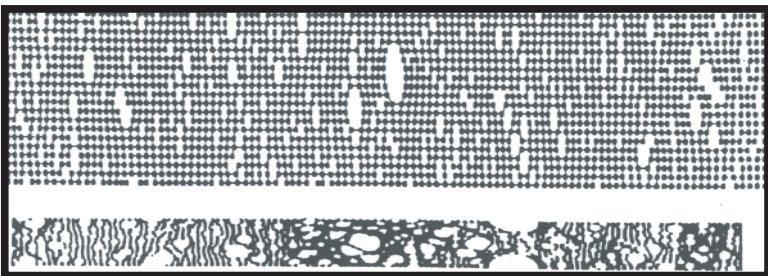


FIGURE 8.30. Fracture pattern in the model of Termonia *et al.* (1985,1986) in which the probability of failure for the primary and secondary bonds were not the same. The top pattern corresponds to 150% strain, while the bottom one is for 300% strain.

was taken into account. Such effects play a prominent role in failure of polymers. In their model the rate of failure of van der Waals bonds and that of chain slippage were both given by Eq. (177), but with different activation energies and volumes. For the chain slippage process, $\sigma$ represents the stress difference between two parts of a chain separated by an entanglement point. Their model used a 2D diamond-like lattice with nodes that represented the entanglement points between pairs of polymers. The lattice was then decorated randomly with polymer molecules that intersect only at the entanglement points until there is an entanglement point associated with every node. The stress $\sigma$ is predicted by the classical theory of rubber elasticity to be

$$\sigma = \beta k_B T \mathcal{L}^{-1} \left( \frac{R}{n_c \ell} \right), \tag{179}$$

when $n_c$ was the number of chain segments of length $\ell$ between a pair of entanglement points separated by a distance $R$, $\mathcal{L}(z) = \coth(z) - 1/z$ is the Langevin function, and $\beta$ is given by

$$\beta = \frac{N_c \sqrt{n_c}}{3}, \tag{180}$$

with $N_c$ being the number of chain strands per unit volume. The predicted fracture patterns were found to be in good agreement with experimental observations. More details can be found in the review by Meakin (1990).

## 8.7   Fracture of Thin Solid Films

The last class of materials for which dynamic lattice models of fracture have been developed is thin solid films. A thin film is attached to a substrate with mechanical properties that are usually very different from those of the film. As a result, the surface layer usually suffers large stresses that are generated by various factors, such as decohesion, buckling, spalling, and in-plane cracking, leading to fracture patterns that resemble those in dried-up mud. Some examples are $Al_2O_3$ sputter deposited onto copper (see, for example, Jarvinen *et al.*, 1984), and chromium metal electrodeposited onto an aluminum alloy (see, for example, Namgoong and Chun, 1984).

Meakin (1987) developed a simple model for this type of fracturing process. In his model the thin layer is represented by a triangular network of Hooke's springs. In addition, each node of the network is connected to a rigid substrate by another Hooke's spring much weaker than those in the triangular network. Only the bonds in the triangular network are allowed to break. The probability of their failure is given by Eq. (178) with $k_B T = 1$. In a typical simulation, the bond length $\ell = 1.0$ at the start of the process, and the equilibrium length $\ell_0 = 0.90$. Figure 8.31 shows the fracture pattern resulting from such a model.

Skjeltrop and Meakin (1988) modified Meakin's original model in order to simulate fracture patterns of polystyrene bead monolayers. The monolayers are constructed by placing an aqueous dispersion of microspheres between two parallel
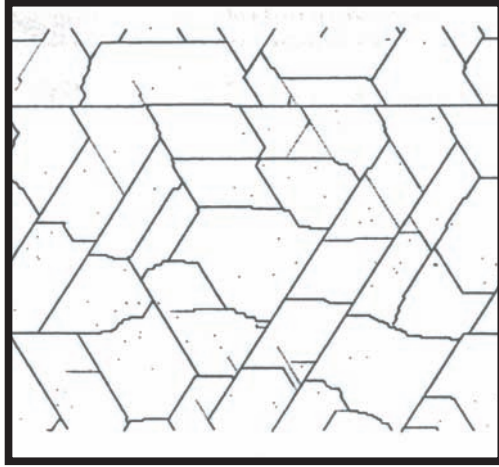
FIGURE 8.31. A typical fracture pattern in thin films, for which 5000 bonds have been broken (after Meakin, 1987).

sheets of glass and allowing the water to evaporate slowly along one edge of the gap between the two sheets. The system resembles a thin film. In the model of Skjeltrop and Meakin (1988) the points of attachment of the network's nodes to the substrate were allowed to move. The force that a weak bond that connects a node $i$ of the triangular network to the substrate exerts on $i$ is given by $\mathbf{F} = \alpha_w(\mathbf{R}_i^0 - \mathbf{R}_i)$, where $\alpha_w$ is the elastic constant of the weak bond, $\mathbf{R}_i$ is the position of the $i$th node, and $\mathbf{R}_i^0$ is its initial position. If $|\mathbf{R}_i^0 - \mathbf{R}_i|$ exceeds a threshold, the point of attachment to the substrate at $\mathbf{R}_i^0$ is moved towards the current position $\mathbf{R}_i$ of the node until $|\mathbf{R}_i^0 - \mathbf{R}_i|$ becomes equal to the threshold. The resulting fracture patterns were in very good agreement with the experimental patterns.

More recent efforts in this area include those of Crosby and Bradley (1997) and Leung and Néda (2000). The latter group used a model in which the grains in the thin films were represented by an array of blocks. Each pair of the neighboring blocks was connected by a bundle of coil springs having a unit spring constant and an equilibrium length $\ell$. Initially, the blocks are randomly displaced by $\mathbf{r} \ll \ell$ about their mean positions on a triangular lattice. For slow cracking on a frictional substrate, the motion of the grains is over-damped. Thus, the system evolves quasi-statically with a driving rate which is much slower than the relaxation rate. This implies that one does not have to solve the equation of motion, but rather update the configuration of the system according to certain criteria.

In a typical cracking of thin films, the film often hardens and/or weakens in time, hence tending to contract which is, however, resisted by friction from the substrate. As a result, as mentioned above, stress is built up and relaxed slowly. To incorporate such effects in any reasonable model, one can pre-strain the array of the blocks. Then, two force thresholds, $F_s$ for slipping and $F_c$ for cracking, are decreased systematically to induce these two modes of motion. For example,

in a drying experiment, the narrowing and breaking of liquid bridges between the grains, due to evaporation, provide an example of this driving. This model generates cracking patterns that are very much similar to those that are developed in thin films during their drying.

## Summary

Lattice models offer a detailed account of most aspects of fracture, including those that had been considered too complex to quantitatively describe, or those that had simply been ignored. At the same time, there is no conflict between the lattice models and conventional fracture mechanics. In addition, lattice models of crystalline materials make the *additional* prediction that a forbidden band of velocities exists for cracks. This means that a fracture can only propagate stably above a *finite* minimum velocity. Molecular dynamics simulations (see Chapter 9) of crystalline materials indicate that this forbidden band of states may disappear at room temperature, but should be observable in low-temperature experiments.

Through precise experiments (most of which were described in Chapters 6 and 7), large-scale simulations of lattice models of quasi-static and dynamic fracture described in this chapter, and molecular dynamics simulations to be described in Chapter 9, we now have a deeper understanding of fracture propagation, and in particular the structure and dynamics of mechanisms of dissipation within the near vicinity of the tip of a propagating fracture in a brittle material. We now know that fracture in brittle materials is governed by a dynamic instability that leads to repeated attempts for fracture branching. Although the instability also appears in dynamic finite-element simulations of fracture (see Chapter 7), it appears to have no analytical explanation in a continuum framework. In fact, many classical models of the cohesive zone have been shown to be ill-posed in that, they admit a *set* of possible states under *identical* conditions. Theories formulated on a lattice, on the other hand, do not exhibit such difficulties. It remains to be seen whether a simple continuum limit exists, or whether a crucial ingredient in understanding fracture is the discreteness of the underlying atoms.

However, despite considerable progress, our understanding of fracture propagation phenomena is still incomplete. In addition to not having a simple continuum limit of the discrete models, we must keep in mind that the most detailed and precise experiments on brittle fracture have so far been carried out in amorphous materials at room temperature, whereas the most detailed theories currently available apply mostly to very low (or zero) temperatures. A general theoretical framework for analyzing fracture propagation over a wide range of brittle materials has not yet been developed.

# Part III

# Atomistic and Multiscale Modeling of Materials

# 9
# Atomistic Modeling of Materials

## 9.0   Introduction

In all the chapters of Volume I, as well as in the present Volume so far, we have uti-
lized continuum mechanics and lattice models to describe modeling and simulation
of morphology of heterogeneous materials and estimating their effective proper-
ties. These models are appropriate for microscopic as well as macroscopic length
scales, but cannot provide any insight into materials' properties at the smallest
length scales, namely, the molecular scale. In this chapter we describe and discuss
modeling and simulation of materials and their properties at the molecular scale.
To achieve our goal we describe three important theoretical and computational
tools that have been developed over the past three decades, namely, the density
functional theory (DFT) and its variants, classical molecular dynamics (MD) simu-
lation, and quantum MD (QMD) technique. The advent of very fast computers and
development of massively-parallel computational strategies have made it feasible
to carry out large scale calculations at the molecular level, and in this endeavor
these three methods have become indispensable tools for predicting the properties
of materials at such length scales.

Prediction of electronic properties and morphology of a material requires, in
principle, quantum-mechanical computation of the total energy of the system and
minimization of this energy with respect to the electronic and nuclear coordinates.
To carry out such computations, one must start with the Hamiltonian of the system
which, for a system of $N$ electrons and $N'$ nuclei with charges $Z_n$, is given by

$$\mathcal{H} = \sum_{i=1}^{N} \frac{p_i^2}{2m} + \sum_{n=1}^{N'} \frac{P_n^2}{2M_n} + \frac{1}{2}\frac{1}{4\pi\varepsilon_0} \sum_{i=1,i\neq j}^{N} \sum_{j=1}^{N} \frac{e}{|\mathbf{r}_i - \mathbf{r}_j|}$$

$$-\frac{1}{4\pi\varepsilon_0}\sum_{n=1}^{N'}\sum_{i=1}^{N}\frac{Z_n e^2}{|\mathbf{r}_i - \mathbf{R}_n|} + \frac{1}{2}\frac{1}{4\pi\varepsilon_0}\sum_{n=1,n\neq n'}^{N'}\sum_{n'}^{N'}\frac{Z_n Z_{n'}\, e^2}{|\mathbf{R}_n - \mathbf{R}_{n'}|}, \tag{1}$$

where $p_i$ and $P_n$ are the momenta of the $i$th electron and the $n$th nucleus, re-
spectively (subscript $i$ refers to the electrons and $n$ to the nuclei), $\mathbf{r}_i$ and $\mathbf{R}_n$ are
their position vectors, $m$ is the electron mass, $M_n$ is the mass of the $n$th nucleus,
$e$ is the electron charge, and $\varepsilon_0$ is the permittivity. The first two terms of Eq. (1)
represent the kinetic energy of the electrons and nuclei, respectively, while the
third and fourth terms are, respectively, the result of Coulomb repulsion between

the electrons and Coulomb attraction between electrons and nuclei. Equation (1) is too complex for use in exact computations, especially when one must deal with a large system, and therefore reasonable approximations must be made in order to make the computations feasible. One obvious simplification can be made by taking advantage of the fact that there is a large difference in mass between the electrons and nuclei, while the forces on the particles are the same. Therefore, the electrons respond essentially instantaneously to the motion of the nuclei. As a result, electronic and nuclear coordinates in the many-body wave function can be separated, and the nuclei can be treated adiabatically. This separation is the well-known *Born–Oppenheimer approximation* which reduces the solution of the many-body problem to that of the dynamics of the electrons in some frozen-in configurations of the nuclei. Thus, the Hamiltonian of the system in the Born–Oppenheimer approximation is given by

$$\mathcal{H} = \sum_{i=1}^{N} \frac{p_i^2}{2m} + \frac{1}{2}\frac{1}{4\pi\varepsilon_0} \sum_{i=1, i\neq j}^{N}\sum_{j=1}^{N} \frac{e^2}{|\mathbf{r}_i - \mathbf{r}_j|} - \frac{1}{4\pi\varepsilon_0}\sum_{n=1}^{N'}\sum_{i=1}^{N} \frac{Z_n e^2}{|\mathbf{r}_i - \mathbf{R}_n|}. \quad (2)$$

Based on such simplifications, other approximate schemes have been suggested. A well-known example of such approximations is the Hartree–Fock theory. Suppose that

$$\mathcal{F}\{\Psi_k\} = \epsilon_k \Psi_k, \quad (3)$$

where $\mathcal{F}$ is called the *Fock operator* defined below, $\epsilon_k$ is the eigenvalue of the operator, and $\Psi_k(\mathbf{x})$ is the spin-orbital (electronic state or wave function), a function that depends on the spatial position $\mathbf{x}$ and spin coordinates of *one* electron. The Fock operator is defined by

$$\mathcal{F}\{\Psi_k\} = \left[ -\frac{1}{2}\nabla^2 - \sum_n \frac{Z_n}{|\mathbf{r}_n - \mathbf{R}_n|} \right]\Psi_k(\mathbf{x}) + \sum_{l=1}^{N} \int |\Psi_l(\mathbf{x}')|^2 \frac{1}{|\mathbf{r} - \mathbf{r}'|}\Psi_k(\mathbf{x})\,dx'$$

$$- \sum_{l=1}^{N} \int \Psi_l^*(\mathbf{x}') \frac{1}{|\mathbf{r} - \mathbf{r}'|}\Psi_k(\mathbf{x}')\Psi_l(\mathbf{x})\,dx' \quad (4)$$

where $*$ denotes a complex conjugate, and the following notation convention has been used: $\int dx' = \sum_{s'} \int d^3\mathbf{r}'$; that is, $\int dx'$ denotes a sum over the spin $s'$ and an integral over the spatial coordinate $\mathbf{r}'$. We have used the standard atomic units (using the Bohr radius, and the electron mass and charge as the basic units) so that, for example, instead of having the usual factor $\hbar^2/(2m)\nabla^2$ we have $\frac{1}{2}\nabla^2$. We use this convention throughout this chapter. Equations (3) and (4) constitute the well-known Hartree–Fock theory, also known as the *self-consistent-field theory*. Note that the fourth term of Eq. (4) is a non-local term, since although the operator acts on $\Psi_k$, its value at $\mathbf{r}$ is determined by the value assumed by $\Psi_k$ at *all* possible positions $\mathbf{r}'$. Note also that the ground state electron density $\rho(\mathbf{r})$ is given by

$$\rho(\mathbf{r}) = \sum_i |\Psi_i(\mathbf{r})|^2. \quad (5)$$

The Hartree–Fock equation has been a pillar of computation of the electronic structures of atoms and molecules. It is a nonlinear equation which is solved numerically by a sort of self-consistent iterative procedure, the essence of which is as follows. Since solving Eq. (4) yields an infinite spectrum, to obtain the ground state, one must take the lowest $N$ eigenvalues of the spectrum as the electronic states (spin-orbitals of the electrons). These constitute the first approximate solution of Eq. (4) which are then utilized for constructing the next iteration of the Fock operator. The operator is diagonalized again to obtain the next approximate solution to the electronic states. The procedure is repeated until convergence has been achieved.

However, computation of properties of solid materials based on the Hartree–Fock equation is an extremely difficult problem, and therefore, over the past several decades, other alternatives have been developed. Notable among them is the DFT developed by Hohenberg and Kohn (1964) and Kohn and Sham (1965), which is now widely used for predicting electronic properties of hard materials. (In recognition of his contributions to this research field, Kohn received the 1998 Nobel Prize in chemistry.) In words, the DFT allows one, in principle, to map *exactly* the problem of a strongly interacting electron gas, in the presence of nuclei, onto that of a single particle moving in an effective *non-local* potential, and provides an expression for the total energy of the system. The effective potential is not known exactly, but often, for reasons that are not completely understood yet, local approximations (see below) to the non-local potential are highly accurate. One can then minimize the total energy of the system that is provided by the DFT, often referred to as the Kohn–Sham total-energy functional, in order to determine its various properties. Beginning in the mid 1970s, the DFT was revitalized because extensive computations indicated that local approximations to the effective non-local potential can predict a variety of ground-state properties of materials that are within a few percent of experimental data (Dreizler and Gross, 1990).

Such a computational strategy, which requires only a specification of the ions present (by their atomic numbers), is usually referred to as an *ab initio* method. Although nearly two decades ago most *ab initio* methods, including the DFT, were capable only of modeling systems of a few atoms, they can now model systems with a large number of atoms. Of all *ab initio* methods, the total-energy pseudopotential technique (see below) based on the DFT, stands alone. This technique, combined with a method developed by Car and Parrinello (1985) (see Section 9.4), or a direct method of minimization of the total energy, such as the conjugate-gradient technique (see Section 9.5.2), have transformed the way in which people view quantum-mechanical *ab initio* computations and hence total-energy pseudopotential calculations because, in addition to their remarkable efficiency, they allow *ab initio* quantum-mechanical computations at *non-zero* temperatures.

Computation of materials' properties based on the DFT involves quantum-mechanical calculations. Over the past two decades, MD simulations, in which atoms and molecules are treated as classical particles and quantum-mechanical effects are neglected, have also become an important tool for investigating and predicting various static as well as dynamical properties of materials. We refer

to this method as the classical MD simulation. The increasing popularity of the classical MD methods is due to the fact that, over the past decade, highly efficient simulation techniques, based on massively-parallel and vectorized computations (see Section 9.6), have been developed that allow one to carry out MD simulations with *billions* of atoms, with the world record at the time of writing this book being a MD simulation with 5,180,116,000 particles (Roth *et al.*, 2000). Moreover, such MD techniques allow us now to simulate much longer time periods than what was believed to be feasible only a decade ago.

In order to increase the efficiency of *ab initio* computations, one can combine quantum electronic structures with a MD simulation to not only calculate the nuclear positions, but also the electronic charge cloud. This method, pioneered by Car and Parrinello (1985), is what we refer to as the QMD. This method describes a system in which the electronic structure does not, in general, completely relax to the true ground state, but follows it rather closely, and has been proven to be highly successful for describing many properties of materials.

Therefore, the purpose of this chapter is twofold.

(1) We wish to provide an overview of the essential concepts and ideas of the *ab initio* quantum-mechanical computations, MD and QMD simulation, and related problems, such as direct methods of minimization of the total energy, and vectorized and parallelized algorithms for MD simulations. Our goal is to describe the advances that have been made in these areas. Our discussion is not, and cannot be, exhaustive, as comprehensive description of all aspects of each subject would require a book by itself. For example, the review by Abraham (1986), and the books by Allen and Tildesley (1987) and Rapaport (1995) describe the classical MD simulations in detail. Instead, we set for ourselves the modest goal of outlining the basic concepts, ideas, and techniques of each of these computational tools which, together with adequate references to the recent literature, should enable the interested reader to pursue them. These techniques, with their ever increasing accuracy and efficiency, have gained popularity as major tools of investigating static as well as dynamical properties of materials, and therefore it is more than appropriate for this book to describe them.

(2) In addition to several examples discussed throughout this chapter, we also discuss, as an application of computations of materials' properties at the molecular scales, MD simulation of dynamic fracture of materials. At first, it may seem strange that computer simulations that involve *atoms* can be used for studying fracture propagation which is a *macroscopic* phenomenon involving *large-scale* structural changes. Even at molecular scales the motion of dislocations over *long distances* is of primary interest. So, why should computer simulations in which one arranges atoms in a crystal lattice and deforms it have anything to do with the true dynamics of fracture propagation in the crystal? The answer is obvious: Fracture of brittle materials is a physical process which *naturally* connects small and large length scales. Although stresses and strains that cause deformation and fracture of a material are applied at macroscopic

scales, fracture itself is the severing of the bonds at atomic scale. This implies that MD simulation of dynamic fracture in materials that account for the phenomena at molecular scales does not have to be very large, although large-scale MD simulations of dynamic fracture with up to about $10^8$ atoms have been carried out (Abraham *et al.*, 1997a,b; see below).

Our discussions in this chapter are also a prelude to Chapter 10 where we describe how atomistic and molecular modeling of materials is integrated with the microscopic and macroscopic methods of the previous chapters in order to develop a *multiscale approach* for investigating materials' properties over several widely disparate length scales.

## 9.1 Density-Functional Theory

Comprehensive reviews of the DFT and its applications are given by Jones and Gunnarsson (1989) and Payne *et al.* (1992). Some of our discussions in this section, and in Sections 9.3 and 9.4 closely follow these reviews. Motivated by the fundamental theorems of Hohenberg and Kohn (1964) of the DFT, Kohn and Sham (1965) developed a set of accessible one-electron self-consistent eigenvalue equations that have provided a practical means of realistic electronic structure calculations on a large array of atoms, molecules and materials (see, for example, Parr and Yang, 1989). The Hohenberg–Kohn theorems were extended to finite-temperature quantum systems (Mermin, 1965) and to purely classical fluids (see, for example, Hansen and McDonald, 1986). In addition, an integral formulation of electronic structure has also been developed in which the one-electron density is obtained directly without the introduction of orbitals (Harris and Pratt, 1985), although to date this theory, despite its promise, has not been used extensively in numerical studies.

The electron density $\rho(\mathbf{r})$ is subject to the constraint that

$$\int \rho(\mathbf{r})\, d^3r = N, \tag{6}$$

where $N$ is the total number of electrons in the system. The Hamiltonian of a many-electron system is given by

$$\mathcal{H} = \sum_i \left[ -\frac{1}{2}\nabla_i^2 + U_c(\mathbf{r}_i) \right] + \frac{1}{2} \sum_{i,i\neq j} \sum_j \frac{1}{|\mathbf{r}_i - \mathbf{r}_j|}, \tag{7}$$

where $U_c(\mathbf{r})$ is an external potential representing the Coulomb attraction by the frozen-in nuclei. The ground state can be determined by a variational approach which is carried out in two steps.

(1) One minimizes, for a given electron density, the energy functional with respect to the wave functions that are consistent with the given density.

(2) The result,

$$E[\rho(\mathbf{r})] = \min_\Psi \langle \Psi | \mathcal{H} | \Psi \rangle, \tag{8}$$

is then minimized with respect to $\rho$, subjected to the constraint imposed by
Eq. (6). If we now separate the Hamiltonian, $\mathcal{H} = \mathcal{H}_0 + U_c(\mathbf{r})$, where $\mathcal{H}_0$
represents the Hamiltonian without the external potential (i.e., the Hamiltonian
of a homogeneous electron gas), then

$$E(\rho) = \min_\Psi [\langle \Psi | \mathcal{H}_0 | \Psi \rangle] + \int U_c(\mathbf{r}) \rho(\mathbf{r}) \, d^3 r = E'(\rho) + \int U_c(\mathbf{r}) \, d^3 r,$$
(9)

where $E'(\rho) = \min_\Psi [\langle \Psi | \mathcal{H}_0 | \Psi \rangle]$ does not depend on the external potential
$U_e$.

The main problem lies in the fact that $E(\rho)$ is not known for both interacting and
non-interacting electron systems. We know, however, that in the non-interacting
case $E(\rho)$ can be written as

$$E(\rho) = E_k(\rho) + \int U_c(\mathbf{r}) \rho(\mathbf{r}) \, d^3 r,$$
(10)

where $E_k(\rho)$ is the kinetic contribution to $E(\rho)$. Variation of $E$ with respect to $\rho$
yields the following equation

$$\frac{\delta E_k(\rho)}{\delta \rho(\mathbf{r})} + U_c(\mathbf{r}) = \lambda \rho(\mathbf{r}),$$
(11)

where $\lambda$ is the Lagrange multiplier associated with the constraint (6). The exact
form of $E_k(\rho)$ is, of course, unknown, but we know that the ground state density is
given by Eq. (5), and that the spin-orbitals satisfy the single-particle Schrödinger
equation:

$$\left[ -\frac{1}{2} \nabla^2 + U_c(\mathbf{r}) \right] \Psi_k(\mathbf{r}) = \epsilon_k \Psi_k(\mathbf{r}),$$
(12)

where $\epsilon_k$ are the associated eigenvalues. The spin-orbitals $\Psi_k$ must be normalized
in order for the constraint (6) to be satisfied. It is clear that if $E_k(\rho)$ can be used in
place of $E'(\rho)$, then it must also be assumed that $E_k$ is independent of the external
potential $U_e(\mathbf{r})$ (just as $E'$ is independent of $U_e$). We now write down the energy
functional for a many-electron system with electronic interactions included:

$$E(\rho) = E_k(\rho) + \int U_c(\mathbf{r}) \rho(\mathbf{r}) \, d^3 r + \frac{1}{2} \int d^3 r \int \rho(\mathbf{r}') \frac{1}{|\mathbf{r} - \mathbf{r}'|} \rho(\mathbf{r}) \, d^3 r' + E_{Xc}(\rho).$$
(13)

The first three terms of Eq. (13) represent the contribution to $E(\rho)$ by the non-
interacting electron gas, whereas $E_{Xc}(\rho)$, which is called the *exchange-correlation
energy*, represents all other contributions to $E(\rho)$ that are not accounted for by the
first three terms. Equation (13) has the advantage that it contains no approximation,
and that all the unknown contributions have been lumped in $E_{Xc}(\rho)$. Varying
Eq. (13) with respect to $\rho$, we obtain

$$\frac{\delta E_k(\rho)}{\delta \rho(\mathbf{r})} + U_c(\mathbf{r}) + \int \frac{1}{|\mathbf{r} - \mathbf{r}'|} \rho(\mathbf{r}') \, d^3 r' + \frac{\delta E_{Xc}(\rho)}{\delta \rho(\mathbf{r})} = \lambda \rho(\mathbf{r}).$$
(14)

In essence, we have an *effective potential* given by

$$U_{eff}(\mathbf{r}) = U_c(\mathbf{r}) + \int \frac{1}{|\mathbf{r} - \mathbf{r}'|} \rho(\mathbf{r}') d^3 r' + \frac{\delta E_{Xc}(\rho)}{\delta \rho(\mathbf{r})}. \qquad (15)$$

Therefore, one may write

$$\left[ -\frac{1}{2} \nabla^2 + U_{eff}(\mathbf{r}) \right] \Psi_k(\mathbf{r}) = \epsilon_k \Psi_k(\mathbf{r}). \qquad (16)$$

Putting everything together, we finally obtain

$$E = \sum_{i=1}^{N} \epsilon_i - \frac{1}{2} \int d^3 r \int \rho(\mathbf{r}) \frac{1}{|\mathbf{r} - \mathbf{r}'|} \rho(\mathbf{r}') \, d^3 r' - \int \frac{\delta E_{ex}(\mathbf{r})}{\delta \rho(\mathbf{r})} \rho(\mathbf{r}) \, d^3 r + E_{Xc}(\rho). \qquad (17)$$

Equations (5) and (15)–(17) constitute the DFT, first derived by Kohn and Sham (1965).

As Eqs. (15)–(17) indicate, the Kohn–Sham equations map the interacting many-electron system onto a system of non-interacting electrons moving in the effective potential $U_{eff}$ generated by all the other electrons. If the exchange-correlation energy functional $E_{Xc}$ were known exactly, then taking the functional derivative with respect to the density would produce an exchange-correlation potential that includes the effects of correlation and exchange *exactly*. However, $E_{Xc}$ is not known exactly, and thus the main task is developing an accurate approximation for it. Once this approximation has been decided upon, the bulk of the work involves finding the solution of the eigenvalue problem described by Eq. (16). Note that the main difference between the Hartree–Fock and DFT approximations is that, the latter replaces the Hartree–Fock exchange term by the exchange-correlation energy. A particularly powerful method for developing accurate approximations for $E_{Xc}$ is the so-called local-density approximation which we will discuss shortly.

Thus, minimization of the total energy functional can be carried out directly using the Kohn–Sham model, Eqs. (16) and (17), subjected to the constraint (6), utilizing various minimization techniques, such as the simulated annealing method described in Chapter 3 of Volume I, or the conjugate-gradient method that will be described in Section 9.5. Usually, the eigenvalues $\epsilon_k$ are interpreted as the excitation energy. When such interpretation has been utilized for predicting many properties of various atoms and molecules, the predictions have been found to be in excellent agreement with the experimental data. For example, Table 9.1 compares the DFT predictions for the cohesive energy and lattice constants of diamond, silicon and germanium, and it is clear that the agreement between the predictions and the data is very good.

Let us mention that the DFT has been extended to time-dependent problems (Runge and Gross, 1984). Time-dependent DFT (TDDFT) makes it possible to apply the DFT to excited states of many-body systems. In this scheme, one must define a dynamic exchange-correlation energy $E_{Xc}(\rho, \mathbf{r}, t)$ which must somehow be calculated. It has been shown that substantial improvements of excitation energies with respect to the original Kohn–Sham eigenvalues are obtained with the

TABLE 9.1. The DFT predictions of lattice constants and cohesive energies for three materials, and their comparison with the experimental data. Atomic units have been used (adapted from Jones and Gunnarsson, 1989).

| Material | Lattice Constant | | Cohesive Energy | |
|---|---|---|---|---|
| | DFT | Expt. | DFT | Expt. |
| Diamond | 6.807 | 6.740 | 7.58 | 7.37 |
| Si | 10.30 | 10.26 | 4.84 | 4.64 |
| Ge | 10.69 | 10.68 | 4.04 | 3.85 |

TDDFT. However, even within this extension, certain problems in solids persisted for some time, such as the wrong Kohn-Sham band gap, regardless of the type of approximation used for $E_{Xc}$. Tokatly and Pankratov (2001) showed that, at excitation frequencies, $E_{Xc}(\mathbf{r}, \mathbf{r}')$ exhibits a highly non-local behavior with a range that is as large as the system itself, and hence it diverges as the system's size becomes very large. They developed a perturbation technique which maintains a correct electron density in every order of the perturbation theory, and therefore removed this unphysical feature of the TDDFT.

However, in general, the DFT has been designed for predicting the *ground-state* properties, and the Kohn–Sham eigenvalues are actually the derivatives of the total energy with respect to the occupation numbers of these states (Janak, 1978). Therefore, it may be appropriate to interpret $\epsilon_k$ and $\Psi_k$ as only auxiliary variables that are used for constructing the ground-state energy and density, because there are many materials for which interpretation of $\epsilon_k$ as the excitation energy is wrong and leads to erroneous results. Examples include band gaps in semiconductors and insulators. This should not be surprising because the DFT scheme has really been designed for computing the ground states only.

Alternatively, one can use an indirect method for carrying out the minimization of the total energy functional. Among the most successful such indirect methods is the QMD method of Car and Parrinello (1985). We will describe this method in Section 9.3 after describing the classical MD technique.

## 9.1.1   Local-Density Approximation

The local-density approximation assumes that the exchange-correlation energy functional is purely local. In an inhomogeneous material, the exchange-correlation potential at any point $\mathbf{r}$ depends not only on the electron density at this point, but also on its variations in the vicinity of $\mathbf{r}$. Thus, one may develop a gradient expansion in which $E_{Xc}$ depends on $\rho(\mathbf{r})$, $\nabla \rho(\mathbf{r})$, $\nabla[\nabla \rho(\mathbf{r})]$, $\cdots$ However, if such an expansion is used in the DFT, the required computations will be very difficult to carry out in a reasonable amount of time. In addition, simple gradient expansions

are rather badly behaved, and therefore one must be careful in using a gradient expansion, an issue that will be discussed in the next section. Such difficulties provided the prime motivation for developing the local-density approximation which ignores all corrections to the exchange-correlation energy at any point $\mathbf{r}$ due to heterogeneities in its vicinity. Instead, one makes the *ansatz* that

$$E_{Xc}[\rho(\mathbf{r})] = \int \epsilon_{Xc}[\rho(\mathbf{r})]\rho(\mathbf{r}) \, d^3r, \tag{18}$$

where $\epsilon_{Xc}[\rho(\mathbf{r})]$ is the exchange-correlation energy per particle of a homogeneous electron gas at density $\rho(\mathbf{r})$. In other words, if the exchange-correlation potential is assumed to be a local *function*, as opposed to be a *functional*, of the density with the same value as for a uniform electron gas, one obtains the local-density approximation. The non-uniform gas at $\mathbf{r}$ is therefore treated as if it were part of a uniform electron gas of constant density.

The local-density approximation is very accurate if $\rho(\mathbf{r})$ does not vary too rapidly, but is also surprisingly accurate when the distribution of electrons is strongly inhomogeneous, such as at surfaces and in molecules. To calculate $\epsilon_{Xc}$, the exchange effects are separated out from the dynamic correlations effects that are due to the Coulomb interaction between the electrons. The exchange part, commonly denoted by $\epsilon_x$, is given by

$$\epsilon_x[\rho(\mathbf{r})] \sim -c\rho^{1/3}(\mathbf{r}), \tag{19}$$

where $c$ is a constant. Equation (19) can be derived based on calculations for a homogeneous electron gas. For open shell systems, the spin-up and -down densities $\rho_\uparrow$ and $\rho_\downarrow$ are usually taken into account as two independent densities in the exchange-correlation energy. In this version of the local-density approximation, which is called the *local spin density approximation*, $E_x$, the exchange part of the energy, is given by (Jones and Gunnarsson, 1989)

$$E_x(\rho_\uparrow, \rho_\downarrow) = -\frac{3}{(3/4\pi)^{1/3}} \int [\rho_\uparrow^{4/3}(\mathbf{r}) + \rho_\downarrow^{4/3}(\mathbf{r})] \, d^3r, \tag{20}$$

which is obtained by inserting (19) into (18). Various schemes have also been proposed for taking into account the effect of the dynamic correlations (see, for example, Ceperley, 1978; Perdew and Zunger, 1981, and references therein).

In general, the local-density approximation gives a single well-defined global minimum for the energy of a non-spin-polarized system of electrons in a fixed ionic potential, implying that any energy minimization scheme (such as simulated annealing described in Chapter 3 of Volume I) will locate the global energy minimum of the electronic system. Magnetic materials, on the other hand, are expected to have more than one local minimum in the electronic energy. In this case, performing total-energy calculations will be very difficult, because the global energy minimum could be found only by sampling the energy functional over a large region of phase space. Considering the nature of local-density approximation, its success in providing quantitative predictions for many materials is remarkable.

## 9.1.2   Generalized Gradient Approximation

Another popular approximation to the exchange-correlation energy that is now widely used is the generalized gradient approximation according to which

$$E_{Xc}[\rho(\mathbf{r})] = \int f(\rho, \nabla \rho) d^3 r, \tag{21}$$

where $f$ is an analytic parameterized function that is fitted in such a way that $E_{Xc}$ satisfies several exact requirements. Many functional forms for $f$ have been suggested, a list of which is too long to be given here. Here, we only describe a relatively simple one due to Perdew *et al.* (1996). To begin with, the total electron density $\rho$ is written as the sum of up and down spin densities, $\rho(\mathbf{r}) = \rho(\mathbf{r})_{\downarrow} + \rho(\mathbf{r})_{\uparrow}$, $E_{Xc} = E_X + E_c$, and $\epsilon_{Xc} = \epsilon_X + \epsilon_c$. Perdew *et al.* (1996) proposed that

$$E_c(\rho) = \int \rho(\mathbf{r})[\epsilon_c(r_c, \zeta) + H(r_s, \zeta, t)] d^3 r, \tag{22}$$

where $r_s$ is called the Seitz radius (such that $\rho = 3/4\pi r_s^3$), $\zeta = (\rho_\uparrow - \rho_\downarrow)/\rho$, and $t = |\nabla \rho|/(2k_s \phi \rho)$ is a dimensionless density gradient. Here, $\phi(\zeta) = \frac{1}{2}[(1 + \zeta)^{2/3} + (1 - \zeta)^{2/3}]$, and $k_s$ is the Thomas-Fermi screening wave number. Moreover, $\epsilon_c(r_s, \zeta) = (e^2/a_0)\phi^3[\gamma \ln(r_s/a_0) - \omega]$, with $\gamma = (1 - \ln 2)/\pi^2 \simeq 0.031091$ and $\omega \simeq 0.046644$.

The function $H$ is selected in such a way that $E_{Xc}$ satisfies several rigorous constraints. The proposed form for $H$ is given by

$$H = \gamma \phi^3 \left( \frac{e^2}{a_0} \right) \ln \left[ 1 + \frac{\beta}{\gamma} t^2 \left( \frac{1 + At^2}{1 + At^2 + A^2 t^4} \right) \right], \tag{23}$$

with

$$A = \frac{\beta}{\gamma} \left[ \exp \left( -\frac{a_0 \epsilon_c}{\gamma \phi^3 e^2} \right) - 1 \right]^{-1}. \tag{24}$$

In these equations, $e$ is the electron charge, $a_0 = \hbar^2/me^2$, and $\beta \simeq 0.066725$.

The $E_X$ portion of $E_{Xc}$ is written as

$$E_X = \int \rho(\mathbf{r})\epsilon_x(\rho) F_X(s) d^3 r, \tag{25}$$

where $\epsilon_X(\rho) = -3e^2(3\pi^2 \rho)^{1/3}/(4\pi)$, and $s = (r_s/a_0)^{1/2}\phi t/c$ is, similar to $t$, a dimensionless density with $c \simeq 1.2277$. The function $F_X(s)$ is given by

$$F_X(s) = 1 + \kappa - \frac{\kappa}{1 + \mu s^2/\kappa}, \tag{26}$$

with $\kappa = 0.804$ and $\mu = \beta \pi^2/3 \simeq 0.21951$. Let us emphasize that, although we only provided here the numerical values of the several constants that appear in these equations, they are in fact related to fundamental physical constants, and do not represent numerical fits to some experimental data. These equations have been shown to provide very accurate predictions for atomization energies of many molecules, and therefore are widely used in the DFT computations.

Figure 9.1. A supercell for a bulk solid with a point defect at its center. The cell is enclosed by dashed lines, with the rest being its periodic images (after Payne *et al.*, 1992).



### 9.1.3 Nonperiodic Systems

While the DFT computations are carried out most conveniently for periodic systems, they run into difficulty when the system contains some sort of a defect. Such systems have attracted wide attention, especially when they contain a large number of atoms, in which case they are referred to as *mesoscopic systems*. Examples of such systems are abundant and include scanning tunneling microscope tip and surface, grain boundaries, quantum dots and wires, and biological macromolecules. To study such systems with the DFT a periodic *supercell* is used, an example of which for a system with one defect is shown in Figure 9.1. The supercell contains the defect which is surrounded by a region of bulk crystal. Periodic boundary conditions are applied to the supercell, so that it is replicated throughout space. Thus, due to periodicity, one actually computes the energy per unit cell of a crystal containing an array of defects, rather than the energy of a crystal containing a single defect. To prevent the defects in the neighboring cells to interact appreciably with each other, one must include enough bulk solid in the supercell.

Another case for which a supercell must be used in the computations is when a surface is only partially periodic. For example, it may have periodicity in its own plane, but not in the direction perpendicular to its plane. The supercell for such systems is shown in Figure 9.2. It contains a crystal slab and a vacuum region, and is repeated over all space, so that the total energy of an array of crystal slabs is calculated. In order to ensure that the results of the computations are true representative of an isolated surface, the vacuum region must be wide enough that faces of adjacent crystal slabs do not interact across the vacuum. The crystal slab must also be thick enough that the two surfaces of each crystal slab do not interact through the bulk crystal. Even molecules can be studied in this fashion; see Joannopoulos *et al.* (1991).

### 9.1.4 Pseudopotential Approximation

For periodic systems, one takes advantage of Bloch's theorem to simplify the computations. According to this theorem (Ashcraft and Mermin, 1976) in a periodic solid material each electronic wave function can be written as the product of a

FIGURE 9.2. A supercell for a surface of a bulk solid. The cell is enclosed by dashed lines, with the rest being its periodic images (after Payne *et al.*, 1992).

cell-periodic part and a wavelike part:

$$\Psi_i(\mathbf{r}) = \exp(i\mathbf{k} \cdot \mathbf{r}) f_i(\mathbf{r}). \tag{27}$$

The cell-periodic part of $\Psi_i$ can be expanded using a basis set that consists of a discrete set of plane waves with wave vectors that are reciprocal lattice vectors of the crystal. Therefore,

$$f_i(\mathbf{r}) = \sum_{\mathbf{G}} c_{i,\mathbf{G}} \exp(i\mathbf{G} \cdot \mathbf{r}), \tag{28}$$

where the reciprocal lattice vectors $\mathbf{l}$ are defined by, $\mathbf{G} \cdot \mathbf{l} = 2\pi m$ for all $\mathbf{l}$, with $\mathbf{l}$ being a lattice vector of the crystal and $m$ an integer. Therefore,

$$\Psi_i(\mathbf{r}) = \sum_{\mathbf{G}} c_{i,\mathbf{k}+\mathbf{G}} \exp[i(\mathbf{k} + \mathbf{G}) \cdot \mathbf{r}]. \tag{29}$$

In practice, this series is truncated to include only plane waves that have kinetic energies less than a cutoff energy. This introduces some error into the computations, but the error decreases with increasing energy cutoff. A suitable value of the energy cutoff can be selected by an optimization technique (Rappe and Joannopoulos, 1991), since the value of the cutoff is not included in the theory itself. Substitution of Eq. (29) in (16) converts the Kohn–Sham equations into a relatively simple set of equations in terms of the coefficients $c_{i,\mathbf{k}+\mathbf{G}}$ and the eigenvalues $\epsilon_k$ in which the kinetic energy is diagonal. Even for systems which contain aperiodicity, such as those with defects, use of a supercell makes the system amenable to this type of analysis. In any event, the solution of the Kohn–Sham equations, when written in terms of the eigenvalues of the coefficients $c_{i,\mathbf{k}+\mathbf{G}}$, can be obtained by diagonalization of the associated matrix. The size of the matrix is dictated by the choice of the cutoff energy, and can be prohibitively large if the system contains both valence and core electrons. The problem can be overcome by use of the *pseudo-potential approximation* which we now describe.

It is well known that many properties of solid materials depend on the valence electrons much more strongly than on the core electrons. This fact is exploited in the pseudo-potential approximation by which one removes the core electrons and replaces them and the strong ionic potential by a weaker *pseudo-potential* that acts on a set of *pseudo-wave functions*, rather than the true valence wave functions (Vanderbilt, 1990). The valence wave functions oscillate strongly in the region occupied by the core electrons, due to the strong ionic potential in this region. Ideally, the pseudo-potential must be constructed in a way that its scattering properties for the pseudo-wave functions are identical to those of the ion and the core electrons for the valence wave functions. But, this must be done in a way that the pseudo-wave functions have no radial nodes in the core region, where the total phase shift produced by the ion and the core electrons will be greater by $\pi$, for each node that the valence functions had in the core region, than the phase shift produced by the ion and the valence electrons. The phase shift generated by the ion core is different for each angular momentum component of the valence wave function, and therefore the scattering from the pseudo-potential must depend on angular momentum. Outside the core region, the two potentials and their scatterings are identical. The most general form for a pseudo-potential is then given by

$$U_{ps} = \sum_{lm} |lm\rangle u_l \langle lm|, \tag{30}$$

where $|lm\rangle$ are the spherical harmonics, and $u_l$ is the pseudo-potential for angular momentum $l$. Acting on the electronic wave function with this operator decomposes the wave function into spherical harmonics, each of which is then multiplied by the relevant pseudo-potential $u_l$.

A pseudo-potential that does not depend on the angular momentum components of the wave function is called a *local pseudo-potential*, which is a function only of the distance from the nucleus. While it is possible to generate arbitrary, pre-determined phase shifts for each angular momentum state with a local potential, in practice there are limits to the amount that the phase shifts can be adjusted for the different angular momentum states, since one must preserve the smoothness and weakness of the pseudo-potential, without which it becomes difficult to expand the wave functions using a reasonable number of plane-wave basis states.

In order for the exchange-correlation energy to be represented accurately, the pseudo and the true wave functions must be identical outside the core region. If one adjusts the pseudo-potential to ensure that the integrals of the squared ampli-tudes of the true and pseudo-wave functions *inside* the core region are identical, the equality of the pseudo and true wave functions outside the core region is guaranteed. Pseudo-potentials that possess this property were first constructed by Starkloff and Joannopoulos (1977), and have been shown to be highly accurate for heavy atoms. Moreover, Hamann *et al.* (1979) pointed out that, a match of the pseudo and real wave functions outside the core region also ensures that the first-order energy dependence of the scattering from the ion core is correct, so that the scattering is accurately described over a wide range of energy. A technique for

constructing pseudo-potentials that corrects even the higher-order energy dependence of the scattering was introduced by Shirley *et al.* (1989). In the best case scenario, one should develop non-local pseudo-potentials that use a different potential for each angular momentum component of the wave function, and efforts in this direction (see, for example, Shirley *et al.*, 1989, and references therein) have also been fruitful.

The following procedure is typically used for constructing an ionic pseudo-potential. The computations are carried out with a periodic unit cell (for non-periodic systems, see Section 9.1.3). All-electron computations are carried out for an isolated atom in its ground state and also some excited states, using a given equation for the exchange-correlation density functional, which result in valence electron eigenvalues and valence electron wave functions for the atom. A pseudo-potential with a few parameters is then selected and the parameters are adjusted in such a way that a pseudo-atom calculation using the same form for exchange-correlation as in the case of the all-electron atom yields both pseudo-wave functions that match the valence wave functions beyond some cutoff radius $r_c$, and pseudo-eigenvalues that are equal to the valence eigenvalues. The ionic pseudo-potential so obtained is then utilized, without further modifications, for any environment of the atom. The electronic density in any new environment of the atom is then determined utilizing both the ionic pseudo-potential obtained in this way and the same form of exchange-correlation functional employed in the construction of the ionic pseudo-potential.

One advantage of pseudo-potential method is that, by replacing the true ionic potential by a weaker pseudo-potential, one expands the electronic wave functions using far fewer plane-wave basis states that would be necessary if the full ionic potential were to be used. Moreover, the rapid oscillations of the valence wave functions in the cores of the atoms are removed, and the small core electron states are no longer present. Another advantage of the pseudo-potential approximation is that fewer electronic wave functions must be calculated.

Total-energy pseudo-potential computations require significant amounts of computer time, even when the number of atoms in the unit periodic cell is small. Moreover, the computational time always increases with the number of atoms in the unit cell, and therefore use of the most efficient computational method is crucial to the success of this method. The quantum molecular dynamics method devised by Car and Parrinello (1985) has improved the efficiency of these computations dramatically, hence allowing one to simulate systems with a large number of atoms. As already mentioned above, since the essence of the total-energy pseudo-potential computations is finding the electronic states that minimize the Kohn–Sham energy functional, one can also attempt to directly minimize the energy functional by a method such as the simulated annealing method described in Chapter 3 of Volume I, or by the conjugate-gradient technique which is used for minimization problems. After we describe the classical MD simulation, the QMD method of Car and Parrinello and the conjugate-gradient technique for directly minimizing the Kohn–Sham energy functional will be described.

## 9.2   Classical Molecular Dynamics Simulation

Historically, the first true MD simulation seems to have been carried out by Alder and Wainwright (1957), who studied a system with only a few hundreds hard-sphere particles and discovered a fluid-solid phase transition. At the same time, Wood and Parker (1957) investigated the properties of simple fluids using the Monte Carlo method. Rahman (1964) was apparently the first to carry out MD simulations using the Lennard–Jones potential (see below for a description of this potential). Unlike Alder and Wainwright, Rahman's work was the first to involve particles with smoothly varying potentials. He computed the diffusion coefficient and pair-correlation function for liquid argon and showed them to be in very good agreement with the experimental data. These three pioneering works opened the way for simulation, and hence understanding, of many-body systems. Later work by Verlet (1967), whose method helped MD simulations to become much more efficient (see below), by Alder and Wainwright (1969) who discovered, unexpectedly, an algebraic long-time tail in the velocity autocorrelation functions of hard sphere (a discovery that intensified further the interest in MD simulations), and by Rahman and Stillinger (1971), who addressed simulations of such complex molecules as liquid water, firmly established the classical MD simulations as an every-day tool of studying fluids and materials.

The 1970s witnessed further improvements in methodologies and algorithms for MD simulations. For example, Evans and Murad (1977) succeeded in developing an algorithm for computing molecular rotations [Ciccotti *et al.* (1982) made further improvement to this method], and Bennett (1976) and Torrie and Valleau (1977) developed efficient methods for measuring the free energies (see also Frenkel and Ladd, 1984), and so on. Much more progress was made in the 1980s when Andersen (1980) and Parinello and Rahman (1981) developed methods for carrying out MD simulations under constant pressure and constant temperature, and Nosé (1984) developed equations for simulating constant-temperature MD simulations by introducing additional degrees of freedom [his equations were simplified by Hoover (1985); see below]. At the same time, the advent of vector computers further motivated the search for methods that could take advantage of vectorization techniques, especially those that could be used with the Verlet algorithm. All of these advances took MD simulations to a stage where, by mid 1980s, they could be used for studying *non-equilibrium* systems. Abraham *et al.* (1984) studied, using MD simulations, the incommensurate phase of Krypton on graphite using more than 160,000 atoms (a "revolution" for its time), and Car and Parinello (1985) succeeded in combining MD simulations and electronic structure calculations (see Section 9.4). The first million-particle MD simulations were carried out by Swope and Andersen (1990) who studied homogeneous nucleation of crystals in a super-cooled atomic liquid (i.e., below its freezing point). Their study was important, not only because of the large number of atoms that had been used, but also because it showed that certain physical phenomena can be reproduced in the MD simulations only when the size of the system is large enough.

Generally speaking, two types of classical MD simulations can be carried out. Equilibrium MD simulations are suited for systems that can, in principle, be treated by statistical mechanics. This type of MD simulations can yield equilibrium properties of materials. Non-equilibrium MD techniques are appropriate for systems that are under the influence of an external driving force, and are most suitable for computing the transport properties of a system. We first discuss general concepts and ideas of MD simulations that are applicable to both the equilibrium and non-equilibrium methods, after which we describe those aspects of non-equilibrium MD technique that are different from the equilibrium methods.

### 9.2.1 Basic Principles

Molecular dynamics simulation of any phenomenon consists of integration of Newton's equation of motion for a system of $N$ particles that represent the material or the system under study. Therefore, the MD method is a way of simulating the behavior of a system as it evolves with time since, unlike the Monte Carlo (MC) method, in the MD simulations the system moves along its physical trajectory. The main advantage of the MD method over the MC technique is that, not only it provides a method for computing the static properties of a system, but also allows one to calculate and study the dynamical properties, including dynamic fracture propagation in materials that are of interest to us in this book.

Suppose that we have a collection of $N$ particles in a simulation cell with dimensions $L_x, L_y$ and $L_z$. The particles interact with each other, and for simplicity we assume for now that the interaction force can be written as a sum over pair forces $\mathbf{F}(r)$, the magnitude of which depends only on $r$, the distance between the particle pairs. Thus, the force acting on any particle $i$ is given by

$$\mathbf{F}_i(\mathbf{r}^N) = \sum_{j=1, \ j \neq i}^{N} F(| \ \mathbf{r}_i - \mathbf{r}_j \ |)\hat{\mathbf{r}}_{ij}, \tag{31}$$

where $\mathbf{r}^N = \{\mathbf{r}_1, \mathbf{r}_2 \cdots, \mathbf{r}_N\}$ is the position coordinates of all the particles, and $\hat{\mathbf{r}}_{ij}$ is a unit vector along $\mathbf{r}_i - \mathbf{r}_j$, pointing from particle $i$ to $j$. Then, the equation of motion for particle $i$ is given by

$$m_i \frac{d^2 \mathbf{r}_i(t)}{dt^2} = \mathbf{F}_i(\mathbf{r}^N) + \mathbf{F}_e, \tag{32}$$

where $m_i$ is the atomic mass of particle $i$, and $\mathbf{F}_e$ represents all the external forces that are imposed on the system, either by nature (for example, the gravitational force) or by the experimentalist (for example, an external pressure gradient applied to the system). Molecular dynamics consists of writing Eq. (32) for all the $N$ particles of the system and integrating them numerically and simultaneously (since the equations are coupled through the force $\mathbf{F}_i$). The solution of this set of equations describes the time evolution of the system. If the forces between the particles depend only on their relative positions, then the system's energy and momentum are *automatically* conserved.

To check whether the system has reached a steady state, if one exists, quantities such as the kinetic energy are computed and their variations with the time are monitored. When these quantities no longer vary with the time appreciably, then attainment of the steady state has been confirmed. The time to reach a steady state depends on the initial state of the system, and how far it is from its steady state.

Although the MD approach is, in principle, a rigorous method, in practice, and similar to almost all computational strategies for studying a phenomenon, it is only an approximate technique. Thus, it should be used with considerable care. Some of the problems that one must pay particular attention to are as follows.

(1) The interaction potentials between the particles are not, in almost all cases, known, and therefore one must use approximate expressions for describing these potentials. In principle, quantum-mechanical calculations can be used for determining these forces, but such computations can be subject to errors. Typically, the interaction potentials or forces are written in terms of several parameters which are determined either by *ab initio* computations or by fitting the results to experimental data (see Section 9.7).

At atomic scale, the interactions are either of intra-molecular or inter-molecular type. We will discuss the intra-molecular interactions separately, and for now briefly describe the inter-molecular interactions. In many MD simulations the interaction potential between a pair of particles, the centers of which are a distance $r$ apart, is represented by the classical Lennard–Jones (LJ) potential. It was thought for a long time that this potential is too simple to mimic the behavior of real materials, particularly brittle ones. However, using what the physics Nobel Laureate R. P. Feynman emphasized nearly 40 years ago as a guide, the value of the LJ potential is in its universal nature, since according to Feynman the single most important statement describing our real world is that, *all things are made of atoms, little particles that move around in perpetual motion, attracting each other when they are a little distance apart, but repelling upon being squeezed into one another,* which is precisely how the LJ and similar potentials have been constructed for representing real atoms, fluids, and materials. For use in MD simulation and modeling of many properties of materials, including their dynamic fracture which is of interest to us in this chapter (see Section 9.8), a relatively general potential can be constructed, part of which is based on the classical LJ potential. This potential is given by (Holian *et al.*, 1991)

$$U(r) = \begin{cases} 4\epsilon \left[ \left(\dfrac{\sigma}{r}\right)^{12} - \left(\dfrac{\sigma}{r}\right)^{6} \right], & r < r_i \\ -a_1(r_c^2 - r^2)^2 + a_2(r_c^2 - r^2)^3, & r_i \leq r \leq r_c \\ 0, & r > r_c \end{cases} \tag{33}$$

where $\epsilon$ is the energy parameter of the potential (the maximum energy of attraction between a pair of molecules), or the LJ well depth, and $\sigma$ is the size parameter (or the distance at which the LJ potential passes through zero),

also called the *collision diameter*. Note that $\sigma$ is not the same as molecular diameter of the molecules, although the two quantities are usually close to each other. The 12-6 part of $U(r)$ is the classical LJ potential in which the $r^{-12}$ term represents a hard-core or repulsive potential, while the $r^{-6}$ term is its attractive part. The second equation (for $r_i \leq r \leq r_c$) has been added for accommodating the difference between brittle and ductile materials. Here $r_i$ is the inflection point in the potential, and

$$r_c^2 = r_i^2 \left\{ 5 - 5 \left[ 1 - \frac{1}{25} \left( 9 - \frac{24U_{LJ}(r_i)}{r_i U'_{LJ}(r_i)} \right) \right]^{1/2} \right\}, \tag{34}$$

$$a_1 = \frac{5r_i^2 - r_c^2}{8r_i^3(r_c^2 - r_i^2)} U'_{LJ}(r_i), \tag{35}$$

$$a_2 = \frac{3r_i^2 - r_c^2}{12r_i^3(r_c^2 - r_i^2)^2} U'_{LJ}(r_i). \tag{36}$$

If the intention is to study the properties of materials via MD simulations (without studying their fracture), then typically only the $12 - 6$ part of the potential is used. The force $\mathbf{F}(r)$ between the particles is then given by, $\mathbf{F}(r) = -\nabla U(r)$ for $r \leq r_c$; of course, $\mathbf{F}=\mathbf{0}$ for $r > r_c$. If the system contains solid walls, then the interactions between the materials' atoms and those of the walls must also be taken into account. A well-known potential due to Steele (1973) has been used in many simulations (although Steel's potential represents some sort of a mean-field approximation, as it assumes that the wall is smooth and structureless):

$$U_w = 2\pi\rho_w\epsilon_w\sigma_w^2\Delta \left[ \frac{2}{5} \left( \frac{\sigma_w}{z} \right)^{10} - \left( \frac{\sigma_w}{z} \right)^4 - \frac{\sigma_w^4}{3\Delta(z + 0.61\Delta)^3} \right], \tag{37}$$

where $\epsilon_w$ and $\sigma_w$ are the energy and size parameters that characterize the interactions between the atoms in the system and those of the walls, $z$ is the vertical distance from the wall, $\rho_w$ is the density of the wall's atoms, and $\Delta$ is the distance between the atomic layers within the wall.

The accuracy of the MD method depends to a large extent on the accuracy of the interaction potentials used in the simulation. Lennard–Jones type potentials are too simple to represent complex atoms and molecules. To remedy the situation, one can fit the size and energy parameters of the LJ potential, so that certain predictions of the MD simulations fit, in some sense, the experimental data. For example, these parameters can be estimated from the properties of fluids at the critical point, liquids at the normal boiling point, or solids at the melting point. This method has been reasonably accurate for relatively simple molecules and might, in some cases, provide qualitative insight into the behavior of materials. However, many materials of technological interest have strong and specific chemical interactions that cannot be described by simple, pairwise additive potentials, such as the LJ potential. Thus, more sophisticated

computations must be used for determining the interaction potentials. We will come back to this issue later in Section 9.7

(2) Molecular dynamics suffers from the same problem that every computer simulation of a physical system suffers from, namely, while the intention is to simulate a real system which, at least in atomic units, is very large, the simulated system is of finite size. The finite-size effect can be particularly severe if there are correlations in the system (as is the case with almost any physical system). If the correlation length is much smaller than the linear size of the system, finite-size effects do not pose any significant problem. If the correlation length is much larger than the linear size of the system, then one can use finite-size scaling (see Chapter 2 of Volume I) for extrapolating the results for a system of finite size to one with an infinite size (although finite-size scaling is usually used for second-order phase transitions, and the MD methods are not used very often for studying such transitions, as the computations would be very intensive). The intermediate case, in which the correlation length is not too large or too small, is usually addressed by using the periodic boundary conditions, since the finiteness of the system is manifested through its boundaries (an infinite system does not have any boundary!). Using periodic boundary conditions means that the finite simulation cell is embedded in an infinite system, such that it is surrounded by replicas of itself on all sides. In that case,

$$\mathbf{F}(\mathbf{r}_i - \mathbf{r}_j) = \sum_{\mathbf{n}} \mathbf{F}\left(\mid \mathbf{r}_i - \mathbf{r}_j + \sum_{k=1}^{3} \mathbf{V}_k n_k \mid\right), \tag{38}$$

where $\mathbf{V}_k$ are vectors along the edges of the rectangular simulation cell. The first sum on the right-hand side of Eq. (38) is over all vectors $\mathbf{n} = (n_1, \cdots, n_k)$. The force $\mathbf{F}$ is along the line connecting particle $i$ and the image particle $\mathbf{r}_j - \sum_{k=1}^{3} \mathbf{V}_k n_k$. In principle, calculating terms of this infinite sum until it converges to a well-defined value is a difficult task, but methods have been developed for computing such sums. Use of periodic boundary conditions also has a negative side effect: In a periodic system the angular momentum is *not* conserved, since the periodic boundary conditions break the spherical symmetry of the interactions.

(3) When time-averaged properties are calculated, the averaging is clearly carried out over a finite-time period. There is, however, a limitation in time as a result of the finite number of integration steps that one can carry out. Finite size of the system also limits the time, especially if the particles (in, for example, simulation of liquids or gases) travel more than half the linear size of the simulation cell.

(4) In any MD simulation, there is always a competition between the speed of the computations and their accuracy. Normally, as the size of the time step increases, so also does the inaccuracy in the simulation results. Therefore, for any MD simulation, there is an optimal choice of the time step.

Having described the strengths and weaknesses of the classical MD simulation, let us now describe some basic issues that arise in such computations.

## 9.2.2  Evaluation of Molecular Forces in a Periodic System

From the computational view point, the most intensive part of any MD simulation is calculation of the forces between the particles which accounts for $70 - 90\%$ of the total time. Periodic boundary conditions create a problem for evaluation of the forces between the particles, since in a periodic system not only does a particle interact with other particles in the simulation cell, but also with those in the images of the system that surround it, and thus, in principle, one must sum over an infinite number of interactions. However, in many cases the force between two particles decays rapidly as the distance between them increases, and therefore the particles that are far from any given particle, whether they are in the simulation cell or in its images, do not contribute significantly. If the force between two particles can be ignored for distances that are larger than half the system's linear size, then one can use the *minimum image convention* according to which, for each particle in the system, one takes into account only the interactions with the nearest copy of each of the remaining particles, implying that each infinite sum over all the images is replaced by a single term. For example, for a cubic simulation cell, the minimum interaction distance is given by

$$(r_{ij})_{\min} = \min|\mathbf{r}_i - \mathbf{r}_j + V_k n_k|, \tag{39}$$

where the notations are the same as in Eq. (38). Of course, the potentials will no longer be analytic, but the discontinuities will not be important if the potential is small for distances that are larger than half the system's linear size.

In practice, what is done in most cases is cutting the interactions off at a distance $r_c$. Typically, $r_c$ is set to be a multiple of the effective molecular diameter of the largest atom in the simulations, and therefore it is usually smaller than half the system's linear size. If this approach is taken, then at every step of the integration one must check, for any particle $i$, the distances of all other particles from $i$ to see whether they are at a distance larger or smaller than $r_c$. This search constitutes one the three important tasks in any MD simulation (the other two being computing the forces, and integrating the equations of motion). There are efficient methods of doing this which we will discuss below. We must point out, however, that cutting off the interaction potentials violates energy conservation, although if $r_c$ is selected carefully, the effect will be small. Moreover, by shifting the interaction potentials one can avoid violation of energy conservation altogether by writing

$$U(r) = \begin{cases} U_o(r) - U_o(r_c) & \text{if } r \leq r_c \\ 0 & \text{if } r > r_c \end{cases} \tag{40}$$

where $U_o(r)$ represents the original interaction potential to be used. However, this shift does not affect the force resulting from the shifted potential; it remains discontinuous at $r_c$. In order to make the force also continuous at the cutoff point, we write

$$U(r) = \begin{cases} U_o(r) - U_o(r_c) - \dfrac{d}{dr}[U_o(r_c)(r - r_c)] & \text{if } r \leq r_c \\ 0 & \text{if } r > r_c \end{cases} \tag{41}$$

This algorithm was first suggested by Stoddard and Ford (1973). The actual number of interacting particles (i.e., those that are within a sphere of radius $r_c$, centered at the center of a given particle) is a function of the molecular density. For example, in the simulation of a typical liquid state using the LJ potential, the computation of a single pair-interaction requires about $30 - 40$ floating-point operations. Therefore, a complete force calculation requires of the order of 2,000 floating-point operations per particle, still computationally intensive, but much more efficient than a full $N$-body calculation. Let us point out that, the cut-and-shift procedure cannot be used if electric and gravitational forces are operative in the system, since they decay only as $1/r$. Such cases must be treated separately; see Section 9.2.8.

## 9.2.3   The Verlet and Leapfrog Algorithms

In a typical MD simulation, the time that the computer program spends for integrating the equations of motion is about 2–3% of the total time. However, accurate integration of the equations of motion is the most important part of the computations. Various methods have been proposed for achieving this goal, a detailed discussion of which can by itself be the subject of a minireview. One heavily-used procedure is due to Verlet (1967). According to his method, the algorithm for integrating the equation of motion for a single particle, which is subjected to a force **F** that depends only on the particle's position, is given by

$$\mathbf{r}(t + \Delta t) = 2\mathbf{r}(t) - \mathbf{r}(t - \Delta t) + \frac{\Delta t^2}{m^2}\,\mathbf{F}(\mathbf{r}), \qquad (42)$$

where $\mathbf{r}(t)$ is the position of the particle at time $t$, and $\Delta t$ is the integration time step (note that $t = n\Delta t$, where $n$ is an integer). The error per time step is usually $O(\Delta t^4)$, but in some cases can be as large as $O(\Delta t^2)$ which is still very accurate. As Eq. (42) indicates, the integration can proceed if the positions of the particles at two previous times ($t$ and $t - \Delta t$) are known. Thus, to begin the integration, the positions at $t = \Delta t$ are first computed from

$$\mathbf{r}(\Delta t) = \mathbf{r}(0) + \Delta t\mathbf{v}(0) + \frac{1}{2m}\Delta t^2\,\mathbf{F}[\mathbf{r}(0)] + O(\Delta t^3), \qquad (43)$$

which, together with the initial positions, provide us with the two previous positions that we need. The particles' velocities at any time $t$ are calculated from

$$\mathbf{v}(t) = \frac{\mathbf{r}(t + \Delta t) - \mathbf{r}(t - \Delta t)}{2\Delta t} + O(\Delta t^2), \qquad (44)$$

which is a standard finite-difference approximation to **v**. If periodic boundary conditions are used, then one must check whether any particle has left the simulation cell in the last integration step, in which case the particle must be translated back over a lattice vector $\mathbf{V}_k$ to keep it inside the cell. Clearly, the velocity evaluation step must be carried out *before* such a translation.

The Verlet algorithm, in its original form, is sometimes susceptible to error. A modified algorithm which is the exact arithmetic equivalent of the original Verlet

algorithm, but is far less susceptible to numerical errors, is the *leap-frog* algorithm, according to which one computes the velocity of a particle at midpoint between $t$ and $t + \Delta t$,

$$\mathbf{v}\left(t + \frac{1}{2}\Delta t\right) = \mathbf{v}\left(t - \frac{1}{2}\Delta t\right) + \frac{\Delta t}{m}\,\mathbf{F}[\mathbf{r}(t)], \tag{45}$$

from which the position of the particle is calculated,

$$\mathbf{r}(t + \Delta t) = \mathbf{r}(t) + \Delta t\,\mathbf{v}\left(t + \frac{1}{2}\Delta t\right). \tag{46}$$

In another modification of the Verlet algorithm, the so-called velocity-Verlet algorithm, one calculates the position and velocity of a particle from

$$\mathbf{r}(t + \Delta t) = \mathbf{r}(t) + \Delta t\mathbf{v}(t) + \frac{\Delta t^2}{2m}\,\mathbf{F}[\mathbf{r}(t)], \tag{47}$$

$$\mathbf{v}(t + \Delta t) = \mathbf{v}(t) + \frac{\Delta t}{2m}\{\mathbf{F}[\mathbf{r}(t + \Delta t)] + \mathbf{F}[\mathbf{r}(t)]\}. \tag{48}$$

This algorithm is most stable with respect to the finite precision arithmetic, and requires no additional computations in order to calculate the velocities.

An important property of the Verlet algorithm and its leap-frog modification is that, the energy that one calculates when using these integration methods does not exhibit any drift in the total energy. This important and desirable stability is due to the fact that the Verlet algorithm is time-reversible, and therefore does not permit steady increase or decrease of the energy for periodic systems. An important time scale in the simulations is the so-called *Poincaré time scale*, which is the time after which a system that starts out with a random configuration returns to its initial configuration. The total time that one can integrate in any MD simulation is, however, much smaller than the Poincaré time, and therefore there is the possibility of having an increasing error in the calculated energy as the equations of motion are integrated for larger times. However, the Verlet algorithm has an additional property called *symplecticity*. This property gives rise to conserved quantities, and in particular Sanz-Serna (1992) showed that, with this property present in the integration procedure, the discrete analogue of the total energy (in numerical integrations one can compute only discrete analogues of the properties of interest) is *rigorously* conserved, and that the discrete analogue of the total energy deviates from its continuum (that is, actual) counterpart by an amount which is of the order of $O(\Delta t^k)$, where $k$ is some relatively large integer. Therefore, the Verlet algorithm and its leap-frog version hold the deviations of the energy bounded. Such desirable properties are the main reason for the popularity of the Verlet algorithm and its modifications. In contrast with the Verlet algorithm, the Gear method (Gear, 1971), which is a predictor-corrector technique, has lost its popularity, despite the fact that it requires only one force evaluation per time step, because it is not symplectic and thus can create major problems for those properties of the system that are supposed to be conserved and bounded.

## 9.2.4   Constant-Energy Ensembles

Since a MD simulation conserves the number $N$ of the particles and the system volume $\Omega$, then, if the energy $E$ and momentum are also conserved, the time averages of any physical property computed from this type of simulations will be equal to the averages in the microcanonical or the (NVE) ensemble. Let us describe how such MD simulations are carried out. One first specifies the number of particles and the interaction potentials between them, and assigns them their initial positions and momenta. Since the temperature of any system can be measured, and because one is usually interested in carrying out the simulations at a specified temperature, rather than a specified energy, one must "push" the system toward the desired temperature. There is an efficient way of achieving the desired temperature, which will be discussed below.

   If the LJ potential is used, then the sites of the FCC lattice are usually used as the initial position, since the FCC lattice represents the ground state of the LJ potential, although any other initial positions can also be utilized. The initial velocities $\mathbf{v}$ are drawn from a Maxwell–Boltzmann distribution with the specified temperature $T$:

$$f(v_x, v_y, v_z) = \left( \frac{m}{2\pi k_B T} \right)^{3/2} \exp(-mv^2/2k_B T), \tag{49}$$

where $v^2 = v_x^2 + v_y^2 + v_z^2$, and $k_B$ is the Boltzmann's constant. This is done by drawing the $x$, $y$ and $z$ velocity components for each particle from a Gaussian distribution. After generating the initial momenta, the average momentum per particle $\langle \mathbf{p} \rangle$ is computed and subtracted from the individual momentum $\mathbf{p}_i$ of each particle. This ensures that the total initial momentum of the system is zero. Once the initial configuration of the system has been prepared, the integration of Newton's equation of motion begins.

## 9.2.5   Constant-Temperature Ensembles

The temperature of a system is a property that can be easily measured and controlled, and therefore it is often desired to perform an experiment under a fixed temperature condition. However, in MD simulations it is often very difficult to fix the temperature of the system at the very beginning. Since the temperature of an infinite system is proportional to the average kinetic energy per degree of freedom, with a proportionality constant $\frac{1}{2}k_B$, this quantity is used in MD simulations (of finite systems) to fix the temperature. One procedure for obtaining a fixed temperature $T_f$ is to calculate a rescaling factor $\alpha$ given by

$$\alpha = \left[ \frac{3k_B T_f (N-1)}{\sum_{i=1}^{N} m_i v_i^2} \right]^{1/2}, \tag{50}$$

where $m_i$ and $v_i$ are the mass and magnitude of velocity of particle $i$, respectively. The velocity $\mathbf{v}_i$ of each particle $i$ is then rescaled by $\alpha$, i.e., $\mathbf{v}_i(t) \to \alpha \mathbf{v}_i(t)$. The rescaling of the velocity also necessitates rescaling of the time $t$, which will be

discussed shortly. Each time this rescaling is done, the actual temperature of the system changes and gradually approaches the given fixed temperature $T_f$. This method can actually be derived by imposing a constant kinetic energy through a Lagrange multiplier term added to the Lagrangian of the isolated system (Haile and Gupta, 1983). However, although the rescaling procedure can be derived, due to certain assumptions used in the derivation, it still represents some sort of an *ad hoc* method. This has motivated the development of more systematic methods for achieving a given temperature and thus being able to perform MD simulations at constant temperature. One particularly effective method is based on introducing an extra force acting on the particles. The force is frictional in nature, is assumed to be proportional to the velocity of the particles, and therefore affects the kinetic energy of the system (and hence the system's temperature) in a direct way. Thus, the equation of motion for the $i$th particle is written as

$$m\frac{d^2\mathbf{r}_i}{dt^2} = F_i(\mathbf{R}) - C_f(R, dR/dt)\frac{d\mathbf{r}_i}{dt}, \tag{51}$$

where the friction parameter $C_f(R, dR/dt)$ is assumed to be the same for all the particles. The sign of $C_f$ depends on whether heat is added to or extracted from the system; in the former case $C_f < 0$ while in the latter case $C_f > 0$. Various equations have been suggested for $C_f$. The best-known equation was proposed by Nosé (1984) and simplified by Hoover (1985), according to which

$$\frac{dC_f}{dt} = \frac{\sum_i m_i v_i^2 - 3Nk_B T_f}{f}, \tag{52}$$

where $f$ is a parameter that must be selected carefully. In effect, in order to reach the desired temperature, the system is connected to a heat bath to exchange heat with it so that it can reach the intended temperature, and the parameter $f$ represents a coupling between the system and the heat bath. The Nosé-Hoover method yields precise canonical distribution for the particles' positions and momenta.

Despite many desirable features, the Nosé-Hoover method does have some shortcomings, the most serious of which is that the coupling parameter $f$ must be chosen. Moreover, although Cho and Joannopoulos (1992) showed that, for LJ fluids at high temperatures, the canonical distribution is reproduced correctly, Holian *et al.* (1995) demonstrated that if the temperature is lowered, it begins to oscillate with an amplitude which is much larger than the standard deviations expected in the canonical ensemble.

Many modifications of the Nosé-Hoover algorithm have been proposed. Chief among them is an approach proposed by Jellinek and Berry (1988) who suggested a generalization of the Nosé Hamiltonian involving multiplicative scaling of coordinates, momenta, and time. They showed that there are *infinitely many distinct* Hamiltonians (i.e., distinct dynamics) that possess *all* the properties of the Hamiltonian dynamics of Nosé-Hoover algorithm. Their model was analyzed in detail by Brańka and Wojciechowsky (2000).

## 9.2.6    Constant-Pressure and Temperature Ensembles

In addition to temperature, it is often desirable to be able to carry out MD simulations at a constant pressure [the (NPT) ensemble]. This is achieved by uniform isotropic volume changes caused by rescaling the atomic coordinates. Andersen (1980) proposed incorporating the volume $\Omega$ of the system into the equation of motion by rescaling the coordinates as

$$\mathbf{r}_i = \mathbf{r}'_i \Omega^{-1/3}, \tag{53}$$

where $\mathbf{r}'_i$ denotes the coordinates of the particles in the rescaled system. Furthermore, the momentum of the particle is also rescaled according to

$$\mathbf{p}_i = \mathbf{p}'_i(s\Omega^{1/3}), \tag{54}$$

where $s$ denotes a new dynamical variable that is equivalent to rescaling the real time, $dt = s(t')dt'$. In effect, the system is connected to a "piston" that can contract or expand it. The Hamiltonian of the system is given by

$$\mathcal{H} = \frac{1}{2}\sum_i \frac{\mathbf{p}_i^2}{ms^2\Omega^{2/3}} + \frac{1}{2}\sum_{ij, i\neq j} E_p(\mathbf{r}_{ij}\Omega^{1/3}) + \frac{p_b^2}{2f} + P\Omega$$
$$+ \frac{p_p^2}{2m_p} + (3N+1)k_BT \ln s, \tag{55}$$

where $p_p$ and $m_p$ are the "momentum" and "mass" of the "piston," respectively, $E_p$ is the potential energy, and $(3N + 1)$ represents the total number of independent momentum degrees of freedom of the system. The term involving $p_b$ represents the coupling to the heat bath (so that the temperature is also held fixed), while $P\Omega$ and $p_p^2/2m_p$ represent the work and kinetic energy arising from the connection of the system to the "piston." The governing equation for the dynamical variable $s(t)$ is simply

$$\frac{ds}{dt} = \frac{p_b}{f} = \frac{\partial \mathcal{H}}{\partial p_b}. \tag{56}$$

Note that if in Eq. (55) we do not rescale the lengths and momenta and delete the "piston" terms, then the modified equation also describes the system for the case in which only the temperature is held fixed, implying that in that case too, the time must be rescaled in the same way as in the present case.

Martyna *et al.* (1992) suggested a Hamiltonian slightly different from (55), and also equations of motion that differed from those used in the above method, in order to remove a minor problem from the above formulation, namely, the fact that the trajectories of the particles produced by the above method depend on the choice of the basis lattice vector.

## 9.2.7    Simulation of Rigid and Semirigid Molecules

Accurate molecular simulations depend largely on realistic representation of the molecules. Although representing atoms or molecules as simple LJ hard spheres

TABLE 9.2. The interaction parameters for $CO_2$.

| | | | |
|---|---|---|---|
| $\sigma_{OO}$(Å) | 3.027 | $\varepsilon_{OO}/k_B$ (K) | 74.8 |
| $\sigma_{CO}$(Å) | 2.922 | $\varepsilon_{CO}/k_B$ (K) | 44.8 |
| $\sigma_{CC}$(Å) | 2.824 | $\varepsilon_{CC}/k_B$ (K) | 26.2 |
| $q_O(e)$ | $-0.332$ | $q_C(e)$ | $+0.664$ |
| $\mathcal{M}$(cm$^2$) | $-14.4 \times 10^{-40}$ | $\ell_{CO}$(Å) | 2.324 |

may be adequate for predicting many qualitative features of experimental data, more sophisticated and realistic representation of the atoms is necessary if MD simulations are to provide quantitative predictions. As a relatively simple but important example, consider molecular modeling of $CO_2$. A molecular model of $CO_2$, which is much more realistic than a simple LJ hard sphere with effective parameters $\sigma$ and $\varepsilon$, was developed by Murthy *et al.* (1983), and further refined by Hammonds *et al.* (1990). In their model $CO_2$ is represented by a rigid linear molecule with quadrupole moments and three LJ sites. Three partial charges, $q_O$, $q_C$ and $q_O$, are used on the O–C–O sites, chosen so as to preserve the quadrupole moments of the molecule. The non-electrostatic interactions are still modeled as site-site LJ potentials, where the interaction sites are located on the three atoms. To fit the interaction parameters, three states of $CO_2$ (namely, solid, liquid, and gas) are used. The LJ parameters are fitted to the phonon frequencies and lattice energy of the solid, the thermodynamic properties of the liquid, and the second virial coefficient of the gas. All the parameters are given in Table 9.2, where $\sigma_{OO}$, $\sigma_{CO}$ and $\sigma_{CC}$ are the LJ size parameters, $\varepsilon_{OO}$, $\varepsilon_{CO}$ and $\varepsilon_{CC}$ are the LJ energy parameters between O–O, C–O, and C–C atoms in a $CO_2$ molecule, respectively, $\ell_{OC}$ is the distance between the O and C atoms, while $\mathcal{M}$ is the quadrupole moment. If the MD simulation involves $CO_2$ and other molecules, then the interaction between molecules $i$ and $j$ is expressed as

$$U(r_{ij}) = \sum_{m=1}^{3} \sum_{n=1}^{3} \left[ U_{LJ}(r_{im,jn}) + U_C(r_{im,jn}) \right], \qquad (57)$$

where $r_{im,jn}$ is the distance between the interacting pairs (site $m$ in molecule $i$ and site $n$ in molecule $j$), $U_{LJ}(r)$ is the-cut-and-shifted LJ potential described above, and $U_C(r)$ is the Coulomb potential given by

$$U_C(r) = \frac{q_{im}q_{jn}}{r}. \qquad (58)$$

During the MD simulations one must keep track of the coordinates of a $CO_2$ molecules. They can be represented by the vector $\mathbf{V}_j$, with $j = 1, 2, \cdots, N_{CO_2}$, where $N_{CO_2}$ is the total number of $CO_2$ molecules. The vector $\mathbf{V}_j$ contains three cartesian coordinates, $(r_x, r_y, r_z)$, determining the position of the molecule's center, and two coordinates determining its orientation. The orientation of $CO_2$ can be determined by a unit vector $\mathbf{e} = (e_x, e_y, e_z)$, directed along the axis of the molecule (where only $e_x$ and $e_y$ are independent, as $e_z = 1 - e_x - e_y$) and the angle that it makes with the surface of the system's walls.

We now describe how MD simulations of more complex materials, such as rigid or semi-rigid molecules, are performed. The motion of a rigid molecule consists of translations of the center of mass and rotations around this point. The force acting between two such molecules consists of atomic pair interactions between atoms that belong to the two different rigid molecules. One can also consider off-center interactions, but we neglect them here. There are several methods for treating rigid molecules, and we describe one of them which is based on imposing constraints on the system that depend on the spatial positions, but are independent of the velocities. Consider the Lagrangian of the system,

$$\mathcal{L}_0 = \int_{t_1}^{t_2} \left[ \frac{1}{2} \sum_i m_i \left( \frac{d\mathbf{r}_i}{dt} \right)^2 - \frac{1}{2} \sum_{i \neq j} E_p(\mathbf{r}_i - \mathbf{r}_j) \right] dt, \qquad (59)$$

where subscript 0 indicates that constraints have not been imposed yet. A constraint is imposed on the system through a Lagrange multiplier $\lambda(t)$, which is a function of time since the constraint should hold for all times. For example, for a rigid molecule that consists of two atoms, we impose the constraint that the distance $\ell$ between the two particles is always fixed. When such a constraint, which is often called the *bond constraint*, is imposed on the system, then the Lagrangian $\mathcal{L}$ of the new system (with the constraint) is given by

$$\mathcal{L} = \mathcal{L}_0 - \int_{t_1}^{t_2} \lambda(t) \left\{ [\mathbf{r}_1(t) - \mathbf{r}_2(t)]^2 - \ell^2 \right\} dt. \qquad (60)$$

$\lambda(t)$ is determined by requiring that the solution must satisfy the constraint. We discuss this shortly.

We are familiar with the concept of the backbone of a disordered medium. The same concept can be applied to a material at atomic scale. In this case some of the material's atoms form the backbone and are fixed by the bond constraints discussed above, while the remaining atoms are fixed by *linear constraints* which we discuss shortly. A good example is provided by a branched polymer in which the backbone is made of the multiply connected atoms. To identify the backbone of any molecular structure, some rules of thumb may be useful. For a planar molecular structure one can consider three non-colinear atoms as a backbone since they satisfy the bond constraint, while the rest of the atoms in the structure are constrained linearly. In a 3D molecular structure, four backbone atoms are subjected to six bond constraints with the remaining ones to a linear vector constraint each. A good example is provided by the linear molecule $CS_2$ (Thijssen, 1999), the motion of which is described by five positional degrees of freedom, two of which define the orientations of the molecules and three define the position of its center of mass. Without any constraint, the three atoms have nine degrees of freedom, but three of them are eliminated by the bond constraints, implying that we still have six degrees of freedom, instead of the required five. The inclusion of the un-needed degree of freedom adds to the computations and makes them inefficient. A better procedure is to fix only the distance between the two sulphur atoms by requiring that, $|\mathbf{r}_{S(1)} - \mathbf{r}_{S(2)}|^2 = \ell^2$, and to fix the position of the carbon atom by a linear

vector constraint which reads

$$\mathbf{r}_C = \frac{1}{2}\left[\mathbf{r}_{S^{(1)}} + \mathbf{r}_{S^{(2)}}\right], \tag{61}$$

which adds up to the four required constraints.

Therefore, let us denote by $\boldsymbol{\mu}$ the linear vector constraint. Then, from the extended Lagrangian, the equations of motion for the three atoms are given by

$$m_S \frac{d^2 \mathbf{r}_{S^{(1)}}}{dt^2} = \mathbf{F}_1 - 2\lambda(t)[\mathbf{r}_{S^{(1)}} - \mathbf{r}_{S^{(2)}}] - \frac{1}{2}\boldsymbol{\mu}, \tag{62}$$

$$m_S \frac{d^2 \mathbf{r}_{S^{(2)}}}{dt^2} = \mathbf{F}_2 + 2\lambda(t)[\mathbf{r}_{S^{(1)}} - \mathbf{r}_{S^{(2)}}] - \frac{1}{2}\boldsymbol{\mu}, \tag{63}$$

$$m_C \frac{d^2 \mathbf{r}_C}{dt^2} = \mathbf{F}_C + \boldsymbol{\mu}, \tag{64}$$

where $m_S$ and $m_C$ denote the mass of the sulphur and carbon atoms. If we now twice differentiate Eq. (61) with respect to time and use Eqs. (62)–(64), we obtain

$$\mathbf{F}_C + \boldsymbol{\mu} = \frac{m_C}{2m_S}(\mathbf{F}_1 + \mathbf{F}_2 - \boldsymbol{\mu}), \tag{65}$$

which helps us eliminate $\boldsymbol{\mu}$ and rewrite the equations of motion for the sulphur atoms as

$$m_S \frac{d^2 \mathbf{r}_{S^{(1)}}}{dt^2} =$$
$$\left[1 - \frac{m_C}{2(2m_S + m_C)}\right]\mathbf{F}_1 + \frac{m_C}{2(2m_S + m_C)}\mathbf{F}_2 + \frac{m_S}{2m_S + m_C}\mathbf{F}_C - 2\lambda(t)[\mathbf{r}_{S^{(1)}} - \mathbf{r}_{S^{(2)}}], \tag{66}$$

$$m_S \frac{d^2 \mathbf{r}_{S^{(2)}}}{dt^2} =$$
$$\left[1 - \frac{m_C}{2(2m_S + m_C)}\right]\mathbf{F}_2 + \frac{m_C}{2(2m_S + m_C)}\mathbf{F}_1 + \frac{m_S}{2m_S + m_C}\mathbf{F}_C + 2\lambda(t)[\mathbf{r}_{S^{(1)}} - \mathbf{r}_{S^{(2)}}]. \tag{67}$$

Equations (66) and (67) govern the motion of the sulphur atoms. The position of the carbon atom is of course fixed by the linear constraint that we have imposed on the system.

We still have one unknown, $\lambda(t)$, that must be determined. Since the bond constraint is quadratic, elimination of $\lambda(t)$ is not as easy as eliminating $\boldsymbol{\mu}$. Therefore, $\lambda(t)$ is determined at each time step by using the constraint equation, i.e., we solve the equations of motion iteratively until their solutions satisfy the bond constraint. Thus, if, for example, we utilize the Verlet algorithm, we can write

$$\mathbf{r}_{S^{(1)}}(t + \Delta t) = 2\mathbf{r}_{S^{(1)}}(t) - \mathbf{r}_{S^{(1)}}(t - \Delta t) + \Delta t^2 \left(1 - \frac{m_C}{2m_S + m_C}\right)\mathbf{F}_1(t)$$
$$+ \Delta t^2 \frac{m_S}{2m_S + m_C}\mathbf{F}_C(t) - 2\Delta t^2 \lambda(t)[\mathbf{r}_{S^{(1)}}(t) - \mathbf{r}_{S^{(2)}}(t)], \tag{68}$$

$$\mathbf{r}_{S^{(2)}}(t + \Delta t) = 2\mathbf{r}_{S^{(2)}}(t) - \mathbf{r}_{S^{(2)}}(t - \Delta t) + \Delta t^2 \left(1 - \frac{m_C}{2m_S + m_C}\right)\mathbf{F}_2(t)$$
$$+ \Delta t^2 \frac{m_S}{2m_S + m_C}\mathbf{F}_C(t) + 2\Delta t^2 \lambda(t)[\mathbf{r}_{S^{(1)}} - \mathbf{r}_{S^{(2)}}]. \tag{69}$$

It can be shown that the error in the numerical values of $\lambda(t)$ so obtained is of the order of $O(\Delta t^4)$.

Our discussion so far has been restricted to totally rigid molecules. We now consider partially rigid molecules that consist of rigid clusters that can move with respect to one another. For this purpose we describe the algorithm due to Rykaert *et al.* (1977), Ciccotti *et al.* (1982), and Rykaert (1985). Their algorithm, which is known as SHAKE (the author does not know why this name was given to this algorithm), is formulated based on the notion that the forces that particles experience are the physical and constraint forces. If $M$ is the number of the constraints, then the constraints are written as, $\mathcal{C}_k(R) = 0$, with $k = 1, 2, \cdots, M$, where $\mathcal{C}_k$ expresses the functional form of the constraint, e.g., restriction that the distance between two particles is always fixed. The constraint forces are given by,

$$\mathbf{F}_c = \sum_{k=1}^{M} \lambda_k \nabla_i \mathcal{C}_k, \tag{70}$$

where $\lambda_k$ is the Lagrange multiplier to be determined, and subscript $i$ signifies the fact that the gradient of constraint $\mathcal{C}_k$ must be taken with respect to $i$. Since we use the Verlet algorithm for integrating the equations of motion, we have the particles' positions at times $t - \Delta t$ and $t$ which satisfy the constraints imposed on the system. Thus, an intermediate position $\tilde{\mathbf{r}}_i$ is first calculated for particle $i$,

$$\tilde{\mathbf{r}}_i(t + \Delta t) = 2\mathbf{r}_i(t) - \mathbf{r}_i(t - \Delta t) + \Delta t^2 \, \mathbf{F}_i[\mathbf{r}_i(t)], \tag{71}$$

where $\mathbf{F}_i$ represents all the physical forces that the particle experiences. The true position is then computed from

$$\mathbf{r}_i(t + \Delta t) = \tilde{\mathbf{r}}_i(t + \Delta t) - \sum_{k=1}^{M} \lambda_k \nabla_i \mathcal{C}_k(\mathbf{r}^N), \tag{72}$$

where, as usual, $\mathbf{r}^N = \{\mathbf{r}_1, \mathbf{r}_2, \cdots, \mathbf{r}_N\}$ represents the position coordinates of all the particles. The Lagrange multipliers are determined by an iterative process. For iteration number $j$, a loop over the constraints is carried out, in each step of which the particles' positions and the Lagrange multipliers are updated. The positions are updated according to the following iterative scheme,

$$\mathbf{r}_i^{(n)} = \mathbf{r}_i^{(o)} - \Delta t^2 \lambda_k^{(j)} \nabla_i \mathcal{C}_k(\mathbf{r}^N), \tag{73}$$

where superscripts $n$ and $o$ denote the new and old values, respectively. To calculate the Lagrange multipliers $\lambda_k^{(j)}$ a first-order expansion of $\mathcal{C}_k(\mathbf{r}^N)$ is carried out which is then required to vanish. Therefore,

$$\mathcal{C}_k(\mathbf{r}^N)^{(n)} = \mathcal{C}_k(\mathbf{r}^N)^{(o)} - \Delta t^2 \lambda_k^{(j)} \sum_i \left\{ \nabla_i \mathcal{C}_k(\mathbf{r}^N)^{(o)} \cdot \nabla_i \mathcal{C}_k[\mathbf{r}^N(t)] \right\} + \cdots = 0, \tag{74}$$

where the sum is over all the particles. The final expression for $\lambda_k^{(j)}$ is given by

$$\lambda_k^{(j)} = \frac{\mathcal{C}_k(\mathbf{r}^N)^{(o)}}{\Delta t^2 \{\sum_i \nabla_i \mathcal{C}_k(\mathbf{r}^N)^{(o)} \cdot \nabla_i \mathcal{C}_k[\mathbf{r}^N(t)]\}}. \tag{75}$$

The iteration continues until the constraints are satisfied numerically to within a fixed acceptable error. This algorithm has turned out to be highly efficient and accurate.

## 9.2.8  Ion–Ion Interactions

As mentioned earlier in this chapter, Coulombic and gravitational interactions are both described by a pair-potential which is proportional to $\ln r$ in 2D and $r^{-1}$ in 3D. Since these are long-range interactions, it is not clear that one can use a finite cutoff distance in order to accurately calculate the corresponding forces. In case of a system with local electrical charges, even the overall neutrality of the system does not guarantee that the screening length is finite, since in most cases of practical interest the screening length is larger than the linear size of the system, and therefore use of a finite cutoff distance $r_c$ is not justified. Another problem arises when periodic boundary conditions are used in simulation of a system with a distribution of local charges. In this case the sums over the image charges in the periodic images of the system do not converge. This problem is solved by subtracting an offset from the potential (since adding or subtracting a constant to the potential does not change the resulting force, as the force is the gradient of the potential), leading to the following expression for the total configurational energy of the system for a set of $N$ particles with charge $q_i$,

$$E = \sum_{\mathbf{R}} \sum_{i<j}^{N} \sum_{j}^{N} \frac{q_i q_j}{|\mathbf{r}_i - \mathbf{r}_j + \mathbf{R}|} - \sum_{i<j} \sum_{j} q_i q_j \sum_{\mathbf{R}\neq 0} \frac{1}{\mathbf{R}}. \tag{76}$$

The sum over $1/\mathbf{R}$ is over the locations $\mathbf{R}$ of the periodic replicas of the system.

To discuss this further, we consider a concrete example—transport of charged particles in a disordered material, especially one with a quenched distribution of charge centers. This problem is relevant to many important phenomena, such as dynamic response of non-metallic materials, e.g., ionic glasses and polymeric and glassy conductors, highly defected crystals, and porous materials that are used for catalytic and separation processes. Although this problem had been studied extensively, until recently no consensus regarding the nature of the transport process had emerged. In particular, if $\langle R^2(t) \rangle$ is the mean square displacement of the mobile charged particles at time $t$, one expects to have

$$\langle R^2(t) \rangle \sim t^{\alpha}, \tag{77}$$

where $\alpha = 1$ for diffusive transport. However, the precise value of $\alpha$ was a controversial subject with some researchers claiming that $\alpha > 1$, while others believing that $\alpha \leq 1$, in which case the transport process is anomalous (see Chapter 6 of Volume I for a discussion of anomalous transport). This controversy also prevented interpretation of experimental data for diffusion of ions in random media with a distribution of charge centers. For example, during diffusion of ions through zeolites, which are porous catalytic materials with a distribution of charge center (ions and cations), it has been observed that, upon changing the charge on the diffusing particles (i.e., making the disorder stronger), the diffusivity decreases *by orders of*

*magnitude*. Despite its great importance, no efficient and reliable computer simulation of this problem was carried out for many years, because (1) the Coulombic interactions between the particles are long-ranged, and (2) the charge centers give rise to deep potential wells that may capture the mobile particles for long periods of time and slow down their motion.

Let us now describe how MD simulations of this problem can be carried and, in particular, how the effect of long-range Coulombic interactions can be taken into account by an efficient and reliable algorithm developed by Mehrabi and Sahimi (1999). They used both a continuum and a lattice representation of the system. The continuum representation was used when the fixed charge centers were distributed randomly in the medium. The lattice, which was simple-cubic, was utilized when a potential-potential correlation function, defined below, was utilized for generating the potentials due to the fixed charge centers. At time $t = 0$ the charged mobile particles are distributed randomly in the system, but as they move correlations develop between them. In addition to the Coulombic interaction, a short-range, LJ-type repulsive interaction (i.e., $\sim r^{-12}$; see Section 9.2.1) was also used to prevent capture of a mobile particle by an immobile one with a charge which has a sign opposite to that of the mobile particle.

The charge centers are either distributed explicitly throughout the medium, or are represented by their potential distribution, generated by the potential-potential correlation function. To make the system neutral, equal numbers of the centers with opposite charges are inserted in the system, and the same is done with the mobile particles. The Coulomb potential $U_i$ acting on the $i$th mobile particle is written as,

$$U_i = U_i^{(fm)} + U_i^{(mm)}, \tag{78}$$

where $U_i^{(fm)}$ is due to the interaction between the mobile particle and the fixed centers, while $U_i^{(mm)}$ is contributed by the interaction between the mobile particles themselves. $U_i^{(fm)}$ can be calculated by two different methods (yielding identical results). In one method, $U_i^{(fm)}$ [and also $U_i^{(mm)}$] is computed by the multipole expansion method described below. In the second method, one can use the fact that diffusion of charged particles in disordered media can be viewed as a transport process in an external potential field generated by the quenched disorder that represents the fixed charge centers. Thus, instead of directly distributing the charge centers with a given density $\rho(\mathbf{r})$, $U_i^{(fm)}$ is formally represented by the solution of the Poisson's equation which, for example in 3D, is given by

$$U_i^{(fm)}(\mathbf{r}) = -\frac{q_f q_m}{4\pi \varepsilon_0} \int \frac{\rho(\mathbf{r}')}{|\mathbf{r} - \mathbf{r}'|} \, d\mathbf{r}' \tag{79}$$

where $q_f$ and $q_m$ are the charges for the fixed and mobile particles, respectively, and $\varepsilon_0$ is the permittivity. The charge density $\rho(\mathbf{r})$ is represented by its correlation function $\chi_{\rho\rho}(\mathbf{r})$ which, in the case of Debye-Hückel statistics, is given by

$$\chi_{\rho\rho}(\mathbf{r}) = \rho_0 \delta(\mathbf{r}) - \frac{\rho_0 \kappa^2 e^{-\kappa |\mathbf{r}|}}{4\pi |\mathbf{r}|}, \tag{80}$$

where $\rho_0 = \langle \rho(\mathbf{r}) \rangle$, and $\kappa^{-1}$ is the spatial correlation or the screening length. The depth of the potential wells in which the mobile particles are captured by the immobile ones, and thus the radius of influence of the traps, is controlled by $\kappa^{-1}$. The larger $\kappa^{-1}$, the deeper is the potential well, and thus the larger the time spent in such traps. In the limit $\kappa^{-1} \to \infty$, the trapping times become infinitely large, and therefore the effective diffusivity is *zero*. The power spectrum $\hat{\chi}_{\phi\phi}(\omega)$ for the potential is calculated from that of the charge density $\hat{\chi}_{\rho\rho}(\omega)$, since Eq. (79) is a convolution integral of the charge density and the Green function for the potential generated by a single charge particle, and therefore in 3D

$$\hat{\chi}_{\phi\phi}(\omega) = \left( \frac{q_f q_m}{\varepsilon_0} \right)^2 \frac{\rho_0}{\omega^2(\omega^2 + \kappa^2)}. \tag{81}$$

Note that the 1D version of Eq. (81) is given by, $\hat{\chi}_{\phi\phi} = \rho_0(q_f q_m/\varepsilon_0)^2/(\omega^2 + \kappa^2)$. Hence, a realization of the potential field is generated as follows. Random numbers, distributed uniformly in $[-\sqrt{3}, \sqrt{3})$ (this ensures that their power spectrum is 1, as it should be), are assigned to the sites of the system. The resulting array is then Fourier transformed and multiplied by $\sqrt{\hat{\chi}_{\phi\phi}(\omega)}$, and then inverse Fourier transformed.

$U_i^{(fm)}$ and $U_i^{(mm)}$ can also be calculated by a multipole expansion method (Ding *et al.*, 1992; Mehrabi and Sahimi, 1999). In this method particle $i$ interacts with the nearby particles through the usual Coulomb potential, and with the far away particles through their *pre-calculated* multipole expansions of the potential. The total potential $U^{(g)}(\mathbf{r}) = \sum_j^N U_j(\mathbf{r})$ produced by a group of $N$ charges is

$$U^{(g)}(\mathbf{r}) = \frac{q}{r} - \mathbf{P} \cdot \nabla \left( \frac{1}{r} \right) + \frac{1}{2}\mathbf{Q} : \nabla\nabla \left( \frac{1}{r} \right) - \frac{1}{6}\mathbf{O} \vdots \nabla\nabla\nabla \left( \frac{1}{r} \right) + \cdots \tag{82}$$

where $q$, $\mathbf{P}$, $\mathbf{Q}$, and $\mathbf{O}$ are, respectively, the monopole, dipole, quadrupole, and octapole moments of the group of charges around the origin. In practice, we write

$$U^{(g)}(r) = \frac{q}{r} + \frac{1}{r^3} \sum_\alpha P_\alpha r_\alpha + \frac{1}{2r^5} \left( \sum_\alpha Q_{\alpha\alpha} r_\alpha r_\alpha + \sum_\alpha \sum_\beta Q_{\alpha\beta} r_\alpha r_\beta - r^2 Q \right)$$

$$+ \frac{1}{6r^7} \left[ 15 \left( \sum_\alpha O_{\alpha\alpha\alpha} r_\alpha r_\alpha r_\alpha + \sum_\alpha \sum_\beta O_{\alpha\alpha\beta} r_\alpha r_\alpha r_\beta + \sum_\alpha \sum_\beta \sum_\gamma O_{\alpha\beta\gamma} r_\alpha r_\beta r_\gamma \right) \right.$$

$$\left. - 9r^2 \left( \sum_\alpha O_\alpha \right) \right] + \cdots, \tag{83}$$

with

$$q = \sum_i q_i, \tag{84}$$

$$P_\alpha = \sum_i q_i R_{i\alpha}, \tag{85}$$

$$Q_{\alpha\alpha} = \sum_i q_i R_{i\alpha}^2, \tag{86}$$

$$Q_{\alpha\beta} = \sum_i q_i R_{i\alpha} R_{i\beta}, \tag{87}$$

$$Q = \sum_i q_i R_i^2, \tag{88}$$

$$O_\alpha = \sum_i q_i R_i^2 R_{i\alpha}, \tag{89}$$

$$O_{\alpha\alpha\alpha} = \sum_i q_i R_{i\alpha}^3, \tag{90}$$

$$O_{\alpha\alpha\beta} = \sum_i q_i R_{i\alpha}^2 R_{i\beta}, \tag{91}$$

$$O_{\alpha\beta\gamma} = \sum_i q_i R_{i\alpha} R_{i\beta} R_{i\gamma}, \tag{92}$$

where $r = |\mathbf{r}|$, $\mathbf{R}_i$ is the position vector of the $i$th charge, $\alpha$, $\beta$, and $\gamma$ stand for the coordinates $x$, $y$, and $z$, and $q_i$ is the charge of the $i$th (fixed or mobile) particle.

A highly efficient simulation technique is fundamental to this study. Hence, in addition to taking advantage of the multipole expansion, the 3D simulation cell is divided into 8 smaller equal boxes, called *children* of the original box [see, for example, Greengard and Rokhlin (1987) and references therein]. Each child box is a *parent* to 8 smaller boxes, with the division continuing up to a certain level which is called the maximum level (maxlevel) of division. The data needed for each particle, i.e., its position and type (mobile or fixed), are stored in a *particle object*. A *cell object* contains a list of its *current* particles. Each particle is also "connected" to the next and previous particle in the list. After setting up the entire data structure, the multipoles of each cell around its center at the maxlevel are calculated using the above expressions. Then, the multipoles of the parent cells are computed by translating and adding the multipoles of their children by a displacement vector $\boldsymbol{\ell} = (\ell_x, \ell_y, \ell_z)$. In terms of the old quantities, the new translated (primed) quantities are given by,

$$P'_\alpha = P_\alpha - q\ell_\alpha, \tag{93}$$

$$Q'_{\alpha\alpha} = Q_{\alpha\alpha} - 2\ell_\alpha P_\alpha + q\ell_\alpha^2, \tag{94}$$

$$Q'_{\alpha\beta} = Q_{\alpha\beta} - \ell_\beta P_\alpha - \ell_\alpha P_\beta + q\ell_\alpha\ell_\beta, \tag{95}$$

$$Q' = Q'_{\alpha\alpha} + Q'_{\beta\beta} + Q'_{\gamma\gamma}, \tag{96}$$

$$O'_{\alpha\alpha\alpha} = O_{\alpha\alpha\alpha} - 3\ell_\alpha + 2\ell_\alpha^2 - q\ell_\alpha^3, \tag{97}$$

$$O'_{\alpha\alpha\beta} = O_{\alpha\alpha\beta} - \ell_\beta O_{\alpha\alpha} + \ell_\alpha[-2Q_{\alpha\beta} + 2\ell_\beta P_\alpha + \ell_\alpha(P_\beta - q\ell_\beta)], \tag{98}$$

$$\begin{aligned} O'_{\alpha\beta\gamma} = {} & O_{\alpha\beta\gamma} + \ell_\gamma(\ell_\beta P_\alpha - Q_{\alpha\beta}) + \ell_\beta(\ell_\alpha P_\gamma - Q_{\alpha\gamma}) \\ & + \ell_\alpha(\ell_\gamma P_\beta - Q_{\beta\gamma}) - q\ell_\alpha\ell_\beta\ell_\gamma. \end{aligned} \tag{99}$$

Each particle's potential energy is divided into $U_i^{\text{near}}$ and $U_i^{\text{far}}$. A particle in a cell at the maxlevel interacts with all other particles in the same cell and in the neighboring cells by the usual Coulomb potential, thus yielding $U_j^{\text{near}}$. It also interacts with its parent's neighbors' children through the corresponding multipole expansions. Computations continue up to the entire simulation cell, hence yielding $U_i^{\text{far}}$. In this way, the number of the cells that interact with each particle is drastically reduced as one moves away from the particle. For example, in 3D with four levels of division the number of the interacting cells is only 415, rather than the original 4069 cells. This method is highly efficient for taking into account the effect of Coulomb and other long-range potentials.

However, even with such an efficient algorithm, MD simulation of this problem requires intensive computations. If the simulations are not carried out for long enough times, one may not be able to obtain the true asymptotic (long time) behavior of the system. As an example, consider the problem in a 1D system, such as a highly anisotropic material so that the motion of the mobile particles is restricted essentially to one direction. In this case, the mobile particles can only travel in the space *between* themselves, since they cannot "jump" over each other. There are many relatively fast diffusive jumps in the mean square displacements of the particles after certain periods of time. In between the jumps one has a slow motion that causes the *overall* transport to be anomalous [i.e., $\alpha < 1$ in Eq. (77)], not only in 1D but also in 2D and 3D. The mobile particles can be trapped in the potential wells (traps) that the quenched distribution of the charge centers creates. The traps have a finite sphere of influence, such that for any particle $i$ within the sphere the potential difference $\Delta U_i$ for a displacement that can take $i$ out of the sphere is very large, and thus the probability of an appreciable jump is small. However, after some time the particles are close to the boundary of the traps and escape with a displacement that takes them out of the traps. They then resume their diffusive motion until they are captured by another trap, and so on.

# 9.3    Nonequilibrium Molecular Dynamics Simulation

As mentioned in Section 9.2, equilibrium MD (EMD) simulations are applicable to systems that, at least in principle, are amenable to treatment by statistical mechanics. Although, as our discussions so far should have made it clear, many assumptions must be made and several approximations must be used, their main purpose is to make the computations tractable. However, if we are to compute the effective transport properties, such as the diffusivity, shear viscosity, and thermal conductivity, then EMD is not an effective tool. For example, as is well-known, by calculating the velocity correlation functions for every distinct pair of species in the system, one can obtain information about the microscopic motion of the molecules. However, since the velocity correlation function decays as the size of the system increases, use of EMD is not feasible for estimating the transport properties of a mixture of molecules in a system which is under the influence of an external potential gradient—a situation which is encountered in a very large

number of practical problems. One can use the velocity autocorrelation function, but this quantity can only be used for predicting the tracer or self-diffusivity $D_s$ of a species (i.e., when the system is very dilute) via the Green-Kubo equation:

$$D_s = \frac{1}{3N} \int_0^\infty \left\langle \sum_{i=1}^{N} [\mathbf{v}_i(t) \cdot \mathbf{v}_i(0)] \right\rangle dt, \qquad (100)$$

where $\mathbf{v}_i$ is the velocity of particle $i$. The self-diffusivity is, however, completely different from the *transport* diffusivity because tracer diffusion ignores the effect of the collective motion of other molecules, especially in systems with a moderate or high density. As such, EMD is not suitable for investigation of a transport process in a system on which an external potential (pressure, voltage, chemical potential, concentration, etc.) gradient has been imposed.

Non-equilibrium MD (NEMD) simulation represents a practical alternative to EMD for those systems for which the velocity correlation function is difficult, or meaningless, to measure. It is particularly ideal for the practical situation in which an external driving force is applied to the system. Several such modeling efforts have been reported. Notable among them are the Grand-Canonical Molecular Dynamics (GCMD) method [see, for example, Cagin and Petit (1991); Lupkowski and van Swol (1995)] in which Monte Carlo (MC) and MD simulations are combined (see below), and a dual control-volume GCMD technique (DCV-GCMD) [see, for example, Heffelfinger and van Swol (1994); MacElroy (1994); Cracknell *et al.* (1995); Xu *et al.* (1998)]. We describe here the DCV-GCMD method which has become an effective tool for studying systems that are under the influence of an external potential gradient.

To describe this method, we consider a concrete example, namely, transport of binary gas mixtures in a carbon nanopore, a problem of considerable importance in separation and purification processes (Xu *et al.*, 1998, 1999). The same method is applicable to transport of a one-component fluid in the pore which is under the influence of an external potential gradient. Other types of carbon pores, such as carbon nanotubes, are also the subject of intense research activity as their mechanical strength appears to be much larger than what is expected of such materials at this length scale. Transport in such nanotubes can also be studied by the method we describe here. For simplicity, we consider a slit pore, a schematic representation of which is shown in Figure 9.3 in which the origin of the coordinates is at the center. The external driving force is either a chemical potential (or, equivalently, pressure) or a concentration gradient applied in the $x$-direction. The system is divided into three regions. The $h$- and $\ell$-regions at the two ends represent two control volumes (CVs) exposed to the bulk fluid at high and low chemical potential, pressure, or concentration, respectively, while the middle region represents the pore in which transport occurs. The pore's length is $nL$ with $n$ being an integer. Periodic boundary conditions are employed in the $y$-direction. The two carbon walls are located at the top and bottom of the $xy$ planes. We consider transport of a binary gas mixture in such a slit pore.

FIGURE 9.3. Schematics of a slit pore, where $h$ and $l$ denote the high and low pressure (or chemical potential) control volumes (after Xu *et al.*, 1998).

The DCV-GCMD method combines integration of the equations of motion with GCMC insertions and deletions of the gas molecules in the two CVs. It is essential to maintain the densities of each gas component in the two CVs at some fixed values, which are in equilibrium with two bulk phases, each at a fixed gas pressure and concentration. The densities, or the corresponding chemical potentials of each component in the CVs, are maintained by carrying out a sufficient number of GCMC insertions and deletions of the particles. The probability of inserting a particle of component $i$ is given by

$$p_i^+ = \min \left\{ Z_i V_{ci} \exp \left( \frac{-\Delta E / k_B T}{N_i + 1} \right), \ 1 \right\}, \qquad (101)$$

where $Z_i = \exp(\mu_i / k_B T)/\Lambda^3$ is the absolute activity at temperature $T$, $\Lambda_i$ the de Broglie wavelength, $k_B$ the Boltzmann's constant, $\mu_i$ the chemical potential of component $i$, $\Delta E$ the potential energy change resulting from inserting or removing a particle, and $V_{ci}$ and $N_i$ the volume and number of atoms of component $i$ in each CV, respectively. The probability of deleting a particle is given by

$$p_i^- = \min \left\{ N_i \exp \left( \frac{-\Delta E / k_B T}{Z_i V_{ci}} \right), \ 1 \right\}. \qquad (102)$$

When a particle is inserted in a CV, it is assigned a thermal velocity selected from the Maxwell–Boltzmann distribution, Eq. (49), at the given $T$. An important parameter of NEMD simulations is the ratio $\mathcal{R}$ of the number of GCMC insertions and deletions in each CV to the number of MD integration steps between successive GCMC steps. This ratio must be selected appropriately in order to maintain the correct density and chemical potentials in the CVs, and also reasonable transport rates at the boundaries between the CVs and the transport region. Its typical value varies anywhere from 50:1 to 400:1. During the MD calculations particles crossing the outer boundaries of the CVs must be removed. In addition, one should allow for a non-zero streaming velocity (the ratio of the flux to the concentration of each component) in the entire transport region of each component, consistent with the presence of bulk pressure/chemical potential gradients along the flow direction. Otherwise, a zero streaming velocity in the transport region will lead to severely underestimated fluxes. Since the two CVs are assumed to be well mixed and in equilibrium with the two bulk phases that are in direct contact with

FIGURE 9.4. Density profiles for two gases components in the slit pore. The pore is defined for $-10 \leq X^* \leq 10$, where $X^* = H/\sigma$, with $H$ being the height of the pore and $\sigma$ a Lennard–Jones molecular diameter (after Xu *et al.*, 1998).



them, there should be no *overall* non-zero streaming velocity in these regions. However, to reduce the numerical instability caused by the discontinuity of the streaming velocities at the boundaries between the CVs and the transport region, a small streaming velocity can be added to the thermal velocity of all the newly inserted molecules within each CV that are located within a distance $0.5\sigma_1$ from the boundaries, where $\sigma_1$ is the LJ size parameter (or the effective molecular size) of the lightest of the two gases. The streaming velocity of each component in the transport region can be obtained by linearly interpolating between its two values in the two CVs. After a few thousands of integration steps, this procedure generates a potential gradient across the pore, an example of which is shown in Figure 9.4. To study the transport of a mixture in a potential gradient, the temperature of the system must be held constant in order to eliminate any contribution of the temperature gradient to the transport; hence special care must be taken to achieve this.

All the quantities of interest are calculated from such simulations. For example, one can calculate the density profiles of component $i$ along the $x$- and $z$-directions, $\rho_i^z(x)$ and $\rho_i^x(z)$, respectively. To compute $\rho_i^x(x)$ the simulation cell is divided in the $x$-direction into grids of size $\sigma_1$, and for each MD integration step the density profiles $\rho_i^z(x)$ are obtained by averaging the number of particles of component $i$ over the distance $\sigma_1$, with a similar procedure for $\rho_i^x(z)$. These quantities are important to understanding adsorption and transport properties of the gases between the two pore walls. For each component $i$ one can also calculate its flux $J_i$ by measuring the net number of its particles crossing a given $yz$ plane of area $S_{yz}$:

$$J_i = \frac{N_i^{LR} - N_i^{RL}}{n_s \, \Delta t \, S_{yz}}, \tag{103}$$

where $N_i^{LR}$ and $N_i^{RL}$ are the number of the gas molecules of type $i$ moving from the left to the right and vice versa, respectively, $\Delta t$ is the MD time step, and $n_s$ is the number of the MD steps over which the average is taken. The system is considered to have reached a steady state when the fluxes calculated at various $yz$ planes are, to within an acceptable error, equal to the averaged values, after which

the fluxes are calculated at the center of the transport region. The permeability $k_i$ of species $i$ is then calculated from

$$k_i = \frac{J_i}{\Delta P_i / nL} = \frac{nLJ_i}{\Delta P_i}, \qquad (104)$$

where $\Delta P_i$ is the partial pressure for species $i$ along the pore. The transport diffusivity $D_e$ is then obtained from the Fick's law, $J_i = -D_e \partial \rho_i / \partial x$, where $\partial \rho_i / \partial x$ is the adsorption density gradient of component $i$ along the $x$-direction. Another important property for the problem that we are describing is the dynamic separation factor $S_{21}$ defined as

$$S_{21} = \frac{k_2}{k_1}. \qquad (105)$$

The NEMD method that we described here has proven to be a practical tool for simulating transport properties of fluid mixtures in not only a carbon nanopore, but also in nanoporous materials, such as a variety of membranes, with an interconnected network of nanopores (Xu *et al.*, 2000a). Its predictions for some properties of interest are in good quantitative agreement with experimental data (Xu *et al.*, 2000b). Other NEMD methods have been discussed by Rapaport (1995).

## 9.4    Quantum Molecular Dynamics Simulation: The Car–Parrinello Method

In our discussion of MD simulation presented so far, we have treated the molecular system as a collection of classical particles for which the interaction potential is known, and have neglected all the quantum-mechanical effects. We now discuss quantum MD simulations in which this restriction is removed and quantum mechanical effects for the electronic degrees of freedom are taken into account. The landmark paper of Car and Parrinello (1985) in which a MD technique was used for minimizing the total energy functional and allowed one to use, with a very high degree of efficiency, local and non-local pseudo-potentials, opened the way for QMD computations. We now describe and discuss their method.

The key idea of Car and Parrinello was to treat the electronic wave functions $\Psi_i$ as dynamical variables, and to define a Lagrangian for the electronic system which is given by

$$\mathcal{L} = \sum_i m_\Psi \langle \dot{\Psi}_i | \dot{\Psi}_i \rangle - E[\{\Psi_i\}, \{\mathbf{R}_I\}, \{\ell_n\}]. \qquad (106)$$

Here $m_\psi$ is a *fictitious* mass associated with the electronic wave functions, giving rise to the kinetic energy term of the Lagrangian that arises due to the fictitious dynamics of the electronic degrees of freedom, $E$ is the Kohn–Sham energy functional, described in Section 9.1, which plays the role of the potential energy, $\mathbf{R}_I$ is the position of the ion $I$, $\ell_n$ defines the size and shape of the periodic unit cell, and $\dot{\Psi}_i = d\Psi_i / dt$. As Eq. (106) suggests, the details of the kinetic energy do not

matter. What is more important is that the mass $m_\Psi$ should be much smaller than the nuclear masses which would prevent transfer of energy from the classical to the quantum degrees of freedom over long periods of the numerical simulations. The electronic wave functions are subjected to the orthonormality constraints:

$$\int \Psi_i^*(\mathbf{r})\Psi_j(\mathbf{r})\, d^3r = \delta_{ij}, \tag{107}$$

where $\delta_{ij}$ is the Kronecker delta. To incorporate these constraints, Lagrange multipliers are introduced, so that the Lagrangian of the system is rewritten as

$$
\begin{aligned}
\mathcal{L} = \sum_i m_\Psi \langle \dot{\Psi}_i | \dot{\Psi}_i \rangle - E[\{\Psi_i\}, \{\mathbf{R}_I\}, \{\ell_n\}] \\
+ \sum_i \sum_j \lambda_{ij} \left\{ \left[ \int \Psi_i^*(\mathbf{r})\Psi_j(\mathbf{r})\, d^3r \right] - \delta_{ij} \right\}.
\end{aligned}
\tag{108}
$$

Mathematically, the Lagrange multipliers $\lambda_{jj}$ ensure that the wave functions remain normalized, while $\lambda_{ij}$ (with $i \neq j$) guarantee that the wave functions remain orthogonal. Physically, the Lagrange multipliers can be thought of as providing additional forces acting on the wave functions, ensuring that, at any given time $t$, they remain orthonormal as they propagate along their MD paths. The Lagrange multipliers are symmetric, $\lambda_{ij} = \lambda_{ji}$, and represent $\frac{1}{2}N(N+1)$ independent values which are determined by the $\frac{1}{2}N(N+1)$ orthonormality conditions. The iterative algorithm SHAKE, described in Section 9.2.7, can be used for determining the Lagrange multipliers, and was in fact utilized by Car and Parrinello.

### 9.4.1   The Equations of Motion

Having defined the Lagrangian of the system, the equations of motion for the electronic states are derived from

$$\frac{d}{dt}\left( \frac{\partial \mathcal{L}}{\partial \dot{\Psi}_i^*} \right) = \frac{\partial \mathcal{L}}{\partial \Psi_i^*}, \tag{109}$$

which yield the following equations of motion

$$m_\Psi \frac{d^2 \Psi_i}{dt^2} = -\mathcal{H}\Psi_i + \sum_j \lambda_{ij}\Psi_j, \tag{110}$$

where $\mathcal{H}$ is the Kohn–Sham Hamiltonian defined in Section 9.1. The similarity between Eq. (110) and Eq. (32), the equations of motion for a classical many-particle system solved by the MD method, is now apparent. To ensure that at any given time during the MD integration the wave functions remain orthonormal, the Lagrange multipliers must vary continuously with the time, and therefore they must be estimated continuously at *infinitely small* time separations. However, doing so would make the computations intractable. To make the computations tractable, it is usually assumed that the Lagrange multipliers remain constant in each time step $\Delta t$ during which the equations of motion are integrated. Although this approximation

makes the computations tractable, it also creates a new problem: At the end of each time step the wave functions will not be *exactly* orthonormal, so that a separate orthonormalization must also be carried out at the end of each time step (see below).

However, since a separate orthonormalization step must be carried out at the end of each time step, one may remove the orthogonality constraints from the equations of motion and use *partially* constrained equations of motion. After these equations have been integrated, the orthogonality constraints are imposed again, and the Lagrange multipliers $\lambda_{ij}$ for the normalization constraints are approximated by the expectation values of the energies of the states given by

$$\lambda_i = \langle \Psi_i | \mathcal{H} | \Psi_i \rangle. \tag{111}$$

With this approximation the equations of motion for the wave functions are given by

$$m_\Psi \frac{d^2 \Psi_i}{dt^2} = -(\mathcal{H} - \lambda_i)\Psi_i, \tag{112}$$

which ensures that the acceleration of an electronic state is always orthogonal to that state, and that the acceleration becomes zero when the wave function is an exact eigenstate. However, more generally, one can start from the estimate

$$\lambda_{ij} = 2\langle \Psi_i | \mathcal{H} | \Psi_j \rangle - m_\Psi \langle \dot{\Psi}_i | \dot{\Psi}_j \rangle, \tag{113}$$

and proceed with integration of Eq. (110). The difference between (111) and (113), aside from the approximate nature of (113), is that the latter depends on $m_\Psi$ which itself must be somehow selected carefully.

## 9.4.2    The Verlet Algorithm

Equations (111) can now be integrated by a Verlet algorithm. In analogy with Eq. (42), one writes

$$\Psi_i(t + \Delta t) = 2\Psi_i(t) - \Psi_i(t - \Delta t) + \Delta t^2 \frac{d^2 \Psi_i}{dt^2}. \tag{114}$$

If we utilize Eq. (112), the Verlet algorithm becomes

$$\Psi_i(t + \Delta t) = 2\Psi_i(t) - \Psi_i(t - \Delta t) - \frac{\Delta t^2}{m_\Psi}(\mathcal{H} - \lambda_i)\Psi_i(t). \tag{115}$$

A similar Velert algorithm can be written down for Eq. (110) if the initial estimates (113) are to be employed. Note that, if an expansion such as (29) is used as the solution for $\Psi_i$ (which is always the case when a pseudo-potential approximation is made), then the coefficients $c_{i,\mathbf{k}+\mathbf{G}}$ must be time-dependent. If we substitute such an expansion into Eq. (110) or (112), we obtain the governing equations for $c_{i,\mathbf{k}+\mathbf{G}}$ which can be integrated by the Verlet algorithm (see also Section 9.4.5).

Generally speaking, the performance of the Verlet algorithm for this type of computations is evaluated in terms of the rate at which it converges to the minimum-energy state. One must carry out the integration for a certain amount of time for the

problem to converge; the amount of computational time is controlled, of course, by the size of the time step $\Delta t$ used in the integration. It can be shown that for the QMD simulations that we are describing here, the largest possible value of $\Delta t$ is approximately given by

$$\Delta t \simeq 2 \sqrt{\frac{m_\Psi}{\epsilon_M - \epsilon_m}}, \tag{116}$$

where $\epsilon_M$ and $\epsilon_m$ are, respectively, the largest and smallest eigenvalues of the problem. Use of any $\Delta t$ which is significantly larger than the estimate (116) will lead to instability and large errors in the numerical solution of the equations of motion. Recall that in pseudo-potential approximation the electronic wave functions are represented by a plane-wave basis set, in which case the largest eigenvalue is determined by the cutoff kinetic energy that is used for truncating the basis set. This implies that if the cutoff energy is increased, the time step used in the Verlet algorithm must decrease. Moreover, the maximum allowed value of $\Delta t$ decreases as the size of the system increases, limiting the maximum system size that can be simulated by this algorithm.

We also note that, in our discussion of the integration algorithm, we have tacitly assumed that the Kohn–Sham Hamiltonian $\mathcal{H}$ is constant during the time evolution of the system. However, in addition to the instabilities that are caused by utilizing a too large value of the time step $\Delta t$, another type of instability arises if $\mathcal{H}$ is not allowed to vary when it must. Since the wave functions evolve under the MD scheme, the contribution of the exchange-correlation potential to the Kohn–Sham Hamiltonian [see Eqs. (12) and (15)] also varies with the time, leading to a new Hamiltonian. Thus, at the end of each time step, the Kohn–Sham Hamiltonian is updated with the new electronic density, so that the final time step leads to a self-consistent solution of the Kohn–Sham Hamiltonian and the determination of the minimum in the total energy.

### 9.4.3  The Kohn–Sham Eigenstates and Orthogonalization of the Wave Functions

As discussed above, the wave functions that are obtained from integrating the partially-constrained equations of motion are not orthogonal, and therefore an orthogonalization procedure must be used. The correct application of the constraints of orthogonality and normalization is actually essential for the success of the Car–Parrinello method. Car and Parrinello (1985) used an iterative method to orthogonalize the wave functions by utilizing the following equation

$$\Psi_i^{(n)} = \Psi_i^{(o)} - \frac{1}{2} \sum_{j \neq i} \langle \Psi_j^{(o)} | \Psi_i^{(o)} \rangle \Psi_j^{(o)}, \tag{117}$$

where the superscripts $n$ and $o$ refer to the new and old estimates of the wave functions, respectively. If Eq. (117) is applied repeatedly to the old estimates, then the electronic wave functions can be made orthogonal to any desired accuracy. For

example, if algorithm (117) is applied to two wave functions, they will be exactly orthogonal after only one application of the algorithm. Moreover, if the estimates of the wave functions are orthonormal up to order $\Delta t^4$ (the accuracy of the Verlet algorithm), then over a given time step algorithm (117) changes them to an extent within the same order. In general, the number of times that algorithm (117) must be applied depends on the number of the wave functions to be computed and their initial departure from orthogonality. However, algorithm (117) does not preserve normalization of the wave functions, and therefore they must be normalized after each application of the algorithm.

The Kohn–Sham energy functional is minimized by *any* set of wave functions that are a linear combination of its lowest-energy eigenstates. Under the MD scheme, the wave functions that are obtained after orthogonalization will be stationary, implying that in the MD method each wave function will converge to a linear combination of the lowest-energy Kohn–Sham eigenstates. This can create severe problems for treating metals, since the ability to converge to Kohn–Sham eigenstates (and not to their linear combinations) is highly important for metallic materials. Several methods have been proposed for addressing this problem (see, for example, Pederson and Jackson, 1991, and references therein). The actual Kohn–Sham eigenvalues can be found by either diagonalization of the matrix $\mathbf{A} = (A_{ij})$, the entries of which are given by

$$A_{ij} = \langle \Psi_i | \mathcal{H} | \Psi_j \rangle, \qquad (118)$$

or by the Gram–Schmidt algorithm (see, for example, Noble and Daniel, 1977) which is similar to (117) but without the factor 1/2 in front of the sum. Thus, instead of using algorithm (117) which generates wave functions that are linear combinations of the Kohn–Sham eigenstates, one can use the Gram–Schmidt method which results in orthogonal wave functions in such a way that all of the higher-energy wave functions are *forced* to be orthogonal to the lowest-energy wave functions. This in turn forces each state to converge to its lowest possible energy while satisfying the constraint that it be orthogonal to all states below it. Therefore, the set of lowest possible single-particle levels under these constraints comprises the Kohn–Sham eigenstates.

## 9.4.4   Dynamics of the Ions and the Unit Cell

In our discussion of QMD simulations so far we have assumed that the ionic positions and the shape of the unit cell remain fixed. However, in practice these represent additional degrees of freedom that have their own dynamics, and therefore must be considered. Fortunately, since we are interested in the final state of the system, which consists of ions and electrons in their minimum energy configurations, the exact path for reaching this state is not very important. This fact provides us with considerable flexibility.

The QMD technique of Car and Parrinello allows us to include in the computations the positions of the ions and the coordinates that define the size and shape of the unit cell. These constitute dynamical variables and including them in the

computations requires writing down a new Lagrangian, commonly referred to as the Car–Parrinello Lagrangian, which is given by

$$\mathcal{L} = \sum_i m_\Psi \langle \dot{\Psi}_i | \dot{\Psi}_i \rangle + \sum_I \frac{1}{2} m_I \left( \frac{dR_I}{dt} \right)^2 + \sum_n \frac{1}{2} m_c \left( \frac{d\ell_n}{dt} \right)^2 - E[\{\Psi_i\}, \{\mathbf{R}_I\}, \{\ell_n\}],$$
(119)

which should be compared with Eq. (106) that was for a system in which the positions of the ions and the shape and size of the unit cell were fixed. Here, $m_c$, similar to $m_\Psi$, is a fictitious mass associated with the dynamics of the coordinates that define the unit cell, namely, $\{\ell_n\}$, and $m_I$ is the mass of the ions. The rest of the notations is the same as those for Eq. (106). Given Lagrangian (119), we can derive two new sets of equations of motion, one for the ions given by

$$m_I \frac{d^2 \mathbf{R}_I}{dt^2} = -\frac{\partial E}{\partial \mathbf{R}_I},$$
(120)

and a second one for the coordinates of the unit cell,

$$m_c \frac{d^2 \ell_n}{dt^2} = -\frac{\partial E}{\partial \ell_n}.$$
(121)

These equations relate the rate of acceleration of the length of the lattice vectors to the diagonal components of the stress tensor, and also relate the accelerations of the angles between the lattice vectors to the off-diagonal components of the stress tensor, both integrated over the unit cell. Although these equations can be integrated at the same time that the equations of motion for the electronic states are integrated, the matter is not as straightforward as it may seem. The reason for this is that *physical* ground-state forces acting on the ions, and integrated stresses acting on the unit cell, cannot be calculated for *arbitrary* electronic configurations. The following discussion establishes this fact.

### 9.4.4.1   The Hellmann–Feynman Theorem

To understand why the ground-state forces acting on the ions and the integrated stresses acting on the unit cell cannot be calculated for arbitrary electronic configurations, recall that the force $\mathbf{F}_I$ acting on an ion is given by

$$\mathbf{F}_I = -\frac{dE}{d\mathbf{R}_I}.$$
(122)

As ions change their positions, the wave functions must change to self-consistent Kohn–Sham eigenstates corresponding to the new positions of the ions. These changes contribute to $\mathbf{F}_I$ since

$$\mathbf{F}_I = -\frac{\partial E}{\partial \mathbf{R}_I} - \sum_i \frac{\partial E}{\partial \Psi_i} \frac{d\Psi_i}{d\mathbf{R}_I} - \sum_i \frac{\partial E}{\partial \Psi_i^*} \frac{d\Psi_i^*}{d\mathbf{R}_I},$$
(123)

which follows from Eqs. (119) and (122), and should be compared with Eq. (120) which states that, $\mathbf{F}_I = -\partial E / \partial \mathbf{R}_I$. This apparent contradiction is due to the fact that in the Lagrange equations of motion for the ions, Eq. (120), the force acting

on the ions is not a physical force, rather it is a force that the ions experience from a particular electronic configuration. However, it is not difficult to show that, when the electronic wave functions are the eigenstates of the Hamiltonian, then the second and third terms of Eq. (123) vanish, and therefore, in this case, $\partial E / \partial \mathbf{R}_I$ yields the true *physical force* on the ions. This important result is known as the *Hellmann–Feynman theorem* (Hellmann, 1937; Feynman, 1939), and in fact holds for *any* derivative of the total energy.

The Hellmann–Feynman theorem greatly simplifies the task of computing the physical forces acting on the ions and the integrated stresses that are exerted on the unit cell, since it allows one to compute these quantities only when the wave functions are very close to exact eigenstates. Once the forces and stresses have been computed, the positions and the shape and size of the unit cell are also calculated using Eqs. (120) and (121), which represent their equations of motion. Each time the positions of the ions or the shape and size of the unit cell change, the electrons must be close to the ground state of the new ionic configuration in order to compute forces and stresses for the new configuration.

The simplest way that the Hellmann–Feynman forces are used is for determining the position of a local energy minimum of the ionic system. In this scheme the ions are moved along the directions of the Hellmann–Feynman forces until the residual forces (i.e., the deviations from $\partial E / \partial \mathbf{R}_I$) on all the atoms are smaller than a given value. These forces cause the ions to fluctuate around their equilibrium positions. The residual forces acting on the ions are never zero. If the system is to approach the minimum energy state, the magnitudes of the errors in the Hellmann–Feynman forces must be reduced, implying that the electronic configuration must be relaxed closer and closer to the instantaneous ground states as the ionic configuration approaches the local energy minimum.

### 9.4.4.2   Pulay Forces and Stresses

The reader may have noticed the absence of a fourth term in Eq. (123), $\partial \phi / \partial \mathbf{R}_I$, the derivative of the basis set $\{\phi\}$ (for representing the wave functions $\Psi_i$) with respect to the ionic positions. This term is called the Pulay force (Pulay, 1969). If a plane-wave basis set, Eq. (29), is used for representing the wave functions, then the derivative of each basis set with respect to $\mathbf{R}_I$ vanishes and the Pulay forces are exactly zero, in which case the Hellmann–Feynman forces are exactly equal to $\partial E / \partial \mathbf{R}_I$, provided that the electronic wave functions are Kohn–Sham eigenstates. If the Pulay forces do not vanish and are not computed either, then the calculated Hellmann–Feynman forces will be in error, with the error being independent of how close is the electronic configuration to its ground state, which means that a local energy minimum of the ionic system cannot be computed by calculating only the Hellmann–Feynman forces acting on the ions. This causes severe problems for the computations and may render them inefficient.

Although when a plane-wave basis set is used, the Pulay forces acting on the ions are exactly zero, the corresponding Pulay stresses acting on the unit cell may not be zero. Changing the size of the unit cell while keeping the number of

plane-wave basis states constant results in a change in the cutoff energy for the basis set. If we increase the number of plane-wave basis states by increasing the cutoff energy for the basis set, it will eventually result in a reduction in the total energy of the system. However, if the cutoff energy is large enough to achieve absolute convergence, the change in the total energy will be zero. In practice, most pseudo-potential computations are carried out with a cutoff energy at which energy differences have converged, but at which the total energies have not converged, in which case the diagonal components of the Pulay stresses acting on the unit cell will have non-zero values. It can be shown that the change in the total energy *per atom* is independent of the details of the ionic configuration, provided that the cutoff energy for the basis set is large enough for the energy differences to have converged. This facilitates computation of the Pulay stresses, since they can be calculated once and for all from the change in the total energy of a small unit cell as the cutoff energy is varied.

### 9.4.4.3   The Structure Factor and Total Ionic Potential

The total ionic potential in a solid material is computed by placing an ionic pseudo-potential at the position of every ion in the material. To do this, one calculates first the structure factor which, at wave vector $\mathbf{G}$ for ions of species $\alpha$, is given by

$$S_\alpha(\mathbf{G}) = \sum_I \exp(i\mathbf{G} \cdot \mathbf{R}_I), \tag{124}$$

where the sum is over the positions of all the ions of species $\alpha$ in a single unit cell. The total ionic potential $U_I$ is computed from

$$U_I(\mathbf{G}) = \sum_\alpha S_\alpha(\mathbf{G})u_\alpha(\mathbf{G}), \tag{125}$$

where $u_\alpha$ is the local pseudo-potential [see Eq. (30)], and the sum is over all the species in the unit cell.

At large distances $r$ the pseudo-potential is purely Coulombic and is of the form $Z/r$, where $Z$ is the atom's valence. Thus, the Fourier transform of the pseudo-potential at large distances is $\sim Z/G^2$, which diverges at $\mathbf{G} = \mathbf{0}$ (i.e., at $r = \infty$), and therefore the ion-electron energy is infinite. However, this poses no particular difficulty as there are similar divergences in the Coulomb energies due to the electron-electron and ion–ion interactions, which cancel the $\mathbf{G} = \mathbf{0}$ divergence. On the other hand, the pseudo-potential is not, in general, purely Coulombic, but contains a constant contribution at small $\mathbf{G}$ which for species $\alpha$ is given by

$$u_{\alpha,\text{core}} = \int 4\pi r^2 \left[ Z/r - u_\alpha^0(r) \right] \, dr, \tag{126}$$

where $u_\alpha^0$ is the local pseudopotential for the $l = 0$ angular momentum state [see Eq. (30)]. Integral (126) is zero outside the core region since, by definition, the potentials are identical outside this region. This non-Coulomb part of the pseudo-potential does contribute to the total energy at $\mathbf{G} = \mathbf{0}$, which is given by $N\Omega_c^{-1} \sum_\alpha N_\alpha u_{\alpha,\text{core}}$, where $N$ is the total number of electrons in the system, $N_\alpha$

is the total number of ions of species $\alpha$, and $\Omega_c$ is the volume of the unit cell. We remind the reader that the Coulomb energy of the ionic system can be calculated by the method described in Section 9.2.8.

## 9.4.5  Computational Procedure for Quantum Molecular Dynamics

Having discussed various aspects of the QMD method of Car and Parrinello, we can now summarize the method.

(1) One starts with an initial set of trial wave functions from which the Hartree potential and the exchange-correlation energy are computed.

   The choice of a reasonable initial set of trial wave functions is crucial to the efficiency and success of the computations, because the trial functions must be such that the electronic configurations converge to the ground state. If one is not careful, the convergence may not be achieved for two reasons.

   (i) If the initial states do not span the ground state, then the final solutions relax to a self-consistent Kohn–Sham eigenstates, but not to the ground state. The most obvious initial guess for the wave functions is a set of plane waves with the lowest kinetic energy. However, one must be careful with using this set because such initial states may not span the ground states of the electronic configurations. For example, if this set is used for computation of the electronic structure of an eight-atom cubic cell of any of the tetrahedral semiconductors, the electronic configuration does not converge to the ground state. Germanium, silicon, and carbon each have four valence electrons. An eight-atom cell of these materials contains 32 electrons, and therefore 16 doubly occupied electronic states are required to accommodate the electrons.

   (ii) The convergence may not be achieved if the QMD conserves any symmetry that is shared by the Hamiltonian and the initial electronic configuration. This symmetry can be broken when the electronic wave functions are orthogonalized. For example, algorithm (111) for orthogonalization does not break this symmetry, whereas the Gram–Schmidt algorithm does. However, the initial electronic states may not converge to the ground state without breaking this symmetry, in which case the Gram–Schmidt algorithm must be used. If one utilizes random initial values for the coefficients $c_{i,\mathbf{k}+\mathbf{G}}$ of the plane-wave basis electronic states [see Eq. (29)], then the initial states span the ground state, and no symmetry is conserved. Since only those plane-wave basis states that have the lowest kinetic energies contribute significantly to $\Psi_i$, it is sensible to assign non-zero random values to only those coefficients that have this property.

(2) The Hamiltonian matrices for each of the points included in the computations are constructed.

(3) The equations of motion for the electronic states are integrated using the Verlet algorithm, and the resulting wave functions are orthogonalized, either by algorithm (117) or by the Gram–Schmidt method; the wave functions are also normalized.

There are several "tricks" that can speed up the computations. For example, if we assume a local pseudo-potential approximation and use Eq. (29) in Eq. (16), then it is straightforward to show that the governing equations for the coefficients $c_{i,\mathbf{k}+\mathbf{G}}$ are given by

$$\frac{d^2 c_{i,\mathbf{k}+\mathbf{G}}}{dt^2} = -\omega_{i,\mathbf{k}+\mathbf{G}}^2 c_{i,\mathbf{k}+\mathbf{G}} - B_{i,\mathbf{k}+\mathbf{G}}, \tag{127}$$

where $\omega_{i,\mathbf{k}+\mathbf{G}}$ and $B_{i,\mathbf{k}+\mathbf{G}}$ are two sets of coefficients that arise when Eq. (29) is inserted into Eq. (16). However, we recognize Eq. (127) as the *oscillator equation*, which means that if the Verlet algorithm is to be used for integrating this equation, the time step must be such that $\Delta t\, \omega_{i,\mathbf{k}+\mathbf{G}} < 1$ for all of the plane-wave basis states. This implies that $\Delta t$ is restricted by the plane-wave basis states that have the *largest* kinetic energies. However, as we discussed in Section 9.1.4, the largest kinetic-energy basis states contribute the least to the solution of $\Psi_i$, Eq. (29), and this creates an unsatisfactory situation. As discussed by Payne *et al.* (1992), a way around this problem is assuming that the coefficients $B_{i,\mathbf{k}+\mathbf{G}}$ are constant over a time step of duration $\Delta t$, in which case Eq. (127) can be easily integrated analytically over the time step, with the solution given by

$$c_{i,\mathbf{k}+\mathbf{G}}(t + \Delta t) = 2\cos[\omega_{i,\mathbf{k}+\mathbf{G}}(t + \Delta t)]c_{i,\mathbf{k}_{\mathbf{G}}}(t) - c_{i,\mathbf{k}+\mathbf{G}}(t - \Delta t)$$
$$- \frac{2}{\omega_{i,\mathbf{k}+\mathbf{G}}^2}\{1 - \cos[\omega_{i,\mathbf{k}+\mathbf{G}}(t + \Delta t)]\}, \tag{128}$$

which can be used in the computations. Payne *et al.* (1992) provided a detailed discussion of many other ways of speeding up integration of the equations of motion.

In most cases, to increase the efficiency of the computations, damping is applied to the equations of motion. For example, one can apply a damping of the form $-\eta\dot{\Psi}_i$ to the equations of motion for the wave functions, or use a damping term $-\eta\dot{c}_{i,\mathbf{k}+\mathbf{G}}$ in (the right-hand side of) the equations of motion for the coefficients, where $\eta$ is the damping factor. The damping helps the coefficients $c_{i,\mathbf{k}+\mathbf{G}}$ evolve to the values that minimize the Kohn–Sham total energy functional.

(4) The electronic density generated by the new set of the wave functions is then calculated [see Eq. (5)], and the corresponding new Hamiltonian is constructed.

(5) A new set of the wave functions is obtained by integrating the equations of motion and orthonormalization of the results. The iterations are repeated until the resultant wave functions are stationary.

(6) The Kohn–Sham energy functional is minimized; its value gives the total energy of the system. If the plane-wave basis set, Eq. (29), has been used (which is always the case if pseudo-potential approximation is used), convergence tests must be performed to ensure that the calculated total energy has converged both as a function of the number of terms included in the set and the value of the cutoff energy that has been used for truncating Eq. (29) after a finite number of terms.

(7) Integration of the equations of motion for the ionic positions and the coordinates of the unit cell is also done along the QMD computations for the electronic configurations. It is sensible in the QMD methods to treat the electronic and ionic systems independently when relaxing the ions to their equilibrium positions, and therefore it is also possible to use different time steps for the two systems. As the integration proceeds, the time step for the ionic system must be progressively reduced as the ionic configuration approaches the local energy minimum. Such a procedure allows the electronic configuration to relax closer to its instantaneous ground-state configuration as the ions approach their equilibrium positions, hence ensuring that the errors in the Hellmann–Feynman forces are always less than the actual forces acting on the ions.

The QMD technique allows one to search large regions of configuration space and locate the deeper energy minima in a very efficient manner. Since the QMD combines a MD method with the DFT, it makes it possible to study temperature-dependent effects by a method that is free of the common assumptions about the nature of the interatomic forces. Figure 9.5 presents the pair correlation function



FIGURE 9.5. Pair correlation function for amorphous (top) and liquid (bottom) Si. Solid and dashed curve represent, respectively, the theoretical predictions and the experimental data (after Car and Parrinello, 1985).

FIGURE 9.6. Deformation of Se$_5$ linear chain (from top to bottom) to the calculated ground-state, stable structure. The time between successive figures corresponds to 500 time steps of $3.4 \times 10^{-16}$ sec (adapted from Jones and Gunnarsson, 1989).



$C(r)$ for both amorphous and liquid Si, calculated by Car and Parrinello, and its comparison with the experimental data. Given that the only piece of information that was supplied to the simulator was the volume of the unit cell, the agreement between the predictions and the data is truly remarkable.

Figure 9.6, adapted from Jones and Gunnarsson (1989), shows the evolution of the structure of Se$_5$ molecule, starting with an almost linear geometry, and obtained by the Car–Parrinello method. The effect of the core region was represented by a pseudo-potential. The time between successive structures is almost 500 time steps of $3.4 \times 10^{-16}$ seconds each. The last structure shown in the figure is stable and agrees with the data. The important point to remember is that, theoretically, there are many possible structures that correspond to the local energy minima, and the QMD method of Car and Parrinello can find the true structure very efficiently.

Let us mention three different and successful applications of the QMD and pseudo-potential methods that we just described. These applications represent only a sample of a very large number of computations that have been reported so far. Marks *et al.* (1996) reported the results of *ab initio* simulations of tetrahedral amorphous carbon based on the Car–Parrinello method. The simulated structure was in good agreement with the experimental data. Pickard *et al.* (2000) reported the results of computations for a variety of lanthanide- and actinide-containing compounds. The simulated structures and the associated structural parameters were in excellent agreement with the experimental data. More significantly, they showed that the pseudo-potential formulation allows a steady march through the Periodic Table, in the sense of calculating reliably the structural parameters of compounds that contain $f$ electrons.

Yoon *et al.* (2001) employed *ab initio* pseudo-potential DFT with a linear combination of atomic orbits and an exchange-correlation functional in the generalized gradient approximation, described in Section 9.1.2, to study structural deformation and intertube conductance of crossed carbon nanotube junctions. They reported good agreement between the results of their simulations and the experimental data.

## 9.4.6   Linear System-Size Scaling

The QMD method becomes increasingly less efficient as the size of the system increases, which is, of course, the problem with all molecular modeling methods. If the size of the system becomes too large, then the choice of an appropriate time step becomes very crucial as the charge density starts to fluctuate strongly. The fluctuations are the result of instabilities in the Kohn–Sham energy Hamiltonian, and reducing or eliminating them may require time steps that are too small, hence making the computations intractable. Thus, for very large systems one must utilize a highly efficient method.

The standard methods for solution of the electronic structure problems scale as $O(N^3)$ for large systems, where $N$ is the number of atoms in a cluster or supercell. This scaling holds for both iterative as well as standard eigensolution methods since, as described above, one must keep the occupied wave functions orthonormal. Thus, a considerable amount of effort has been devoted to development of methods that can, in principle, provide $O(N)$ scaling for the computations. We mention two of such methods that seem to achieve linear scaling. In the variational method proposed by Li *et al.* (1993), one takes advantage of the locality of the density matrix in real space to achieve $O(N)$ scaling. There is only one approximation which is controlled by the real-space radius $R_c$ that is used to truncate off-diagonal elements of the density matrix. The method is, of course, exact when $R_c \rightarrow \infty$. The solution of the variational problem is obtained by an *unconstrained* minimization, and can be incorporated into the Car–Parrinello method.

In the method of Ordejón *et al.* (1995) (see also Mauri *et al.*, 1993) an energy functional is used that possesses a global minimum for which, (1) the electronic wave functions are orthonormal, and (2) the correct electronic ground-state energy is obtained. Linear scaling is then achieved by introducing a spatially-truncated Wannier-like representation for the electronic states.

## 9.4.7   Extensions of the Car–Parrinello Quantum Molecular Dynamics Method

The QMD simulation method that was described above has become an almost every-day tool of studying materials' properties. However, despite all of its success, the Car–Parrinello method had certain limitations. For example, the QMD of Car and Parrinello has the same shortcoming of the classical MD methods, namely, the atomic nuclei are propagated according to the laws of *classical* mechanics on a single potential energy surface. In addition, since nuclear forces are obtained from the standard DFT, only the electronic ground state can be accessed.

Over the past few years, several extensions of the Car–Parrinello QMD have been proposed that have made the method an even more powerful technique. For example, Marx *et al.* (1999) incorporated the nuclear quantum effects in the method through the use of path integrals. In addition, Doltsinis and Marx (2002) developed a QMD that allows efficient simulations of electronically non-adiabatic processes, by coupling the restricted Kohn–Sham excited state to the ground state using a

surface hopping scheme. Evaluation of the non-adiabatic coupling is achieved by an efficient method that exploits the available wave function derivatives.

### 9.4.8   Tight-Binding Methods

An efficient method for *direct* minimization of the Kohn–Sham total energy functional is tight-binding MD (TBMD), a promising but, to our knowledge, not heavily used, method for electronic structure computations. In this method, the system is described by the following Hamiltonian

$$\mathcal{H} = \sum_n \frac{P_n^2}{2M_n} + \sum_n^{\text{occupied}} \langle \Psi_n | \mathcal{H}_{\text{TB}} | \Psi_n \rangle + E_r, \qquad (129)$$

where the first term is similar to the corresponding term of Eq. (1), $\mathcal{H}_{\text{TB}}$ is the tight-binding Hamiltonian, and $E_r$ is a short-range repulsive energy. The repulsive energy is written as, $E_r = \sum_i f[\sum_j U(r_{ij})]$, where $U(r_{ij})$ is a pairwise repulsive interaction and $f(x)$ is a functional which must be specified for the computations to proceed. In this sense, the Car–Parrinello method is superior because it requires specification of no such terms.

In TBMD simulations, the covalent bonding of a material is incorporated into the computations from its underlying electronic structure, rather than through an $N$-body potential. As discussed above, the Car–Parrinello approach relies on the expansion of the electronic wave functions by plane waves. In a TBMD simulation, electronic calculations require a few atomic orbitals for each atom, hence allowing a larger number of atoms to be used and longer simulation times to be utilized. The second term of Eq. (129) represents the electronic-band structure energy $E_{\text{TB}}$, which is calculated from a parameterized TB Hamiltonian. The difficulty of TBMD simulations is that if one is to compute $E_{\text{TB}}$ from a direct diagonalization of the tight-binding Hamiltonian $\mathcal{H}_{\text{TB}}$, the computations scale as (size of the system)$^3$, and hence the simulations may be limited to small systems. However, there are approximate methods of calculating $E_{\text{TB}}$ without direct diagonalization of $\mathcal{H}_{\text{TB}}$, e.g., based on a variational approach (see, for example, Qiu *et al.*, 1994). This method seems to provide accurate predictions for some carbon-based materials.

Note that the parameters of the TB Hamiltonian must be estimated empirically. However, this approach keeps some of the fundamental physics of the problem through a quantum-mechanical description of the electronic degrees of freedom, while the minimal basis and the sparse Hamiltonian make it very efficient. Mercer (1996) extended this method to compounds, while Bernstein and Kaxiras (1997) utilized it for defects and interfaces in silicon.

## 9.5   Direct Minimization of Total Energy

The need for direct minimization of the total-energy functional arises because when one attempts to find the electronic states that minimize the Kohn–Sham functional by an indirect minimization technique, certain instabilities in the evolution of the

electronic states may be encountered. Such instabilities will not arise in a direct method. In this section, we describe two of such methods and their application to minimization of the Kohn–Sham total-energy functional.

## 9.5.1    The Steepest-Descent Method

One of the best-known methods for minimizing a function $f(\mathbf{x})$ of a 2D or 3D variable $\mathbf{x}$ is the *steepest-descent* method, which is the oldest and most straightforward minimization technique. The algorithm proceeds by making an initial guess $\mathbf{x}_1$ and then improving it by moving in the direction where the functional $f(\mathbf{x})$ appears to change most rapidly. The steepest-descent direction is aligned with the vector

$$\mathbf{v}_1 = -\nabla f(\mathbf{x} = \mathbf{x}_1). \tag{130}$$

To reduce the value of $f(\mathbf{x})$, one travels in the direction of $\mathbf{v}_1$ from $\mathbf{x}_1$ to $\mathbf{x} = \mathbf{x}_1 + b\mathbf{v}_1$ (where $b$ is a scalar parameter), where the function is minimum. Thus, one can, for example, sample $f(\mathbf{x})$ at a number of points along $\mathbf{x}_1 + b\mathbf{v}_1$ in order to determine that value of $b$ which minimizes $f(\mathbf{x})$. This procedure minimizes only the value of the function along a particular line, and thus finds a *local* minimum. To find the absolute minimum of $f(\mathbf{x})$, one carries out a series of such line minimizations by using $\mathbf{x}_1 + b\mathbf{v}_1$ as the starting point for the next iteration and obtaining $\mathbf{x}_2$. Thus, a series of such vectors $\mathbf{x}_i$ is obtained such that the value of $f(\mathbf{x}_i)$ decreases with increasing $i$, the iteration number.

However, although this method is guaranteed to converge to the true minimum, the rate of convergence can be prohibitively slow. For example, if $f(\mathbf{x})$ has narrow valleys, successive approximations bounce off opposite sides, slowly approaching the bottom. Moreover, after a minimization is performed along a given gradient direction, a subsequent minimization along the new gradient re-introduces errors that are proportional to the previous gradient. As such, the steepest-descent method is not an attractive technique for minimization of the total energy.

## 9.5.2    The Conjugate-Gradient Method

Consider now the following symmetric and positive-definite functional form,

$$f(\mathbf{x}) = \frac{1}{2}\mathbf{x} \cdot \nabla \cdot \mathbf{x}. \tag{131}$$

Suppose that we wish to minimize $f(\mathbf{x})$ along a direction $\mathbf{d}_1$, starting from $\mathbf{x}_1$. The minimum will be at a point $\mathbf{x}_2 = \mathbf{x}_1 + b_1\mathbf{d}_1$, where $b_1$ satisfies the following equation

$$(\mathbf{x}_1 + b_1\mathbf{d}_1) \cdot \nabla \cdot \mathbf{d}_1 = 0. \tag{132}$$

A subsequent minimization along a direction $\mathbf{d}_2$ yields $\mathbf{x}_3 = \mathbf{x}_2 + b_2\mathbf{d}_2$, with $b_2$ satisfying

$$(\mathbf{x}_1 + b_1\mathbf{d}_1 + b_2\mathbf{d}_2) \cdot \nabla \cdot \mathbf{d}_2 = 0. \tag{133}$$

The best choice of $b_1$ and $b_2$ for minimizing $f(\mathbf{x})$ along $\mathbf{d}_1$ and $\mathbf{d}_2$ is obtained by differentiating Eq. (131) with respect to $b_1$ and $b_2$ and evaluating the result at $\mathbf{x}_3$. This procedure yields two equations,

$$(\mathbf{x}_1 + b_1\mathbf{d}_1 + b_2\mathbf{d}_2) \cdot \nabla \cdot \mathbf{d}_1 = 0, \tag{134}$$

$$(\mathbf{x}_1 + b_1\mathbf{d}_1 + b_2\mathbf{d}_2) \cdot \nabla \cdot \mathbf{d}_2 = 0. \tag{135}$$

However, in order for Eqs. (132) and (133) to be consistent with (134) and (135), one must have

$$\mathbf{d}_1 \cdot \nabla \cdot \mathbf{d}_2 = \mathbf{d}_2 \cdot \nabla \cdot \mathbf{d}_1 = 0, \tag{136}$$

implying that the directions $\mathbf{d}_1$ and $\mathbf{d}_2$ must be conjugate to each other, hence the name conjugate-gradient (CG) method. More generally, the directions $\mathbf{d}_i$ and $\mathbf{d}_j$ must be such that

$$\mathbf{d}_i \cdot \nabla \cdot \mathbf{d}_j = 0, \quad i \neq j. \tag{137}$$

Hence, in the CG method one takes the initial direction to be $-\nabla f(\mathbf{x}_1)$, and the subsequent directions are constructed from a linear combination of the new gradient and the previous direction that minimized $f(\mathbf{x})$. In practice, the new direction $\mathbf{d}_i$ in the $i$th iteration is obtained from

$$\mathbf{d}_i = \mathbf{v}_i + \omega_i\mathbf{d}_{i-1}, \tag{138}$$

where

$$\omega_i = \frac{\mathbf{v}_i \cdot \mathbf{v}_i}{\mathbf{v}_{i-1} \cdot \mathbf{v}_{i-1}}, \tag{139}$$

and $\omega_1 = 0$. Note that $\mathbf{v}_i = -\nabla f(\mathbf{x}_i)$. It has also been observed that in some cases a better estimate of $\omega_i$ is given by

$$\omega_i = \frac{(\mathbf{v}_i - \mathbf{v}_{i-1}) \cdot \mathbf{v}_i}{\mathbf{v}_{i-1} \cdot \mathbf{v}_{i-1}}. \tag{140}$$

In the CG method a function is *guaranteed* to converge to its true minimum. The reason is that since minimization along the conjugate directions are independent, each iteration reduces the dimensionality of the vector space by 1. Thus, one reaches the point at which the dimensionality of the function space is zero, i.e., there are no new directions left along which one can minimize the function, and therefore the trial vector must be at the position of the true minimum. The number of the iterations needed to reach the true minimum is therefore *at most* equal to the dimensionality of the vector space, although in practice it usually takes far fewer iterations to converge to the true solution.

### 9.5.3 *Minimizing the Total Energy by the Conjugate-Gradient Method*

The CG method can be used for simultaneous updating of all the electronic wave functions. Although, from a computational view point, this is done very efficiently, it does have the drawback that it requires a large amount of memory for storing data

between the iterations in order to ensure the conjugacy of the search directions.
The ideal method is one that takes advantage of the efficiency of the CG method
without requiring much more computer memory. This can be achieved by updating
a single band at a time. The steepest-descent direction for a single band is given
by

$$\zeta_i^m = -(\mathcal{H} - \lambda_i^m)\Psi_i^m, \tag{141}$$

where $m$ denotes the iteration number and $i$ labels the band. Here [see Eq. (111)],

$$\lambda_i^m = \langle \Psi_i^m | \mathcal{H} | \Psi_i^m \rangle. \tag{142}$$

However, since the electronic wave functions must be orthogonal, one must en-
sure that the steepest-descent vector is orthogonal to all the other bands. The
components of the steepest-descent direction vector that ensures this is given by

$$\zeta_i'^m = \zeta_i^m - \sum_{j \neq i} \langle \Psi_j | \zeta_i^m \rangle \Psi_j. \tag{143}$$

Note that $\Psi_j$ does not vary during iteration for band $i$.

Successive steps along the CG directions reduce the magnitude of the error in
the wave function. It can be shown that the components $\zeta_i$ of the steepest-descent
vector are only a multiple of the components $\Delta\Psi_i$ of the error vector if all the
unoccupied eigenstates of the Kohn–Sham Hamiltonian are degenerate. However,
since the Kohn–Sham Hamiltonian has a broad spectrum of eigenvalues, extending
up to the cutoff energy for the plane-wave basis set (see Section 9.1.4), it can lead
to poor convergence in the CG calculations. Each step tends to remove components
of the error vector that correspond to eigenstates in a particular energy range. The
rate of convergence can be improved significantly by using a *pre-conditioning
technique*. In this method, one multiplies the steepest-descent vector $\boldsymbol{\zeta}$ by a pre-
conditioning matrix $\mathbf{K}$ to generate a pre-conditioned steepest-descent vector $\boldsymbol{\eta}$
that represents more accurately the error vector $\boldsymbol{\Delta\Psi}$. Although, there exists, in
principle, an exact pre-conditioning matrix that pre-conditions the steepest-descent
vector such that it would be parallel to the error vector, its construction is very
expensive, since it is a full matrix. Therefore, one must do the pre-conditioning in
an approximate way.

To carry out an approximate pre-conditioning, we note that the higher-energy
eigenstates of the Kohn–Sham Hamiltonian are dominated (in a pseudo-potential
calculation) by plane-wave basis states, the kinetic energies of which lie close to
the eigenvalues of the state. Therefore, to make those states, with eigenvalues that
are dominated by their kinetic energy, nearly degenerate, one must remove the
effect of the kinetic energy operator from the Hamiltonian, which can be achieved
by using a diagonal pre-conditioning matrix which is essentially the inverse of the
kinetic-energy operator. It has been found that *pre-conditioned* steepest-descent
vectors that accurately represent the errors in the wave functions can be obtained
by multiplying the steepest-descent vectors by a pre-conditioning matrix $\mathbf{K}$ with

entries $K_{G,G'}$ that are given by

$$K_{G,G'} = \frac{8x^3 + 12x^2 + 18x + 27}{16x^4 + 8x^3 + 12x^2 + 18x + 27} \delta_{G,G'}, \qquad (144)$$

with

$$x = \frac{(\hbar^2|\mathbf{k} + \mathbf{G}|^2)}{2mT_i^m}, \quad T_i^m = \langle \Psi_i^m|(-\hbar^2/2m)\nabla^2|\Psi_i^m\rangle,$$

where $T_i^m$ is the kinetic energy of the state $\Psi_i^m$. $K_{G,G'}$ has the property that it approaches unity as $x \to 0$, with its first, second and third derivatives all being zero, hence ensuring that the small wave-vector components of the steepest-descent vector remain unchanged. Moreover, for $x > 1$, $K_{G,G'}$ asymptotically approach $[2(x-1)]^{-1}$ with an asymptotic expansion correct to fourth order in $1/x$. Therefore, this factor causes all of the large wave-vector components to converge at nearly the same rate.

The pre-conditioned steepest-descent vector $\eta^m$ is now given by

$$\eta^m = \mathbf{K}\zeta'^m, \qquad (145)$$

which is *not* orthogonal to all the bands. The components of the pre-conditioned steepest-descent vector which is orthogonal to all the bands is then calculated as

$$\eta_i'^m = \eta_i^m - \langle \Psi_i^m|\eta_i^m\rangle - \sum_{j \neq i}\langle \Psi_j|\eta_i^m\rangle \Psi_j. \qquad (146)$$

The components of the pre-conditioned conjugate directions $d_i^m$ are now given by [see Eq. (138)]

$$d_i^m = \eta_i'^m + \omega_i^m d_i^{m-1}, \qquad (147)$$

with [see Eq. (139)]

$$\omega_i^m = \frac{\langle \eta_i'^m|\zeta_i'^m\rangle}{\langle \eta_i'^{m-1}|\zeta_i'^{m-1}\rangle}, \qquad (148)$$

with $\omega_i^1 = 0$. The resulting conjugate directions will be orthogonal to all the other bands, except the wave function of the present band. Thus, a further orthogonalization to the present band is done via

$$d_i''^m = d_i^m - \langle \Psi_i^m|d_i^m\rangle \Psi_i^m, \qquad (149)$$

$$d_i'^m = \frac{d_i''^m}{\langle d_i''^m|d_i''^m\rangle^{1/2}}. \qquad (150)$$

Having constructed the pre-conditioned conjugate direction, the search for the minimum energy begins. A vector with the following components

$$\Psi_i^m \cos\theta + d_i'^m \sin\theta \qquad (151)$$

is a normalized vector orthogonal to all the other bands $\Psi_j$ (with $j \neq i$ and $\theta$ real), and thus satisfies the constraints of orthonormality required for the electronic wave

functions. The value of $\theta$ that minimizes the Kohn–Sham energy functional is then required by the CG method. It has been found that the following approximation to the Kohn–Sham energy,

$$E(\theta) = \langle E \rangle + A \cos(2\theta) + B \sin(2\theta), \tag{152}$$

is sufficient for locating its minimum. There are three constants, $A$, $B$, and $\langle E \rangle$, and therefore one needs three data points to evaluate them. One piece of information is provided by $E(\theta = 0)$ which is already known. The second data point is supplied by the fact that since at $\theta = 0$

$$\frac{\partial E}{\partial \theta} = \langle d_i'^m | \mathcal{H} | \Psi_i^m \rangle + \langle \Psi_i^m | \mathcal{H} | d_i'^m \rangle = 2\mathrm{Re}\left( \langle d_i'^m | \mathcal{H} | \Psi_i^m \rangle \right), \tag{153}$$

and that, in order to determine the components $\eta_i^m$ of the steepest-descent vector, $\mathcal{H} | \Psi_i^m \rangle$ has been determined, computing $E'(0) = \partial E / \partial \theta$, evaluated at $\theta = 0$ is cheap. The third data point can be taken to be the value of the Kohn–Sham energy at a point other than $\theta = 0$. This point should be selected to be far enough from $\theta = 0$ to avoid round-off errors, yet not so far from $\theta = 0$ that the estimate of $E'(0)$ becomes inaccurate. It has been found that $\theta = \pi/300$ gives very reliable results. With these three data points, the three constants are found to be

$$A = \frac{2E(0) - 2E(\pi/300) + E'(0)}{2[1 - \cos(2\pi/300)]}, \tag{154}$$

$$B = \frac{1}{2} E'(0), \tag{155}$$

$$\langle E \rangle = \frac{2E(\pi/300) - E'(0) - 2E(0) \cos(2\pi/300)}{2[1 - \cos(2\pi/300)]}. \tag{156}$$

We can now determine straightforwardly the value of $\theta_{\min}$ that minimizes the Kohn–Sham energy, since the stationary point of Eq. (152),

$$\theta_s = \frac{1}{2} \tan^{-1}\left( \frac{B}{A} \right), \tag{157}$$

that lies in the range $(0, \pi/2)$ is the required $\theta_{\min}$ for minimizing the Kohn–Sham energy. To start the next iteration of the CG method, a new wave function, given by

$$\Psi_i^{m+1} = \Psi_i^m \cos(\theta_{\min}) + d_i'^m \sin(\theta_{\min}), \tag{158}$$

is utilized. Note that each new wave function generates a new charge density, and therefore the electronic potentials in the Kohn–Sham Hamiltonian must be updated before starting the next iteration.

In practice, only a few iterations are performed on any given band before moving to the next band, since there is little to gain from converging a single band exactly while there are still errors in the estimates of the remaining bands. Once all the bands have been updated, the CG iterations are started again on the lowest band. One can also perform the CG iterations on one band until the total energy changes

by less than a given fraction of the change of energy in the first CG iteration, and then start iterating on the next band.

## 9.6 Vectorized and Massively-Parallel Molecular Dynamics Simulation

Molecular dynamics simulations do not usually use large amounts of computer memory because one must record only the vectors that contain information about the atoms. Computationally, the MD simulations become very intensive and large scale when one wishes to use a large number of atoms and simulate "long" periods of time. Such large-scale simulations are necessitated by the fact that, since the effective size of atoms is typically several angstroms, to approach even sub-micron length scales, one must use millions of atoms in the simulations. Moreover, the typical time steps that are used in the MD simulations are of the order of femtoseconds, and therefore the equations of motion must be integrated over hundreds of thousands of time steps in order to simulate picosends of real time. Such length and time scale constraints conspire together to make the MD simulations very time consuming. They have also provided the impetus for developing efficient MD algorithms that are optimized for vector supercomputers (see, for example, Grest *et al.*, 1989; Schöen, 1989; Morales and Nuevo, 1992), and to special-purpose computers for carrying out MD computations (Auerbach *et al.*, 1987; Bakker *et al.*, 1990). Let us first briefly discuss the vectorized algorithms for the MD simulations.

### 9.6.1   *Vectorized Molecular Dynamics Algorithms*

As discussed in Section 9.2, when the forces are short-ranged, one usually uses a cutoff distance $r_c$ such that any two particles that are apart a distance larger than $r_c$ do not interact with each other. Thus, the key to efficient MD simulations with short-range forces is minimizing the number of neighboring atoms that must be checked for possible interactions. An efficient algorithm for doing this was first suggested by Verlet (1967), and is commonly referred to as the *neighbor lists* method. In this method, a list is constructed for every atom in the simulation cell that contains the nearby atoms that are at an extended distance $r_e = r_c + \delta$. Relative to $r_c$, the value of $\delta$ is small, but its optimal value depends on such parameters as the temperature and particle density. These lists are updated after every few integration steps, before any atom has moved from a distance $r > r_e$ to $r < r_e$.

Hockney *et al.* (1974) proposed another method, usually referred to as the *link-cell* method, in which at every time step all the atoms are binned into 3D cells of linear size $l$ with $l = r_c$ or slightly larger. Therefore, for every atom one must check only 27 bins—the bin the atom is in and the 26 bins surrounding it. A very efficient algorithm results when one combines the two methods such that the atoms are binned only every few time steps in order to construct (or update) the neighbor lists. The size of the cells is $l > r_e$. At intermediate time steps the neighbor lists

are used in the usual way in order to identify neighbors within a distance $r_c$ of each atom. The combined algorithm is made even more efficient by taking advantage of Newton's third law which allows one to compute a force for each pairs of atoms, instead of once for each atom in the pair. This reduces the necessary searches to only half the surrounding bins of each atom to form its neighbor list. In this way an atom $j$ is stored in atom $i$'s list, but not vice versa, hence halving the number of force computations that must be done.

Another factor that can increase the efficiency of a vectorized MD algorithm is careful data and loop structures, without which the computer program cannot be completely vectorized. For example, Grest *et al.* (1989) combined neighbor list/link-cell method to create long lists of pairs of neighboring atoms. At each time step, they updated the lists to keep only those particle pairs that were within the cutoff distance $r_c$. They also organized the lists into packets in order to prevent any atom from appearing twice in the lists. In this way, the force computations for each packet was completely vectorized, resulting in an algorithm that was about one order of magnitude faster than an algorithm without such data structures.

Another way of developing highly efficient algorithms for MD simulations is by taking advantage of the natural parallelism that exists in such simulations when short-range forces are involved. We now discuss such algorithms in detail.

## 9.6.2   Massively-Parallel Molecular Dynamics Algorithms

If the forces that act on the atoms are short-ranged, then MD computations can benefit greatly from parallel algorithms. Such algorithms began to emerge in the late 1980s and early 1990s, and include those of Raine *et al.* (1989), Bruge and Fornili (1990), Mel'cuk *et al.* (1991), Lin *et al.* (1992), Essenlink *et al.* (1993), Form *et al.* (1993), Rapaport (1993), Lomdahl *et al.* (1993), Beazley and Lomdahl (1994), Smith and Forester (1994), Plimpton (1995), Stadler *et al.* (1997), and many more, with the latest (at the time of writing this book) being that of Roth *et al.* (2000). Beazley *et al.* (1995) reviewed and compared many of these algorithms. Most of the work in this direction has been for single-instructor/multiple-data (SIMD) parallel machines, or for multiple-instruction/multiple-data (MIMD) parallel computers. The latter machines with several thousands of processors possess the computational power that is comparable with the fastest vector computers. Here, we describe three parallel algorithms for MD simulations that were developed by Plimpton (1995). His algorithms and their variants have been used in large-scale simulations of materials and have proven to be highly efficient. These algorithms are for systems in which the forces acting on the molecules are short-ranged. For such short-range MD simulations the computational effort per time step scales as $N$, the number of atoms used in the simulations. Another feature of these algorithms is that the atoms can undergo large displacements over the duration of the simulation. Finally, the performance of these algorithms is optimal with respect to both the numbers $N$ of the atoms and $N_p$ of the processors. As such, they are very flexible and efficient.

The reason that MD simulations with short-range forces are amenable to parallel computations is that calculations of the forces and updating of the positions can be done *simultaneously* for all the atoms. Thus, the main goal is to divide the force computations evenly among the processors so as to achieve maximum parallelism. There are at least three ways of achieving this, and what follows is a brief description of each, which are patterned closely after Plimpton (1995).

### 9.6.2.1   Atom-Decomposition Algorithms

In the atom-decomposition (AD) algorithm, at the beginning of the MD simulations each of the $N_p$ processors is assigned $N/N_p$ atoms. Atoms in a group do not have to have any special relation with each other. We assume, for the sake of simplicity, that $N$ is a multiple of $N_p$, although the algorithm is general. Each processor computes the forces on its $N/N_p$ atoms and updates their positions and velocities, *regardless* of where they move in the physical space. In general, one has an $N \times N$ force matrix $\mathbf{F}$, such that the entry $F_{ij}$ of $\mathbf{F}$ represents the force on atom $i$ due to atom $j$. Since the forces are short-ranged, $\mathbf{F}$ is sparse and, moreover, due to Newton's third law, $F_{ij} = -F_{ji}$. Suppose that $\mathbf{r}$ and $\mathbf{f}$ are vectors of length $N$ which store the position and total force on each atom, so that for a 3D simulations $\mathbf{r}_i$ would store the three coordinates of atom $i$. The AD algorithm assigns each processor a sub-block of $\mathbf{F}$ consisting of $N/N_p$ rows of the matrix. We number the processors from 0 to $N_p - 1$, so that processor $P_m$ (with $m = 0, 1, \cdots, N_p - 1$) computes matrix entries in the $\mathbf{F}_m$ sub-block of rows, and is also assigned the corresponding sub-vectors $\mathbf{r}_m$ and $\mathbf{f}_m$, each of length $N/N_p$.

A most important aspect of the computations is the communication between the processors, since to compute all the entries in $\mathbf{F}_m$, processor $P_m$ needs the positions of many atoms that belong to *other* processors. This gives rise to *all-to-all* communication, an operation that, at every time step, supplies the updated positions of the particles to all the processors. Various algorithms have been developed for performing this operation efficiently. We describe the algorithm due to Fox *et al.* (1988), utilized by Plimpton, who referred to the all-to-all communication procedure as an *expand* operation. To store the entire vector $\mathbf{r}$, every processor allocates memory of length $N$. At the beginning of the expand operation, processor $P_m$ has vector $\mathbf{r}_m$, the updated segment of $\mathbf{r}$ of length $N/N_p$. As the updated information is supplied to the processors, they must store them in the right place in their own copy of $\mathbf{r}$. For example, such steps for an eight-processor machine are illustrated in Figure 9.7, where the processors have been mapped consecutively onto the sub-pieces of the vector. In the first step of the communication, each processor exchanges sub-pieces of information. For example, processor 2 does this with processor 3. After this step, each processor has a piece of $\mathbf{r}$ with length $2N/N_p$ (factor 2 is caused by the new information that has just arrived). In the second step, each processor exchanges its new piece with another processor two positions away, after which each processor has a piece of $\mathbf{r}$ with length $4N/N_p$. The last step consists of having each processor exchanging a piece of $\mathbf{r}$ of length $\frac{1}{2}N$

FIGURE 9.7. *Expand* (left) and *fold* (right) operations among 8 processors, each requiring three steps. During the expand, processor 2 receives successively longer sub-vectors from processors 3, 0, and 6. In the fold, processor 2 receives successively shorter sub-vectors from processors 6, 0, and 3 (after Plimpton, 1995).

with a process or $\frac{1}{2} N_p$ positions away, after which each processor has the complete vector **r**.

The inverse of the expansion operation is also a useful procedure and is commonly called a *fold* operation. Suppose that each processor has stored new force values in its copy of **f**. Therefore, processor $P_m$ needs to know the $N/N_p$ values in $\mathbf{f}_m$, where each of the values is summed across all $N_p$ processors. This is also shown in Figure 9.7. Thus, each processor exchanges half the vector with a processor that is $\frac{1}{2} N_p$ positions away, such that it receives the half that it *is* a member of, and sends the half that it is *not* a member of, and then sums the received values with its corresponding retained sub-vector. This operation is repeated such that at each step the length of the exchanged data is halved. Note that the expand and fold operations are optimal, since each processor carries out $\log_2 N_p$ sends and receives and exchanges $N - N/N_p$ additions in the fold operation. The main disadvantage of this algorithm is that, it requires $O(N)$ storage on every processor.

The AD algorithm can be implemented in two different ways. In the first one, which we refer to as the AD1 algorithm, one assumes that, at the beginning of each time step, every processor has an updated copy of the vector **r**, and thus "knows" the positions of all the $N$ atoms. Then, the step-by-step implementation of the AD1 algorithm is as follows.

(1) This step consists of constructing the neighbor lists for all the pairwise interactions that must be computed in block $\mathbf{f}_m$. This is not done at every time step, but after every few steps. It is more efficient for a processor to inspect all the $N^2/N_p$ pairs in its $\mathbf{f}_m$, if the ratio of the physical size of the system to $r_c$ is small. For large simulations, however, in which four or more bins can be created in each dimension, it is more efficient for a processor to bin all the $N$ atoms and then inspect the 27 surrounding bins, the length of each is $N/N_p$. The overall scaling of this inspection process is about $N + N/N_p$.

(2) The neighbor lists are now utilized for calculating the non-zero matrix entries in $\mathbf{F}_m$. Since each pairwise interaction force is calculated and the force components are summed into $\mathbf{f}_m$, $\mathbf{F}_m$ is never actually stored as a matrix, hence saving a lot of computer memory. At the end of this step, each processor knows the total force $\mathbf{f}_m$ on each of the $N/N_p$ atoms that it owns.

(3) The information in (2) is now used to update the positions and velocities of the particles (i.e., by using them in the equations of motion and integrating them over one time step).

(4) The updated positions of the atoms obtained in the previous step are now shared among all the $N_p$ processors, using the expand operations defined above, in preparation for the next time step.

This algorithm does not take advantage of Newton's third law. Algorithm AD2 does use this law and thus reduces the time of the computations by decreasing the cost of communication between the processors. In order to do this, another matrix $\mathbf{G}$ is used such that $G_{ij} = F_{ij}$, except that $G_{ij} = 0$ if $i > j$ and $i + j$ is even, or when $i < j$ and $i + j$ is odd. Thus, for example, step 1 of AD2 is similar to that of AD1, except that only half as many neighbor list entries are constructed by each processor because $\mathbf{G}_m$ has only half the non-zero entries of $\mathbf{F}_m$. Similarly, if the neighbor lists are constructed by binning, then although all the $N$ atoms must be binned, each processor needs to inspect only half the surrounding bins of each of its $N/N_p$ atoms. In step 2 of the AD2 algorithm, the neighbor lists are utilized for computing the entries of $\mathbf{G}_m$. For the interaction between atoms $i$ and $j$, the forces acting on $i$ and $j$ are summed into both $i$ and $j$ locations of force vector $\mathbf{f}$, implying that each processor must store a copy of the *entire* force vector, as opposed to storing only $\mathbf{f}_m$ done in the AD1 algorithm. After all the matrix entries are computed, $\mathbf{f}$ is folded across all the $N_p$ processors, hence enabling each processor $P_m$ to have $\mathbf{f}_m$, the forces acting on its atoms. Steps 3 and 4 of the AD2 algorithm are the same as that of AD1.

## 9.6.2.2    Force-Decomposition Algorithms

The force-decomposition (FD) algorithm is based on a block-decomposition of the force matrix $\mathbf{F}$, rather than a row-wise decomposition used in the AD algorithm. We assume, for the sake of ease of exposition, that $N_p$ is a power of 2 and that $N$ is a multiple of $N_p$, although the algorithm is general. Sub-blocks of the force matrix $\mathbf{F}$ are assigned to the processors. Figure 9.8 shows the decomposition of the matrix among 16 processors. The block-decomposition shows there is actually a permuted force matrix $\mathbf{F}'$, formed by rearranging the columns of $\mathbf{F}$ in a particular way. If we arrange the $\mathbf{r}_\alpha$ pieces in row-order (where $\alpha = 1, 2, \cdots, \sqrt{N_p} - 1$), they form the usual position vector $\mathbf{r}$ which is shown as a vertical bar in Figure 9.8. One spans the columns with a permuted position vector $\mathbf{r}'_\beta$ (where $\beta = 1, 2, \cdots, \sqrt{N_p} - 1$), shown as a horizontal bar in the figure. Thus, in the example of Figure 9.8, $\mathbf{r}$ stores each processor's piece in the usual order $(0, 1, 2, \cdots, 15)$, whereas $\mathbf{r}'$ stores them as $(0, 4, 8, 12, 1, 5, 9, 13, 2, 6, 10, 14, 3, 7, 11, 15)$. In this case, $F'_{ij}$ is the force acting on atom $i$ in vector $\mathbf{r}$ due to atom $j$ in permuted vector $\mathbf{r}'$.

The size of sub-block $\mathbf{F}'_m$ is $(N/\sqrt{N_p}) \times (N/\sqrt{N_p})$. Thus, to compute the entries of $\mathbf{F}'_m$, processor $P_m$ must know $\mathbf{r}_\alpha$ and $\mathbf{r}'_\beta$ of the $\mathbf{r}$ and $\mathbf{r}'$ vectors, each with a length $N/\sqrt{N_p}$. These elements are computed and accumulated into corresponding force sub-vectors $\mathbf{f}_\alpha$ and $\mathbf{f}'_\beta$. Thus, for example, for processor 6 of Figure 9.8, $\mathbf{r}_\alpha$ consists of the $\mathbf{r}$ sub-vectors $(4 - 7)$ and $\mathbf{r}'_\beta$ is made of $\mathbf{r}'$ sub-vectors $(2, 6, 10, 14)$. There are actually two ways, FD1 and FD2, by which this algorithm can be implemented. Algorithm FD1 consists of the following steps.

FIGURE 9.8. The division of the permuted force matrix $F'$ among 16 processors in the force-decomposition algorithm. Processor 6 is assigned a sub-block $F'_6$ with size $N/\sqrt{P} \times N/\sqrt{P}$ which, to compute its matrix elements, must know the corresponding $N/\sqrt{P}$−length pieces $x_\alpha$ and $x'_\beta$ and permuted position vector $\mathbf{x}'$ (after Plimpton, 1995).

(1) Neighbor lists are constructed. If $N$ is small, construction of the lists is most efficiently done by checking all the $N^2/N_p$ possible pairs in $\mathbf{F}'_m$. For large $N$ the $N/\sqrt{N_p}$ atoms in $\mathbf{r}'_\beta$ are binned and the 27 surrounding bins of each atom in $\mathbf{r}_\alpha$ are checked.
(2) The neighbor lists are utilized for computing the entries of $\mathbf{F}'_m$. As they are computed, the entries are summed into a local copy of $\mathbf{f}_\alpha$, and therefore $\mathbf{F}'_m$ does not have to be stored in matrix form.
(3) By a fold operation within each of the rows, processor $P_m$ obtains the total force $\mathbf{f}_m$ acting on its $N/\sqrt{N_p}$ atoms.
(4) Processor $P_m$ utilizes $\mathbf{f}_m$ to update the positions of the $N/N_p$ atoms in $\mathbf{r}_m$.
(5) The processors in row $\alpha$ perform an expand operation on their $\mathbf{r}_m$ sub-vectors, so that each of them obtains the entire $\mathbf{r}_\alpha$. Similarly, the processors in each column $\beta$ perform an expand operation on their $\mathbf{r}_m$. It is in this step that using the permuted force matrix $\mathbf{F}'$ saves extra communications and hence computer time.

However, the FD1 algorithm does not take advantage of Newton's third law, and therefore calculates each pairwise interaction force twice. In algorithm FD2, this is avoided by using the modified force matrix $\mathbf{G}$, the construction of which was discussed for the AD2 algorithm. Then, $\mathbf{G}$ is permuted in the same way as $\mathbf{F}$ to form $\mathbf{G}'$. Step 1 of FD2 is the same as in FD1 (except that half as many interactions are stored in the neighbor lists). In step 2, for each $ij$ entry, the computed components of the force are now summed into two sub-vectors instead of one, such that the force acting on atom $i$ is summed into $\mathbf{f}_\alpha$ in the location corresponding to row $i$, while the force acting on atom $j$ is summed into $\mathbf{f}'_\beta$ in its proper location. Step 3 consists of three stages. First, the $\sqrt{N_p}$ processors in column $\beta$ fold their local

copies of $\mathbf{f}'_\beta$, resulting in $\mathbf{f}'_m$, each element of which is the sum of an entire column of $\mathbf{G}'$. Next, the same type of operations are performed for the row contributions, resulting in $\mathbf{f}_m$, each element of which is the sum of an entire row of $\mathbf{G}'$. Finally, the column and row contributions are subtracted element by element in order to obtain the total forces $\mathbf{f}_m$ acting on the atoms that belong to processor $P_m$, which can now update the positions and velocities of its atoms. Steps 4 and 5 of FD2 are identical to those of FD1.

### 9.6.2.3    Spatial-Decomposition Algorithms

In the spatial-decomposition (SD) algorithm, the simulation cell is divided into small 3D boxes. Each box is assigned to one processor which, at each time step, computes focres on and updates the positions and velocities of all atoms that it contains. The atoms that are assigned to the processors can change; they are assigned as they move through the simulation cell. To compute the forces acting on its atoms, a processor only needs to know the positions of atoms in the nearby boxes, and therefore, unlike the AD and FD algorithms, the communication in the SD algorithm is local. The size and shape of the boxes depend on $N$, $N_p$, and the aspect ratio of the simulation cell (assumed to be a 3D rectangular parallelepiped). The number of the processors in each dimension is selected so as to make the boxes as cubic as possible, since such configurations minimize the cost of communications, when $N$ is very large, which is proportional to the boxes' surface. The linear size of the boxes can be larger or smaller than the cutoff distances $r_c$ and $r_e$ defined above.

Each processor maintains two data sets, one for the $N/N_p$ atoms that have been assigned to it, and one for the atoms in the nearby boxes. The first data set stores all the information, such as the positions, velocities, neighbor lists, etc. These data are stored in a linked list to allow insertions and deletions as the atoms move to new boxes. The second data structure maintains and updates, through communication with other processors, only the atoms' positions. The communication scheme for the SD algorithm is shown in Figure 9.9. In the first step each processor exchanges information with its neighbors in the east-west direction. Thus, processor 2 (see Figure 9.9) fills a message buffer with atoms' positions that it owns that are within



FIGURE 9.9. The spatial decomposition algorithm. In six data exchanges all atom positions in adjacent boxes in the east/west, north/south, and up/down directions are communicated (after Plimpton, 1995).

a distance $r_e$ of processor 1's box, which will be all of its atoms if its dimension $l$ in the east/west direction is less than $r_e$. Otherwise, it will be those that are nearest to processor 1. Then, each processor sends its message in the westward direction and receives a message from the eastward direction which puts it in its second data structure. Then, the process is reversed with each processor sending a message to the east and receiving one from the western processor. If $l > r_e$, all needed atomic positions in the east-west direction are acquired with one such exchange. If, however, $l < r_e$, the east-west steps are repeated, with each processor sending more needed atomic positions to its adjacent processor. This process is repeated until each processor knows all atoms' positions within a distance $r_e$ of its box (dashed boxes in Figure 9.9).

The same procedure is then repeated in the north/south direction. The only difference with the east-west step is that the messages that are now sent to the adjacent processor contain not only the atoms that belong to the processor in its first data structure, but also any atom positions in its second data structure that are needed for the neighboring processor. If, for example, $l = r_e$, this has the effect of sending three boxes worth of atom positions in one message. In the final step, the process is repeated in the up-down direction in which the atoms' positions from an entire plane of boxes are sent in each message. This algorithm has several advantages, some of which are as follows.

(1) If $l \geq r_e$, all the needed atom positions from the 26 surrounding boxes are obtained in only 6 data exchanges. Moreover, if the topology of the parallel machine is that of a simple-cubic lattice, the processors can be mapped onto the boxes in such a way that all six of these processors will be directly connected to the center processor, as a result of which the passage of information will be very efficient.
(2) If $l < r_e$, then atom information is needed from more distant boxes, but the information is obtained by only a few data exchanges, all of which are still with the 6 immediate neighbor processors.
(3) The amount of data communication is minimized, since every processor obtains only the atom positions that are within a distance $r_e$ of its box.
(4) No time is spent for rearranging the data structures. All the newly-arrived information are placed as contiguous data directly into the processor's second data structure. The only time spent is for creating the buffered messages that must be sent.
(5) Message creation is done quickly. The two data structures are scanned only once after every few time steps, when the neighbor lists are created, in order to decide which atom positions to send in each message.

Suppose that box number $m$ is assigned to processor $P_m$ which stores the positions of its $N/N_p$ atoms in $\mathbf{r}_m$, the first data structure, and the forces acting on these atoms in $\mathbf{f}_m$. The details of the procedure for implementing the SD algorithm is as follows.

(1) The positions, velocities, and any other information about atoms that are no longer inside box $m$ are moved from $\mathbf{r}_m$ to a message buffer. These atoms are

then exchanged with the six adjacent processors (see above), during which processor $P_m$ checks for new atoms that are now inside its box, which are then added to $\mathbf{r}_m$. Then, all atom positions that are within a distance $r_e$ of box $m$ are obtained by the communication scheme described above. Since the different messages are buffered by scanning through the two data structures, lists of included atoms are constructed, after which the two data structure of the processor are updated. Then, neighbor lists of the $N/N_p$ atoms of the processor are constructed. If two atoms $i$ and $j$ are in the same box, then the pair $(ij)$ is stored once in the neighbor list. If $i$ and $j$ are in different boxes, their corresponding processors both store them in their respective neighbor lists. If $l < 2r_e$, it is more efficient to find neighbor interactions by checking each atom inside box $m$ against all the atoms in both data structures of the processor, an operation that scales as $O(\sqrt{N/N_p})$. If $l > 2r_e$, it is more efficient to construct the list by binning described above, by mapping all the atoms in both data structure onto bins of size $r_e$, and checking the surrounding bins of each atom in box $m$ for possible neighbors.

(2) Processor $P_m$ uses the neighbor lists to compute all the forces that act on its atoms. If both atoms $i$ and $j$ are inside the same box, then the computed force between them is stored twice in $\mathbf{f}_m$, once each for $i$ and $j$. If the two atoms are in different boxes, only the force acting on the processor's own atom is stored.

(3) After computing $\mathbf{f}_m$, the atoms' positions are updated.

(4) The updated positions are communicated to the surrounding processors in preparation for MD simulation in the next time step.

The cost of these steps scales as the volume of the data exchanged. For example, for the last step, assuming that we have uniform density in the simulation cell, the cost is proportional to the physical volume of the shell of thickness $r_e$ around box $m$, which is $(l + 2r_e)^3 - l^3$. Thus, three cases must be considered. If $l < r_e$, data from many neighboring boxes must be exchanged and the operation cost scales as $8r_e^3$. If $l \simeq r_e$, the data in all the 26 surrounding boxes are exchanged and the cost of operation scales as $27N/N_p$ (note that $N/N_p$ is roughly the number of atoms in volume $l^3$). Finally, if $l \gg r_e$, only atoms' positions near the six faces of box $m$ will be exchanged, and therefore the cost of the communications scales as the *surface area* of the box, namely, $6r_e(N/N_p)^{2/3}$. As before, one can take advantage of Newton's third law and make the algorithm still more efficient.

### 9.6.2.4   Load Balance in Massively-Parallel Molecular Dynamics Simulation

We must point out that all the above algorithms have their optimal efficiency only if there is load-balance, which means that each processor must have an equal amount of work to perform. For the AD algorithms this will be the case if each $\mathbf{F}_m$ or $\mathbf{G}_m$ block has roughly the same number of non-zero entries, which will be the case if the atom density is uniform across the simulation cell. Typically, non-uniformities do arise, but they will not pose any problem if the atoms ordering at the beginning of the simulations is randomly permuted, which is equivalent to permuting rows and columns of $\mathbf{F}$ or $\mathbf{G}$.

There will be load-balance for the FD algorithms only if $\mathbf{F}'_m$ and $\mathbf{G}'_m$ are *uniformly* sparse. Thus, if the atoms are ordered geometrically, these two matrices will not be uniformly sparse, since geometrical ordering generates a force matrix with diagonal bands of non-zero entries. The way to achieve uniform sparsity is to randomly permute the atoms ordering.

Finally, the SD algorithms will be load-balanced only if there is roughly the same number of atoms both in the boxes and in their surroundings. This will not be the case if there is non-uniformity in the density of the particles, or if the physical domain is not a rectangular parallelepiped. Hendrickson and Leland (1995) developed a method for load-balancing the SD algorithms by partitioning an irregular domain or a system with non-uniformly-distributed clusters of atoms, although such algorithms add to the computational cost of the MD simulations.

### 9.6.2.5    Selecting a Massively-Parallel Molecular Dynamics Algorithm

How does one select the most efficient MD algorithm for a given problem? Plimpton (1995) provides some guidelines.

(1) If the communications cost is expected to be negligible, which is the case when $N_p$ is small, then algorithm AD is preferable.

(2) For all other cases algorithm FD will be faster than the AD, both of which scale linearly with $N$ for a fixed $N_p$. If we double $N_p$, the communications cost of the AD algorithm remains the same, while that of the FD algorithm decreases by a factor of $\sqrt{2}$. Thus, as $N_p$ increases and becomes large, the FD algorithm becomes much faster than the AD algorithm.

(3) For a given $N_p$, the scaling of the SD algorithm is not linear with $N$. For large $N$, the efficiency of the algorithm is optimal and nearly 100%, whereas for small $N$ the efficiency is poor. Thus, compared to a FD algorithm, there must be a crossover point when, with increasing $N$ (holding $N_p$ fixed), the efficiency of the SD algorithm exceeds that of the FD algorithm.

Plimpton (1995) tested the efficiency of these algorithms on a benchmark problem, namely, simulation of $N$ atoms, assumed to interact through a LJ potential. The atoms were placed in a 3D parallelepiped with periodic boundary conditions, at a reduced density 0.8442 and reduced temperature 0.72. This represents a liquid state near the LJ triple point. The simulations were begun by placing the atoms on an FCC lattice. Note that the unit cell of this lattice contains 4 atoms. The particles were given initial random velocities (see Section 9.2.4), and after some time the solid (crystal lattice) melted quickly, as it should, and the system of atoms evolved toward its natural liquid state. The cutoff distance was taken to be $r_c = 2.5\sigma$ ($\sigma$ is the size parameter of the LJ potential), the integration time step, in reduced units, was $4.62 \times 10^{-3}$, and the simulations were carried out in the (NVE) ensemble. Table 9.3 compares the performance of the SD1 algorithm on several machines. In this table the dashed sign indicates that the size of the problem was so large that the machine's memory was not sufficient for carrying out the computations. The $10^8$ atoms problem nearly filled the 30 Gbytes of memory on the 1904-processor

TABLE 9.3. CPU seconds/time step for the SD1 algorithm. $N$ is the number of atoms used in the MD simulations, while $N_p$ is the number of the processors of the machine (adapted from Plimpton, 1995).

| Problem Size | | Cray T3D | | Intel iPSC/860 | | Intel Paragon | |
|---|---|---|---|---|---|---|---|
| $N$ | Lattice | $N_p = 256$ | $N_p = 512$ | $N_p = 32$ | $N_p = 64$ | $N_p = 1024$ | $N_p = 1904$ |
| $5 \times 10^2$ | 5 x 5 x 5 | 0.00432 | 0.00446 | 0.0129 | 0.0106 | 0.00564 | 0.00634 |
| $5 \times 10^4$ | 20 x 25 x 25 | 0.0289 | 0.0167 | 0.420 | 0.224 | 0.0174 | 0.0125 |
| $5 \times 10^6$ | 100 x 100 x 125 | 1.86 | 0.994 | – | – | 0.914 | 0.504 |
| $10^8$ | 250 x 250 x 400 | – | – | – | – | – | 9.11 |

Paragon with neighbor lists consuming the majority of the space. It is clear that the algorithm is extremely efficient. However, use of optimized version of the same algorithm increases its efficiency by a factor $2 - 3$. For example, on the Intel Paragon machine with 1840 processors and an assembler version of the algorithm (as opposed to, for example, the Fortran version), the CPU seconds/time step was 5.5, a factor of nearly 2 more efficient than what is listed in Table 9.3. Similar computations were performed for the other algorithms discussed above.

Let us point out that, similar to the classical MD simulations, the electronic structure calculations and QMD computations can also benefit tremendously from parallelization. In particular, parallel QMD can be a powerful method for computing realistic interaction potentials for use in the classical MD simulations (see, for example, Clarke *et al.*, 1992). This is the subject of the next section.

## 9.7    Interatomic Interaction Potentials

Although quantum MD simulation has proven to be a highly successful method for investigating and predicting materials properties at very small scales, due to the huge computations that are necessary, they cannot be used yet for materials at scales that involve several thousands atoms. For this reason, classical MD simulations are a valuable tool for studying various properties of materials using millions of atoms. However, as discussed in Section 9.2, accurate interaction potentials are critical to the success of MD computations. Although simple interaction forces, such as those derived from the LJ potentials, have provided us with much qualitative insight into the properties of various materials, quantitative predictions require much more realistic representation of the interactions between the atoms. Moreover, only for noble gases can one represent the interactions between atoms by density-*independent*, pairwise-additive forces, and the repulsive and attractive forces are due to spherical electron clouds that are close to the nuclei. The alternative to the LJ and similar potentials are semi-empirical expressions designed for accurately describing small distortions from the ground state in more complex systems, a famous example of which is the Keating potential described in Chapter 8 of Volume I for covalent bonds. Such potentials, which can be viewed as a sort of Taylor series expansions of the energy about its minimum, are useful for describing phonons and elastic deformations, but they are incapable of describing the energy

of states which differ significantly from tetrahedral ground state, or when one must deal with large deformations. For these reasons, developing accurate representation of the interaction forces between various atoms has been, for many years, an active research field. We summarize in this section some of the most significant results that have emerged over the last two decades.

### 9.7.1   The Embedded-Atom Model

The embedded-atom model (EAM), which is intended for metals, was developed by Daw and Baskes (1983, 1984). In metals, electrons are not all localized around the nuclei, rather the valence electrons are often shared among many ions, similar to a nearly free-electron gas. This implies that the energy depends upon the local electron density, resulting in many-body, rather than pairwise, forces between ions, hence allowing one to represent the interactions between ions in metals by a relatively-simple approximate functional form, commonly referred to as the embedded-atom potential. In this approximation, the total potential energy $E$ of $N$ ions in an arbitrary volume $\Omega$ is given by

$$E = \sum_{i=1}^{n} \left[ \frac{1}{2} \sum_{j \neq i} U_{ij}(r_{ij}) + E_i^e(\rho_i) \right], \tag{159}$$

where $U(r_{ij})$ is a density-*independent*, pairwise-additive and short-range interaction potential that depends only on distance $r_{ij}$ between particles $i$ and $j$, and $E_i^e$ is the embedding energy that depends on the local embedding density $\rho_i$ at atom $i$. In effect, each atom is viewed as an impurity embedded in a host consisting of all other atoms, such that the embedding energy depends on the electron density. In this sense, the basic idea of the EAM is, on one hand, similar to the effective-medium approximation described in the previous chapters, and, on the other hand, similar to the DFT described in this chapter. If one makes a further simplification by assuming that the density $\rho_i$ can be approximated by

$$\rho_i = \sum_{j \neq i} \rho_j^a(r_{ij}), \tag{160}$$

where $\rho_i^a$ is the atomic density of the constituents, then the energy would be a simple function of the atoms' positions. Note the difference between $\rho_j$ and $\rho_j^a$: Whereas $\rho_j^a$ is the contribution to the density *from* atom $j$, $\rho_j$ is the total electron density *at* atom $j$.

As shown by Daw and Baskes (1984), the ground state properties of the solid can be computed from Eq. (159). For example, consider the case of a perfect, homonuclear crystal. In this case, $E_i^e = E$, $U = U_{ij}(r_{ij})$, and $\rho = \rho^a$. If $\rho_e$ is the equilibrium density, then $\rho_e = \sum_m \rho(l^m)$, where $l^m$ are the distances between neighbors, and the sum is over neighbors. Moreover, one has, for every $i$, $\rho_i = \rho_e$. Then, the lattice constant is obtained from the equilibrium condition:

$$\frac{1}{2} \sum_m \left( \frac{U_m' l_i^m l_j^m}{l^m} \right) + E'(\rho_e) \sum_m \left( \frac{\rho_m' l_i^m l_j^m}{l^m} \right) = 0, \tag{161}$$

where $l_i^m$ is the $i$th component of the position vector to the $m$th neighbor, $U_m' = dU/dr$, and $\rho_m' = d\rho/dr$, with the subscript $m$ indicating that the derivatives are to be evaluated at $r = l^m$.

The elastic constants of the crystal at equilibrium are given by

$$C_{ijkl} = \frac{1}{\Omega_0}[B_{ijkl} + E'(\rho_e)W_{ijkl} + E''(\rho_e)A_{ij}A_{kl}], \qquad (162)$$

where

$$B_{ijkl} = \frac{1}{2}\sum_m \left[\frac{(U_m'' - U_m'/l^m)l_i^m l_j^m l_k^m l_l^m}{(l^m)^2}\right], \qquad (163)$$

$$W_{ijkl} = \sum_m \left[\frac{(\rho_m'' - \rho_m'/l^m)l_i^m l_j^m l_k^m l_l^m}{(l^m)^2}\right], \qquad (164)$$

$$A_{ij} = \sum_m \left(\frac{\rho_m' l_i^m l_j^m}{l^m}\right), \qquad (165)$$

and $\Omega_0$ is the volume of the undeformed crystal. In particular, for a cubic crystal, the three independent elastic constants are given by

$$C_{11} = \frac{1}{\Omega_0}[B_{11} + E'(\rho_e)W_{11} + E''(\rho_e)A_{11}^2], \qquad (166)$$

$$C_{12} = \frac{1}{\Omega_0}[B_{12} + E'(\rho_e)W_{12} + E''(\rho_e)A_{11}^2], \qquad (167)$$

$$C_{44} = \frac{1}{\Omega_0}[B_{12} + E'(\rho_e)W_{12}]. \qquad (168)$$

These equations nicely demonstrate the effect of the interplay between the pair potential $U_{ij}$ and the embedding energy $E^e$. Clearly, if we remove $U_{ij} = U$, we obtain $C_{11} = C_{12}$ and $C_{44} = 0$, which, for real solid materials, are wrong. On the other hand, if we remove the embedding energy $E^e = E$, we obtain the well-known Cauchy relation, $C_{12} = C_{44}$, which does not hold for all materials, but only for a certain class of them.

To use the EAM one must have the embedding energy, the pair potential, and the atomic densities. In their original work, Daw and Baskes (1984) and Foiles *et al.* (1986) used semi-empirical correlations for evaluating these quantities. To develop such correlations, one takes advantage of known properties of these quantities. For example, the embedding energy, when defined relative to the free-atom energy, must vanish for vanishing electron density, and must have negative slope and positive curvature (second derivative) for the background electron densities found in typical metals. On the other hand, the pair-interaction term in Eq. (159) is purely repulsive, and its origin is Coulombic. These observations lead to the following equation

$$U_{ij}(r) = \frac{Z_i(r)Z_j(r)}{r}, \qquad (169)$$

where $Z_i(r)$ is the effective charge of atom $i$. Note that $Z_i(r)$ must be positive and decrease monotonically with increasing $r$. A particularly simple, yet accurate, expression is given by

$$Z(r) = Z_o(1 + a_1 r^{a_2}) \exp(-a_3 r), \tag{170}$$

where $Z_o$ is the number of the outer electrons in the atom (for example, $Z_o = 10$ for Ni, Pd, and Pt, and $Z_o = 11$ for Cu, Ag, and Au). The parameters $a_1$, $a_2$ and $a_3$ must be determined empirically, although $a_2 = 1$ is accurate for Ni, Pd, and Pt, while $a_2 = 2$ leads to accurate elastic constants for Cu, Ag, and Au. Using these observations and properties, Daw and Baskes (1984) and Foiles *et al.* (1986) obtained very accurate empirical correlations for the embedding energy and the effective charges for a variety of metals, from which they computed their various properties, such as their elastic constants and surface energy; see also Johnson (1988) who utilized the EAM to study FCC metals.

Holian *et al.* (1991) developed the following EAM for use in MD simulation of deformation of materials under high stresses. The local embedding density $\rho_i$ was assumed to be given by a pairwise sum over all neighboring particles, weighted by a spherical localization function $w(r_{ij})$, such that

$$\rho_i = \sum_{j \neq j} w(r_{ij}), \tag{171}$$

where

$$w(r) = \frac{1}{ed(d+1)} \left( \frac{r_c^2 - r^2}{r_c^2 - r_e^2} \right)^2, \tag{172}$$

where $d$ is the dimensionality of the system, $e = \exp(1)$, $r_c$ is the cutoff distance given by Eq. (34), and $r_e$ is the equilibrium nearest-neighbor distance. The pairwise interaction potential $U(r)$ was taken to be that given by Eq. (33), except that the LJ part of the potential was written as

$$U_{LJ}(r) = 4\chi\varepsilon \left[ \left( \frac{\sigma}{r} \right)^{12} - \left( \frac{\sigma}{r} \right)^6 \right], \tag{173}$$

where $\chi$ is the fractional pair-potential contribution to the total cohesive energy ($\chi = 1/3$ is a reasonable value for many metals). Finally, the embedding energy $E_i^e$ was taken to be

$$E_i^e(\rho_i) = \frac{1}{2} e\varepsilon d(d+1)(1-\chi)\rho_i \ln \rho_i. \tag{174}$$

One can also replace the LJ potential by the more flexible Morse-like potential,

$$U_{\text{Morse}} = \varepsilon\{\exp[-\alpha(r/r_e - 1)] - 1\}^2 - \varepsilon, \tag{175}$$

where $\alpha$, the steepness of the repulsive well, is related to $\Omega_0$, the volume of the undeformed system, and the bulk modulus of the system at equilibrium.

However, the atomic densities can also be computed by the Hartree–Fock approximation, and the embedding energy and the pair potentials can be calculated by the *ab initio* method. The *ab initio* electronic structure calculations are utilized

in order to generate an accurate database of configurational energies, which are then used for determining the necessary parameters of the EAM (see, for example, Wagner *et al.*, 1992), which take advantage of analytical functional forms obtained from theoretical analysis. The *ab initio* electronic structure calculations are formulated in the context of the quantum mechanical DFT described in Sections 9.1, 9.4 and 9.5.

Baskes *et al.* (1989) proposed a potential based on the EAM by introducing an angular dependence in the host electron energy, which could then adequately account for the bond-bending forces in materials such as diamond. The potential could fit the energies of the high-density polymorphs, with the goodness of the fit being comparable with those of the other potentials described below. It could also describe exactly the static properties of cubic diamonds, and provide a relatively accurate description of bulk defects and vacancies. However, it provides high values for the energy of the intrinsic stacking faults in silicon.

Most recently, Web and Grest (2001) utilized the EAM and its modification to compute the surface tension of liquid-vapor interface of metals. The agreement between the predicted values and the experimental data was very good, hence indicating that the EAM can also describe fluid interface properties of materials.

### 9.7.2   The Stillinger–Weber Potential

Stillinger and Weber (1985) proposed a semi-empirical potential for silicon which appears to be relatively accurate. In their model the total interatomic potential involves two- and three-body contributions and is written as

$$E = \sum_{i<j} U_2(r_{ij}) + \sum_{i,j<k} U_3(\mathbf{r}_{ij}, \mathbf{r}_{ik}), \qquad (176)$$

where $U_2$, the two-body term, can include such effects as the steric repulsion, charge transfer between atoms, charge-dipole and van der Waals interactions, and therefore

$$U_2(r_{ij}) = \frac{H_{ij}}{r_{ij}^{n_{ij}}} + \frac{Z_i Z_j}{r_{ij}} \exp(-r_{ij}/a) - \frac{D_{ij}}{r_{ij}^4} \exp(-r_{ij}/b) - \frac{W_{ij}}{r_{ij}^6}. \qquad (177)$$

The first term of Eq. (177) represents a two-parameter representation of the steric repulsion; the second term is the Coulombic interaction due to charge transfer and contains the effective atomic charges $Z_i$ and $Z_j$; the third term takes into account the charge-dipole interaction due to large polarizability of negative ions, while the last term corresponds to the induced dipole-dipole interactions. Covalent effects are taken into account through three-body bond-bending and bond-stretching terms (similar to the Keating model described in Chapter 8 of Volume I), and includes the Si-C as well as C-Si-C bond angles. It is given by

$$U_3(\mathbf{r}_{ij}, \mathbf{r}_{ik}) = B_{ijk} \exp\left(\frac{\gamma}{r_{ij} - r_0} - \frac{\gamma}{r_{ik} - r_0}\right) \frac{(\cos\theta_{jik} - \cos\bar{\theta}_{jik})^2}{1 + C(\cos\theta_{jik} - \cos\bar{\theta}_{jik})^2} \Theta(r_0 - r_{ij})\Theta(r_0 - r_{ik}).$$
$$(178)$$

TABLE 9.4. Comparison of the MD predictions of the properties of cubic SiC in zinc-blende structure, obtained by using the Stillinger–Weber potential, with the experimental data. The elastic constants are in GPa (adapted from Shimojo *et al.*, 2000).

| Property | MD | Experiment |
|---|---|---|
| Lattice constant (Å) | 4.36 | 4.36 |
| $C_{11}$ | 390 | 390 |
| $C_{12}$ | 144 | 142 |
| $C_{44}$ | 179 | 150–256 |
| Bulk modulus | 225 | 225 |
| Melting temperature (°C) | 3000 | 2830 |



FIGURE 9.10. Volume-pressure relation for zinc-blend structure. Solid circles shows the MD results, while open circles are the experimental data (after Shimojo *et al.*, 2000).

Here, $B_{jik}$ is the strength of the interaction, $\Theta(r_0 - r_{ij})$ is the step function, $\theta_{jik}$ is the angle formed by $\mathbf{r}_{ij}$ and $\mathbf{r}_{ik}$, and $\bar{\theta}_{jik}$ is a constant. The constant $C$ plays an important role if the material undergoes structural transformation. In the original Stillinger–Weber formulation, $\cos\bar{\theta}_{jik}$ was assumed to have a value of $-1/3$, but more generally one can treat this term as an adjustable parameter. Note that the "ideal" tetrahedral angle $\theta$ is such that $\cos\theta = -1/3$, so that the trigonometric part of Eq. (178) clearly discriminates in favor of pairs of bonds emanating from $i$ with the desired geometry. It is clear that the Stillinger–Weber potential contains many parameters which must be estimated by fitting the predictions obtained with the potential to certain properties of the material. Shimojo *et al.* (2000) employed this potential in their MD simulations of cubic SiC under isothermal-isobaric conditions (see Section 9.2), using only 1728 atoms. Table 9.4 compares the predictions of the MD simulations with the experimental data, and it is clear that the agreement between the two sets is excellent. In addition, the volume-pressure relation was computed for SiC in the zinc-blende structure, i.e., in a configuration with 4-fold coordination. Figure 9.10 compares the MD simulation results with the experimental data, and it is clear that the agreement is again excellent, hence demonstrating the significance and utility of an accurate interatomic potential: If the interatomic potential is accurate, then MD simulations provide quantitative predictions for materials' properties.

The Stillinger–Weber potential is, by far, the most widely used potential. It has been utilized in the study of clusters, lattice dynamics, bulk point defects, the liquid and amorphous states, surface diffusion and reconstructions, Si(100) stepped surfaces, the liquid-vapor and crystal-melt interfaces, pulsed melting of surfaces, epitaxial growth from the vapor, liquid-phase epitaxy, and growth of amorphous films via atom deposition, as well as calculation of mechanical properties that was mentioned above. It has also been extended to Ge, sulfur, fluorine, and Si-F materials.

Mistriotis *et al.* (1989) modified the Stillinger–Weber potential in order to correctly describe clusters with more than 6 atoms. The angular dependence of the three-body term was modified, and a four-body term was also added.

## 9.7.3   The Tersoff Potentials

It can be shown, by simple quantum-mechanical arguments, that the more neighbors an atom has, the weaker the bond to each neighbor will be. In general, the bond strength, or bond order, depends in a complex way on the geometry of the material. For example, even-membered rings might be favored over odd ones. However, the most important single variable is the coordination number—the number of neighbors close enough to form bonds. If the energy per bond decreases sufficiently rapidly with increasing coordination, then the diatomic molecule will be the most stable arrangement of atoms. Low coordination numbers are common for atoms at the far right of the Periodic Table (especially near the top). However, if the bond order depends only weakly on the coordination number, then close-packed structures form so as to maximize the number of bonds. This limit corresponds, roughly speaking, to metallic rather than covalent bonding, and is found for atoms at the left and bottom of the Periodic Table, with a trend in between from low coordination numbers at the upper right to high coordination number at the lower left. Thus, bond order is a monotonically decreasing function of the coordination number, and a trade off between this property and number of bonds determines the equilibrium conditions.

Silicon, aside from its technological importance, is a remarkable material since, even with modest changes of pressure, it can take on structures with a large range of coordination. This is due to the fact that a decrease in bond strength with increasing coordination number essentially cancels the increase in the number of bonds, over a large range of coordination number. As such, silicon provides a stringent test of our ability for describing the dependence of bonding upon coordination number, and hence our ability for developing potentials that can accurately describe its structure.

In an important paper, Tersoff (1988) suggested such a potential for silicon. Because of the critical role of bond order and its dependence upon local geometry, one must include an environment-dependent bond order in the potential. Thus, in Tersoff's formulation the total interatomic potential energy is taken to be

$$E = \sum_i E_i = \frac{1}{2} \sum_i \sum_{j \neq i} E_{ij}, \tag{179}$$

where

$$E_{ij} = U_c(r_{ij})[a_{ij}U_r(r_{ij}) + b_{ij}U_a(r_{ij})]. \tag{180}$$

Here $E_i$ and $E_{ij}$ are, respectively, a site and a bond energy, the sums are over the atoms of the system, and $r_{ij}$ is the distance between atoms $i$ and $j$. The function $U_r$ represents a repulsive pair potential, which includes the orthogonalization energy when atomic wave functions overlap (see Section 9.4), while $U_a$ is the attractive pair potential associated with bonding. As already emphasized throughout this chapter, in many applications short-ranged functions allow a tremendous reduction in computational effort, and therefore a cutoff function $U_c$ has been introduced to limit the range of the potential. The *function* $b_{ij}$ is a measure of the bond order, and is assumed to be a monotonically decreasing function of the coordination of atoms $i$ and $j$. In addition, terms that act to limit the range of interaction to the first neighbor shell are included in $b_{ij}$, and the function $a_{ij}$ consists solely of such range-limiting terms.

Ferrante *et al.* (1983) showed that a large number of calculated binding-energies for solid cohesion and chemisorption can be mapped onto a single dimensionless function using a three-parameter scaling, and Abell (1985) showed that this universal behavior can be well-explained by use of a Morse or Morse-like pair potential, Eq. (175). Therefore, Tersoff (1988) proposed the following expressions for $U_r$ and $U_a$:

$$U_r(r) = A \exp(-\lambda_1 r), \tag{181}$$

$$U_a(r) = -B \exp(-\lambda_2 r), \tag{182}$$

whereas the cutoff function $U_c$ was taken to be

$$U_c(r) = \begin{cases} 1, & r < R - D \\ \frac{1}{2} - \frac{1}{2}\sin[\pi(r-R)/(2D)], & R - d < r < R + D \\ 0, & r > R + D \end{cases} \tag{183}$$

The cutoff function (and its derivative) is continuous for all $r$, and varies between 0 and 1 in a small range around $R$, which is selected so as to include only the first-neighbor shell for most structures. In effect, the potential has the form of a Morse pair potential, ignoring the three-body and higher-order effects, but with a bond-order parameter $b_{ij}$ that depends on the local environment. The function $b_{ij}$ is given by

$$b_{ij} = (1 + \beta \zeta_{ij}^n)^{-1/2n}, \tag{184}$$

with

$$\zeta_{ij} = \sum_{k \neq i, j} U_c(r_{ik}) g(\theta_{ijk}) \exp[\lambda_3^3 (r_{ij} - r_{ik})^3], \tag{185}$$

$$g(\theta) = 1 + c^2\{d^{-2} - [d^2 + (h - \cos\theta)^2]^{-1}\}, \tag{186}$$

where $\theta_{ijk}$ is the angle between bonds $ij$ and $ik$. The function $\cos\theta_{ijk}$ is used to ensure the proper analytic behavior for the dependence of $b_{ij}$ on $\theta_{ijk}$. Note that

$b_{ij} \neq b_{ji}$ which, however, does not have any significant consequence. If, however, one insists on a more symmetric form, the sum over pairs of atoms in Eq. (180) can be replaced with a sum over bonds $(i > j)$, and then $b_{ij}$ can be replaced with a symmetric function, $\bar{b}_{ij} = \frac{1}{2}(b_{ij} + b_{ji})$. Tersoff (1988) also proposed the following equation for the function $a_{ij}$,

$$a_{ij} = (1 + \alpha^n \eta_{ij}^n)^{-1/2n}, \tag{187}$$

$$\eta_{ij} = \sum_{k \neq i, j} U_c(r_{ik}) \exp[\lambda_3^3 (r_{ij} - r_{ik})^3], \tag{188}$$

where $\alpha$ is typically small enough that $a_{ij} \simeq 1$, unless, of course, $\eta_{ij}$ is exponentially large.

Subsequently, Tersoff (1989) modified his proposed potential in order to describe multicomponent mixtures, and more specifically C-Si and Si-Ge mixtures. Silicon carbide, in particular, has a wide range of applications, ranging from optoelectric devices and engineering materials to the basic substrate for membranes that must operate at high-temperatures (Suwanmethanond *et al.*, 2000). The reason for its popularity is that it has excellent chemical stability, good electronic properties, and high stiffness and hardness. In Tersoff's generalization, Eqs. (179) and (180) remain the same, but the remaining expressions are modified to account for the multicomponent nature of the system. Thus,

$$U_r(r_{ij}) = A_{ij} \exp(-\lambda_{ij} r_{ij}), \tag{189}$$

$$U_a(r_{ij}) = -B_{ij} \exp(-\mu_{ij} r_{ij}), \tag{190}$$

$$U_c(r_{ij}) = \begin{cases} 1, & r_{ij} < R_{ij} \\ \frac{1}{2} + \frac{1}{2}\cos[\pi(r_{ij} - R_{ij})/(S_{ij} - R_{ij})], & R_{ij} < r_{ij} < S_{ij} \\ 0, & r_{ij} > S_{ij} \end{cases} \tag{191}$$

where the various parameters are given by

$$b_{ij} = \chi_{ij}(1 + \beta_i^{n_i} \zeta_{ij}^{n_i})^{-1/2n_i}, \tag{192}$$

$$\zeta_{ij} = \sum_{k \neq i, j} U_c(r_{ik}) \omega_{ik} g(\theta_{ijk}), \tag{193}$$

$$g(\theta_{ijk}) = 1 + c_i^2 \{d_i^{-2} - [d_i^2 + (h_i - \cos\theta_{ijk})^2]^{-1}\}, \tag{194}$$

$$\lambda_{ij} = \frac{1}{2}(\lambda_i + \lambda_j), \quad \mu_{ij} = \frac{1}{2}(\mu_i + \mu_j), \tag{195}$$

$$A_{ij} = \sqrt{A_i A_j}, \quad B_{ij} = \sqrt{B_i B_j}, \tag{196}$$

$$R_{ij} = \sqrt{R_i R_j}, \quad S_{ij} = \sqrt{S_i S_j}. \tag{197}$$

Note that the parameter $\chi_{ij}$ strengthens or weakens the heteropolar bonds, and therefore represents in some sense the chemistry of the mixture. Note also that $\chi_{ii} = 1$ and $\chi_{ij} = \chi_{ji}$, and that $\omega_{ii} = 1$, although experience has indicated that $\omega_{ij} = 1$ is also a reasonable estimate. Compared with the case of a pure component,

TABLE 9.5. Estimates of the parameters for carbon, silicon and
germanium to be used in the Tersoff potential. Except for $R$ and $S$, all
the parameters have been optimized (adapted from Tersoff, 1989).

| Parameter | C | Si | Ge |
|---|---|---|---|
| $A$ (eV) | 1393.6 | 1830.8 | 1769.0 |
| $B$ (eV) | 346.7 | 471.18 | 419.23 |
| $R$ (Å) | 1.8 | 2.7 | 2.8 |
| $S$ (Å) | 2.1 | 3.0 | 3.1 |
| $c$ | 38049.0 | 100390.0 | 106430.0 |
| $d$ | 4.384 | 16.217 | 15.652 |
| $h$ | $-0.57058$ | $-0.59825$ | $-0.43884$ |
| $n$ | 0.72751 | 0.78734 | 0.75627 |
| $\beta$ | $1.5724 \times 10^{-7}$ | $1.1000 \times 10^{-6}$ | $9.0166 \times 10^{-7}$ |
| $\lambda$ (Å$^{-1}$) | 3.4879 | 2.4799 | 2.4451 |
| $\mu$ (Å$^{-1}$) | 2.2119 | 1.7322 | 1.7047 |
|  |  | $\chi_{C-Si} = 0.9776$ | $\chi_{Si-Ge} = 1.00061$ |

TABLE 9.6. Comparison of the MD
predictions of properties of cubic SiC,
obtained with the Tersoff potential,
with the experimental data. The elastic
constants are in Mbar (adapted from
Tersoff, 1989).

| Property | MD | Experiment |
|---|---|---|
| Lattice constant (Å) | 4.32 | 4.36 |
| $C_{11}$ | 4.2 | 3.6 |
| $C_{12}$ | 1.2 | 1.5 |
| $C_{44}$ | 2.6 | 1.5 |

the potential for mixtures is somewhat simpler, as the parameter $D$ of Eq. (183)
has been eliminated.

All the parameters of these potentials have been estimated by fitting them to heat
of formation and the properties of the respective elements. For Si the parameter $D$
is about 0.15Å. The rest of the parameters are listed in Table 9.5. To provide the
reader with some sense of the accuracy of these potentials, we compare in Table
9.6 the predicted lattice constant and the three elastic constants of cubic SiC with
the experimental data. It is clear that, except for $C_{44}$, the agreement between the
theoretical predictions and the data is quite good. In addition, Mura *et al.* (1998)
studied the properties of $Si_{1-x}C_x$ compounds in the range $0.125 \leq x < 0.875$ by
*ab initio* QMD method of Car and Parrinello (1985) (see Section 9.4) and by MD
simulations using the Tersoff potentials, and found very good agreement between
the predictions of the two methods.

However, because of the large difference in bonding characteristics between
hydrogen and carbon (recall that H is monovalent whereas C has a valency of
up to 4), a set of parameters cannot be found that can adequately describe bond
energies for a large number of hydrocarbons. Furthermore, the Tersoff potentials

cannot describe radicals and non-conjugated double bonds. These deficiencies motivated the development of several modified Tersoff potentials which we now briefly mention.

Chelikowsky *et al.* (1989) developed an interatomic potential similar in form to Tersoff's, which was intended for describing the metallic to covalent transition that occurs in clusters when their size reaches a critical size. To describe the transition, an angular-dependent bond-bending force was incorporated in the potential. The resulting potential provides very accurate description of perfect diamond structure, as well as the high-density polymorphs of silicon. To model clusters, a so-called dangling-bond vector was introduced that describes the transfer of bond strength from a dangling bond to a backbone bond. The energies of $Si_n$ clusters with $n <$ 10 are, however, underestimated. Moreover, the ground-state structures for such values of $n$ are also not well-described. The potential was not intended for bulk point defects and for surfaces.

Following Tersoff, Kohr and Das Sarma (1989) developed a series of interatomic potentials for tetrahedrally-bonded semiconductors. The bond-bending term was modified to deal with the larger angular distortions from the tetrahedral angle encountered on the various (111) surfaces. Moreover, since the bonds of a given atom can be of different nature, the value of the effective coordination number was fixed in an *ad hoc* way. The potential provides accurate description of various (111) surfaces.

Bolding and Andersen (1990) generalized Tersoff's potential by expressing the attractive term of the potential as a sum of $\sigma$- and $\pi$-bonding terms. The functional form is complex, and the potential contains over 30 parameters. Bolding and Andersen used a very large data base for fitting the parameters, including the static properties of the cubic diamond phase, the fact that the first pressure-induced phase transformation from cubic diamond is the $\beta$-tin phase, and the energies and geometries of global and local minima for clusters of 2-10 atoms. For small clusters, the Bolding–Andersen potential generates a surface that has most of the local and global minima predicted by the *ab initio* computations. Moreover, its predictions for the ground-state structures and energies are in excellent agreement with those predicted by quantum-mechanical computations. The static properties of cubic diamond silicon are also well described by the potential. For bulk point defects, only the vacancy is predicted to have a formation energy that is in good agreement with the DFT results.

## 9.7.4   The Brenner Potentials

In addition to the shortcomings of the Tersoff potentials for certain materials, there is another impetus for developing more sophisticated empirical or semi-empirical potentials. Chemical vapor deposition (CVD) of diamond films is a process of tremendous technological importance. It is, however, a complex process in which diamond grows under apparently metastable conditions. Various factors, such as addition to the surface of such species as acetylene, methyl radicals, or a mixture of hydrocarbon molecules, the substrate temperature, and initiation of defects, all

may play important roles. In order to obtain a better understanding of this complex process, atomic-scale simulations of this phenomenon is of course desirable. However, its QMD simulation is still too costly (in terms of the computation time), and thus the classical MD computations with an accurate interatomic potential is desirable. The potential developed by Brenner (1990) is intended for this purpose. Its main use is for hydrocarbons, and it can take into account the effect of intramolecular chemical bonding (which the Tersoff potentials are incapable of doing) in many small hydrocarbon molecules as well as graphite and diamond lattices.

In its spirit, Brenner's formulation is similar to Tersoff's. The binding energy $E_b$ for the hydrocarbon potential is given by

$$E_b = \sum_i \sum_{j>i} [U_r(r_{ij}) - \bar{b}_{ij} U_a(r_{ij})], \tag{198}$$

where, similar to the Tersoff potentials, $U_a$ and $U_r$ are the attractive and repulsive part of the energy, and are given by

$$U_a(r_{ij}) = f_{ij}(r_{ij}) \frac{D_{ij}^e S_{ij}}{S_{ij} - 1} \exp\left[-\sqrt{\frac{2}{S_{ij}}} \beta_{ij} \left(r_{ij} - R_{ij}^e\right)\right], \tag{199}$$

$$U_r(r_{ij}) = f_{ij}(r_{ij}) \frac{D_{ij}^e S_{ij}}{S_{ij} - 1} \exp\left[-\sqrt{2S_{ij}} \beta_{ij} \left(r_{ij} - R_{ij}^e\right)\right], \tag{200}$$

where $D_{ij}^e$ is the well depth, and $R_{ij}^e$ is the equilibrium distance, the value of which is the same as $r_e$ in the Morse potential, Eq. (175). In fact, for $S_{ij} = 2$ the attractive and repulsive potentials are essentially identical with the Morse potential. $f_{ij}(r_{ij})$ are cutoff functions given by

$$f_{ij}(r) = \begin{cases} 1, & r \leq R_{ij}^{(1)} \\ \frac{1}{2} + \frac{1}{2} \cos[\pi (r - R_{ij}^{(1)})/(R_{ij}^{(2)} - R_{ij}^{(1)})], & R_{ij}^{(1)} < r < R_{ij}^{(2)} \\ 0, & r \geq R_{ij}^{(2)} \end{cases} \tag{201}$$

It should be clear that the cutoff functions explicitly restrict the interactions to nearest neighbors. The function $\bar{b}_{ij}$ is given by

$$\bar{b}_{ij} = \frac{1}{2}(b_{ij} + b_{ji}) + F_{ij}[N_i^{(t)}, N_j^{(t)}, N_{ij}^{\text{conj}}], \tag{202}$$

implying that, similar to the Tersoff potentials, $\bar{b}_{ij}$ depends on the environment around atoms $i$ and $j$ and implicitly contains many-body information. The function $F_{ij}$ is a correction term which is used only for carbon-carbon bonds, and

$$b_{ij} = \left\{1 + \sum_{k \neq i,j} g(\theta_{ijk}) f_{ik}(r_{ik}) \exp\left(\alpha_{ijk}[(r_{ij} - R_{ij}^e) - (r_{ik} - R_{ik}^e)]\right) + H_{ij}[N_i^{(H)}, N_i^{(C)}]\right\}^{-\delta_i}. \tag{203}$$

Here $N_i^{(C)}$ and $N_i^{(H)}$ are the number of carbon and hydrogen atoms bonded to atom $i$, $N_i^{(t)} = N_i^{(C)} + N_i^{(H)}$, $N_{ij}^{conj}$ depends on whether a bond between carbon atoms $i$ and $j$ is part of a conjugated system, and $\theta_{ijk}$ is the angle between bonds $ij$ and $ik$. Similar to $F_{ij}$, the function $H_{ij}$ is a correction term, and both functions are used only for hydrocarbons.

The cutoff functions $f_{ij}(r)$ are used for defining the various quantities. One has

$$N_i^{(H)} = \sum_{j=\{H\}} f_{ij}(r_{ij}), \tag{204}$$

$$N_i^{(C)} = \sum_{j=\{C\}} f_{ij}(r_{ij}), \tag{205}$$

where, for example, {C} denotes the set of the carbon atoms. Values of $N_i^{(t)}$ for neighbors of two carbon atoms involved in a bond are used for determining whether the bond is part of a conjugated system. For example, if the neighbors are carbon atoms that have a coordination number of less four (i.e., $N_i^{(t)} < 4$), the bond is defined as part of a conjugated system. For a bond between carbon atoms $i$ and $j$,

$$N_{ij}^{conj} = 1 + \sum_{k(\neq i,j)=\{C\}} f_{ik}(r_{ik}) F(x_{ik}) + \sum_{l(\neq i,j)=\{C\}} f_{jl}(r_{jl}) F(x_{jl}), \tag{206}$$

with

$$F(x_{ik}) = \begin{cases} 1, & x_{ik} \leq 2 \\ \frac{1}{2} + \frac{1}{2}\cos[\pi(x_{ik}-2)], & 2 < x_{ik} < 3 \\ 1, & x_{ik} \geq 3 \end{cases} \tag{207}$$

and

$$x_{ik} = N_k^{(t)} - f_{ik}(r_{ik}). \tag{208}$$

The function $F(x_{ik})$ yields a continuous value of $N_{ij}^{conj}$ as bonds form and break and as second-neighbor coordinations change. If $N_{ij}^{conj} = 1$, a bond is not part of a conjugated system and the function yields appropriate values, while for $N_{ij}^{conj} \geq 2$ the bond is part of a conjugated system and parameters fitted to conjugated bonds are used. Finally, to make the potential continuous, 2D and 3D cubic splines are utilized for the functions $H_{ij}$ and $F_{ij}$, respectively, for interpolating between values at discrete numbers of neighbors. The function $g(\theta)$ is very similar to that in the Tersoff potentials, Eqs. (186) and (194). For example, for carbon one has

$$g_C(\theta) = a_0 \left[ 1 + \frac{c_0^2}{d_0^2} - \frac{c_0^2}{d_0^2 + (1 + \cos\theta)^2} \right]. \tag{209}$$

It is clear that the Brenner potential has a large number of parameters that must be fitted to experimental data. The procedure for doing this is to first fit to systems consisting of only carbons and hydrogen. Parameters are then selected for the mixed hydrocarbon system that produce additive bond energies.

Because the pair terms are first fitted to solid-state carbon structures, the equilibrium carbon-carbon distances and the stretching force constants for hydrocarbons are completely determined by fitting to bond energies. To determine appropriate energies for hydrocarbons with carbon-carbon bonds, additive bond energies for single, double, conjugated double and triple carbon-carbon bonds, and carbon-hydrogen bonds are determined from molecular atomization energies. Values of the parameter $\delta_i$ of Eq. (203) for carbon and hydrogen turn out to be identical and equal to 0.80469. Values of the other parameters are listed by Brenner (1990), a list too long to be given here. The Brenner potential has been utilized successfully in studying many phenomena, including reaction of hydrocarbon species on diamond surfaces, mechanical properties of graphite sheet and nanotubes, and CVD of diamond films. Moreover, excellent agreement was found (Robertson *et al.*, 1992) between the predictions of MD simulations of energetics of nanoscale graphitic tubules using the Brenner potential, and those of first-principle electronic structure calculations using local-density functional described in Sections 9.1 and 9.4.

Despite its success, the Brenner potential does have certain limitations. For example, it does not include the resonance effect of aromatics. Most importantly, the long-range van der Waals and Coulombic interactions are not included explicitly in the model, although such interactions play an important role in many materials. Che *et al.* (1999) modified the Brenner potential to take such effects into account. In their formulation, the total energy of the system is written as

$$E = \sum_i \sum_{j>i} \left[ U_{ij}^{V}(r_{ij}) + P_{ij}(r_{ij}) U_{ij}^{NB}(r_{ij}) \right]. \tag{210}$$

Here superscript V denotes the valence short-range terms (for example, those in the Brenner potential), NB indicates the long-range non-bond part of the energy (for example, the contribution of van der Waals or Coulombic forces), while $P_{ij} = P_{ji}$ is a screening function that properly weights the NB contributions to the total energy, and is given by

$$P_{ij} = f\left(U_{ij}^{V}, U_{ij}^{V}\right) \prod_{k \neq i,j} f\left(U_{ik}^{V}, U_{kj}^{V}\right), \tag{211}$$

with

$$f(x, y) = \begin{cases} \exp(-x^2 y^2), & \text{if } x < 0 \text{ and } y < 0 \\ 1, & \text{otherwise.} \end{cases} \tag{212}$$

As pointed out by Che *et al.* (1999), the screening function eliminates NB interactions between two atoms $i$ and $j$ when they are directly bonded (i.e., $U_{ij}^{V} < 0$) or when they are both bonded to a common atom $k$ (i.e., $U_{ik}^{V} < 0$ and $U_{kj}^{V} < 0$). In both cases the screening function $P_{ij}$ is negligibly small, and therefore the NB interactions do not make improper contribution to the total energy. Using the relation, $F_{\alpha\beta} = -\partial E/\partial r_{\alpha\beta}$, it is now straightforward to show that the force $F_{\alpha\beta}$

between atoms $\alpha$ and $\beta$ is given by

$$F_{\alpha\beta} = F_{\alpha\beta}^{V} + P_{\alpha\beta} F_{\alpha\beta}^{NB} - 4 \sum_{i>j} P_{ij} U_{ij}^{VDW} (U_{ij}^{V})^3 F_{ij,\alpha\beta}^{V}$$

$$- \sum_{ijk} U_{ik}^{V} U_{kj}^{V} P_{ij} U_{ij}^{VDW} \left( U_{ik}^{V} F_{kj,\alpha\beta}^{V} + U_{kj}^{V} F_{ik,\alpha\beta}^{V} \right), \tag{213}$$

with

$$F_{ij,\alpha\beta}^{V} = -\frac{\partial U_{ij}^{V}}{\partial r_{\alpha\beta}}, \quad F_{\alpha\beta}^{NB} = \frac{\partial U_{\alpha\beta}^{NB}}{\partial r_{\alpha\beta}}. \tag{214}$$

In Eq. (213), $F_{\alpha\beta}^{V}$ represents the valence forces, such as bond stretching, bending, and torsion, and is contributed by the Brenner-type potential. The second term represents the NB forces between two properly screened atoms. The third and fourth terms are due to forces arising from correlations between screened bonds which, in most cases, are negligible. Therefore, if atoms $i$ and $j$ do not form a valence bond and do not also form a bond with the same atom $k$, then usually $F_{ij,\alpha\beta}^{V} = 0$, and either $U_{ik}^{V} = 0$ or $U_{kj}^{V} = 0$, leading to zero contribution. However, even if both atoms $i$ and $j$ bind to a common atom $k$ or if they form a valence bond directly, these terms will still make a negligible contribution to the total energy because of the exponential screening factor, Eq. (212), since in this case either $U_{i\alpha}^{V} < 0$ or $U_{i\beta}^{V} < 0$. Note that there is no restriction on this formulation. That is, the specifics of the NB terms do not alter the general formulation of this potential. Che *et al.* (1999) used this generalized interatomic potentials to study the energetics and structures of a variety of materials, including graphite and molecular crystals and bucky tubes, and Lim *et al.* (2003) utilized it for generating porous amorphous carbon structure and investigating transport, chemisorption and separation of gases in such materials.

### 9.7.5  Other Interaction Potentials

Johnson (1964) developed a potential for representing iron which is given by,

$$U(r_{ij}) = -b_1(r_{ij} - b_2)^3 + b_3 r_{ij} - b_4, \tag{215}$$

where the $b_i$s are parameters of the potential. The interaction force that results from Eq. (215) decays much faster past its maximum than the LJ potential. One advantage of the Johnson potential is that, when used in MD simulation of fracture of a material or metal, it can support an atomically-sharp equilibrium fracture, leaving it stable up to the critical Griffith load (see Chapters 6 and 7). Thus, one can make a meaningful comparison between the MD results and predictions of the continuum mechanics.

Kane (1985) listed several older four-parameter interaction potentials that are related to the Keating model (see Chapter 8 of Volume I) and are intended for describing diamond-structure compounds involving C, Si, Ge, and Sn, and also the zinc-blende-structure compounds GaP, GaAs, and ZnS. These potentials also appear to provide accurate predictions for various properties of these materials.

Kaxiras and Pandey (1988) constructed a potential in order to specifically simulate processes in the bulk diamond lattice. The potential was fitted to the entire energy surface of atomic exchange obtained from an accurate DFT computation. It correctly predicts the static properties of the perfect diamond lattice and reproduces the energy of the concerted exchange path to better than 0.1 eV. However, the energies of bulk point defects in their unrelaxed configuration appear to be too low. Because the potential describes very well a large range of local distortions from the perfect tetrahedral configuration, it can be useful in simulations of materials such as amorphous structures where the coordination remains predominantly fourfold.

Development of empirical or semi-empirical potentials for use in MD simulations remains an active research field. It is neither possible nor necessary to list all the interaction potentials that have been proposed in the past. Several such empirical potentials, in addition to those discussed here, are discussed by Günes *et al.* (2000). In addition, Balamane *et al.* (1992) and Bazant *et al.* (1997) compared the performance of many interatomic potentials.

In addition, we should caution the reader that, despite their success in predicting many properties of a wide variety of materials, none of the potentials discussed so far is without problems. For example, many of these potentials do a relatively poor job of modeling the energetics of small clusters, as well as the various reconstruction of the Si(111) surface. Another example of typical shortcoming of such potentials is provided by the MD simulation of fracture propagation in silicon (Holland and Marder, 1998). In these simulations the three-body part of the Stillinger–Weber potential had to be manipulated in order to obtain physical results, which then of course introduced some unwanted changes in the bulk properties of the material. The same type of difficulties were observed with the Tersoff potentials.

## 9.8   Molecular Dynamics Simulation of Fracture Propagation

In the last section of this chapter we describe application of the MD simulation to fracture dynamics of materials. Our goal is threefold.

(1) We would like to describe how dynamic fracture of materials, the practical consequences of which are at *macroscopic* length scales, is studied by the MD method which considers a system at *atomic* scales.

(2) We aim to demonstrate that, despite the disparity between the length scales at which MD simulations can be carried out, and the practical length scales of interest, not only do the simulation results provide deep insight into the fracture phenomena, but they are in fact in agreement with the experimental observations described in Chapters 6 and 7.

(3) Finally, the discussions of this section are in fact a prelude to what we will be describing in Chapter 10 where we consider *multiscale* modeling of materials.

The fundamental problem facing MD simulation of dynamic fracture (and, more generally, atomistic description of *any* phenomenon in materials) is one of length and time scales. Consider, as an example, fracture of silicon, a material with tremendous technological significance. Its single crystals with a variety of orientations are inexpensive, thus making numerical predictions amenable to experimental verification. Silicon is also very brittle, its crystal structure is well-known and, as discussed in Sections 9.7.2 and 9.7.4, considerable effort has been expended in developing classical interatomic potentials suitable for use in the MD simulation of silicon; we have already described many of such potentials. Therefore, silicon provides an ideal testing ground for MD simulation of dynamic fracture of a material and direct comparison of its predictions with experimental observations. Suppose, for example, that we wish to consider a silicon sample with a length and width that are a few centimeters each and a thickness of only one millimeter. Such a sample contains of the order of $10^{22}$ atoms. The duration of an actual fracture experiment is around 50 $\mu$s, whereas the largest simulations that are currently feasible (see below) allow one to follow what happens to a sample of about $10^8$ atoms for about $10^{-9}$ seconds. Therefore, direct MD simulations of such a sample would require *eighteen* orders of magnitude increase in computer power over what is currently available, a truly daunting, if not impossible, task.

The question therefore is, how does one compare the simulation results with the experimental observations? An appealing approach for doing this would be to merge atomistic simulations with continuum modeling. For example, one may use an atomistic description of the area in the vicinity of the fracture tip, where most if not all the interesting phenomena take place, and combine that with the continuum elasticity (see Chapter 7, and also Chapter 7 of Volume I) for describing the material everywhere else. This approach, which can potentially solve the problem of length scales, but not that of the time scales, will be described in Chapter 10. We describe below what has been accomplished so far without using such a multiscale approach.

Given the enormity of the problem of computer simulation of fracture dynamics of materials at the atomic scale, the goal in the MD simulations should *not* be performing the largest simulation possible, but constructing the *smallest* one capable of providing insightful answers to specific physical questions. In fact, certain features of brittle fracture may profitably be studied by comparatively small simulations, involving only thousands or tens of thousands of atoms. Up until recently, this line of thinking was the dominating factor behind most of the MD simulations of fracture dynamics. On the other hand, there are other aspects of dynamic fracture that, if they are to be studied by MD simulations, one must need a very large number of atoms in the computations.

In any case, before carrying out MD simulations of dynamic fracture of any material, three important points regarding the limitations of such simulations must be kept in mind.

(1) Given the severely non-equilibrium nature of dynamic fracture, it is possible that no classical interatomic potential can provide a realistic description of the material that is investigated by the MD simulations. Moreover, in view

of the emission of electrons and light that is observed in the vicinity of the fracture tip (see Chapters 6 and 7), it is entirely possible that even the DFT would fail as well. Therefore, only a detailed and patient comparison of theory and experiment, which, as of the time of writing this book, has not yet been performed, will be able to settle such doubts.

(2) As discussed in Chapters 6 and 7, in steady state, the energy consumed by the fracture per unit length must equal the energy stored, per unit length to the left. This statement, which remains true even for strains large enough that the applicability of linear elasticity could be called into question, relies on symmetry rather than the matched asymptotics of fracture mechanics, and must therefore be reproduced by any MD simulation.

(3) The MD simulation must contain a complete description of the cohesive zone. However, as more energy is fed to the fracture tip, as temperature rises, or if one studies heterogeneous or ductile materials, the size of the cohesive zone increases, and therefore the size of the system used in the MD simulation must increase accordingly. One should not expect MD simulations to provide "easy" (i.e., without much effort, patience, and efficient computational strategy) predictions for materials where the cohesive zone is of the order of microns, let alone millimeters, as in the silicon sample described above. In this regard, MD simulation of fracture of amorphous polymers on which many experiments have focused (see Chapters 6 and 7) provide a particularly great challenge.

In what follows we first discuss the early MD simulations that involved only a few hundred or thousand of atoms. This will give the reader an idea about the long road that has been travelled in order to arrive at the present state-of-the-art of the MD simulation of dynamic fracture. We then describe and discuss the recent advances and compare the results of large-scale MD simulation of dynamic fracture with the relevant experimental data.

## 9.8.1  Early Simulations

The idea that some sort of a thermodynamic approach (which can be related to MD simulations) may be used for investigating fracture of solids was probably first hinted in a paper by Max Born (1939), who was interested in developing a first-principle criterion for melting. He made the observation that, "the difference between a solid and a liquid is that the solid has elastic resistance against shearing stress, while the liquid has not." He developed a thermodynamic approach to this problem, and also proposed that a generalization of his approach which includes anisotropic stress should be capable of accounting for breaking of crystals. Over 40 years later Born's suggestion was taken up by Nishioka *et al.* (1980,1981) who developed a variational formulation for a solid under uniaxial stress, and wrote down the free energy of the system as a function of the lattice constants parallel and perpendicular to the loading direction. By fitting a Gaussian-type pair potential to the zero-temperature Young's modulus, they calculated the fracture strength as the maximum tension allowing positive free energy curvature.

The statistical thermodynamic approach to fracture of solids was further developed in a series of papers by Blumberg Selinger and co-workers. In an interesting paper, Blumberg Selinger *et al.* (1991a; see also Englman and Jaeger, 1990) developed such an approach in which fracture at failure threshold corresponds to a metastability limit, or spinodal. In their formulation, the role of non-equilibrium defects, such as macroscopic fractures, dislocations, and impurities, in lowering the fracture strength of the material is similar to that of dust particles in lowering the nucleation barrier. Rundle and Klein (1989) developed a similar theory using a field-theoretic approach but, typical of such theories, theirs was a coarse-grained theory without any reference to the structural details of the materials. Their theory was tested by MD simulation of an ideal solid by Wang *et al.* (1991), who showed that the solid remains in metastable equilibrium all the way to the critical stress or force for its fracture, at which point it fails irreversibly by nucleation of small defects. Building on their formulation, Blumberg Selinger *et al.* (1991b) proposed that the onset of fracture in a defect-free material is associated with the loss of a metastable minimum in its free energy at the critical stress.

We must, however, point out that the statistical thermodynamic approach to fracture of a solid is useful only when thermal fluctuations play the dominant role in its failure. This will be the case if the solid is perfectly periodic (without any defect) or it contains very little disorder. As soon as the heterogeneity of a material is even "mild," it begins to play the dominant role in its fracture and failure, and therefore thermal fluctuations will no longer be important. As we have emphasized throughout this book, most real materials are at least to some extent heterogeneous, and therefore an approach to their fracture based solely on statistical thermodynamics is inadequate.

The first time that the word "atomistic" was used in the study of cracks was, to our knowledge, in a paper by Sinclair and Lawn (1972). The first "molecular" simulation of fracture dynamics was probably carried out by Weiner and Pear (1975) who used a square lattice of atoms, inserted a crack in its middle, and solved the equation of motion for the atoms. They assumed that if the distance between two atoms becomes too large, they can be considered as disconnected, an assumption similar to what has been used in the quasi-static and dynamic lattice models of fracture described in Chapter 8, which is also typically used in the MD simulation of dynamic fracture as well. Simulations were performed both at zero and non-zero temperatures. Weiner and Pear found that, except at very high applied stresses, the velocity of the crack reaches a steady subsonic and stress-dependent value, which is in agreement with the prediction of continuum fracture mechanics described in Chapter 7.

The first MD-like simulation of fracture was probably carried out by Ashurst and Hoover (1976). They used a triangular lattice in which the atoms interacted with each other by a truncated Hooke's-law force. The most important finding of this study was that, the velocity of the crack never reaches the Rayleigh wave speed $c_R$, consistent with the troubles that linear continuum mechanics of fracture dynamics already had for explaining the experimental data on the speed of fracture propagation (see Chapter 7).

The first attempt for comparing the predictions of the continuum fracture mechanics with the results of MD simulations was made at the beginning of the 1980s. Thomson *et al.* (1971) had presented evidence for *lattice trapping*, a phenomenon in which a crack neither propagates nor heals, rather it remains stable until external loads somewhat larger than the Griffith threshold (see Chapters 6 and 7) are imposed on the system. The magnitude of the trapping range depends strongly on the characteristics of atomic bonding of materials. Lattice trapping may also depend on the direction in which the crack tip bonds are broken, and may therefore be different for fracture propagation along different crystallographic directions. Moreover, as described in Section 7.12, Rice and Thomson (1974) had developed a criterion for the degree of brittleness of a material according to which a material can be considered as brittle if a dislocation in the neighborhood of the fracture tip cannot escape from the tip region. These predictions were put to test in the first truly MD simulations that were carried out by Paskin *et al.* (1980,1981). We call their computations "true" MD simulations because, unlike Ashurst and Hoover (1976), they utilized the LJ potential for representing the interactions between the atoms in a triangular lattice. In their simulation, a crack was inserted in the middle of the lattice in order to initiate fracture propagation. The smallest crack in a MD simulation is represented by a pair of atoms the bond between which has been cut, so that the two atoms do not directly interact with each other. An external force was then applied to the lattice, and Newton's equations of motion were solved in order to calculate the atomic positions, velocities, and forces. The cutoff $r_c$ for the LJ potential was assumed to be slightly smaller than two lattice bonds at equilibrium. Paskin *et al.* showed that the Griffith energy criterion (see Chapters 6 and 7) is incorrect for large cracks. Their MD simulations also indicated that lattice trapping is a negligible effect, which they attributed to the long range of the interaction potentials. However, in general one should expect lattice trapping to disappear at temperatures much lower than room temperature, and therefore, in order to observe this phenomenon, experiments and MD simulations must be carried out at very low temperatures. The necessity of a low temperature explains why no lattice trapping has yet been observed experimentally in either crystalline or amorphous materials. Paskin *et al.*'s simulations also indicated that the Rice–Thomson criterion for brittleness is valid at low temperatures (see below for more discussion of the Rice–Thomson theory).

In addition to Paskin *et al.*'s simulations, interesting MD computations were carried out by Soules and Busbey (1983) to study fracture of sodium silicate fiber glass. Instead of interatomic forces that result from the LJ potentials, these authors used the following semi-empirical equation for computing the interatomic forces,

$$F_{ij}(r) = \left[ 3.448 a_{ij} \exp(-3.448r) + \frac{e^2 Z_i Z_j}{r^2} \right] \left[ 1 - \left( \frac{r}{r_c} \right)^4 \right], \qquad (216)$$

where $a_{ij}$ is a parameter of the model in erg, $Z_i$ is the charge of atom $i$, $e$ is the electron charge, and $r_c$ is the cutoff distance at which the interatomic forces vanish. The exponential term of this equation represents a repulsive force, while the $1/r^2$

term obviously represents the Coulombic contribution. Simulations of Soules and Busbey indicated that the glass breaks when it is suddenly subjected to a large biaxial expansion. Moreover, when the temperature of the system was raised by about one order of magnitude, the strength of the material decreased by a factor of about 2, a result that was claimed to be in agreement with experimental data.

To our knowledge, Ray and Chakrabarti (1985a,b) and Chakrabarti *et al.* (1986) were the first to carry out MD simulations of fracture of a model of materials with quenched disorder. In their model the disorder was percolation-type, i.e., the heterogeneity was generated by randomly breaking some of the bonds between the atoms before the simulations were commenced. The atoms interacted with each other through a LJ potential for which the cutoff distance $r_c$ was set to 1.6 lattice bonds. During deformation of the lattice a bond was considered broken if the distance between its end atoms was larger than $r_c$. Chakrabarti and co-workers found that the stress needed for fracture vanishes, and the time to complete fracture diverges, both at the bond percolation threshold of the lattice, $p_c^B \simeq 0.347$, whereas the elastic moduli of the lattice vanish at the rigidity percolation threshold, $p_{ce}^B \simeq 0.65$ (see Chapter 8 of Volume I). The latter result was of course an artifact of the model, because the lattice used was effectively a central-force model and, as discussed in Chapter 8, the percolation threshold of central-force lattices is much larger than the connectivity threshold. Other issues of interest in dynamic fracture using small-scale MD simulations were investigated by several research groups, whose results up to 1987 was reviewed by Dienes and Paskin (1987).

Sieradzki *et al.* (1988) utilized MD simulations to study dynamics of crack extension, using a triangular lattice of atoms which interacted with each other through the Johnson potential, Eq. (215). Their simulations indicated that the terminal fracture velocity was about 1/4 of the Rayleigh wave speed $c_R$, and that this terminal velocity depended on the configuration of the fracture tip. This was probably the first time that, in addition to the ample experimental evidence, reasonably-accurate MD simulations had also indicated the breakdown of the linear continuum fracture mechanics. Hoagland *et al.* (1990) employed MD simulations to investigate the configuration of the fracture tip in aluminum, using an embedded-atom potential described in Section 9.7.1. The pair potential used was a Morse function given by Eq. (175), and the embedding energy was determined numerically. As described in Section 9.7.1, the advantage of an embedded-atom potential is that it intrinsically incorporates many-body contributions, and thus dynamic fracture may be studied in a more realistic manner. Hoagland *et al.*'s MD simulations provided evidence for several interesting phenomena, including the existence of *two* singular fields for an atomically sharp fracture. One was an outer field with a strength that was equal to the applied Griffith threshold with its origin at the fracture tip, while the other was behind the fracture with a strength less than the threshold. This was attributed to the nonlinearity arising out of elastic softening of the material near the tip.

Cheung and Yip (1990), employing the embedded-atom potentials, studied the response of a crystal containing a sharp fracture to varying stress and temperature. Over a limited range of temperature, a transition from brittle to ductile fracture

was observed, caused by dislocation emission from the tip of the fracture. This result indicated the existence of an energy barrier for nucleation of the dislocation. Cheung and Yip (1990) showed by detailed calculations that this energy barrier could not be predicted by the continuum theory of Rice and Thomson (1974) mentioned above. This issue was also investigated by Zhou *et al.* (1994) using MD simulations, who proposed that the Rice–Thomson theory should be modified to include the effect of tensile broken-bonds, if it is to correctly predict dislocation emission.

### 9.8.2   *Large Size and Scalable Molecular Dynamics Simulation of Fracture*

As discussed above and also in Chapter 7, in order to investigate certain issues in dynamic fracture by MD simulations, using a large number of atoms is important, because fracture phenomena are sensitive to the sample size, and thus one needs significant computational resources and efficient computational strategies in order to study the size effect in the MD simulations. To this end, to our knowledge, the first MD simulation of dynamic fracture using a massively-parallel computational strategy was carried out by Wagner *et al.* (1992) who utilized $10^6$ particles, by far larger than all the previous MD simulations of dynamic fracture. They used a LJ potential with a spline cutoff together with an analytic embedded-atom potential [see Section 9.7.1 and Eqs. (33)–(36)], and investigated the phenomenon of spallation which occurs at very high strain rates. They demonstrated that an adiabatic expansion can cause spallation, and that the spall strength is proportional to the logarithm of the applied strain rate. However, the LJ material exhibited brittle fracture, whereas the embedded-atom material produced ductile-like fracture, so that the effect of the potentials used in MD simulations of fracture is non-trivial.

Over the past few years, a number of large-scale MD simulations of dynamic fracture have been performed (Abraham *et al.*, 1994; Holian and Ravelo, 1995; Zhou *et al.*, 1996; Omeltchenko *et al.*, 1997; Gumbsch *et al.*, 1997). One of the largest of such simulations was carried out by Zhou *et al.* (1997) who used 35 million atoms to study fracture of a 3D solid. They studied ductile failure and observed dislocation loops emitting from the tip of the fracture. One important result of this study was that, the sequence of dislocation emission events depends strongly on the crystallographic orientation of the fracture front, a result that the previous theories had not predicted. Other MD simulations of dynamic fracture were performed by Zhang and Wang (1996), Machová (1996), Español *et al.* (1996), and Hua *et al.* (1997), investigating various aspects of the problem, such as the effect of boundary conditions, and use of an *N*-body potential, instead of the usual two- or three-body potentials.

The largest MD simulations of dynamic fracture that we are aware of were carried out by Abraham *et al.* (1997b) who studied the response of a 3D notched solid under tension using more than $10^8$ atoms. In their simulations, the interatomic interactions were modeled by the LJ potential, a parallel computational strategy

FIGURE 9.11. Notched solid, a FCC crystal, used in the simulation of fracture. Due to the boundary conditions, the front and back faces are not exterior surfaces, and are thus transparent (after Abraham *et al.*, 1997b; courtesy of Dr. Farid F. Abraham).

based on a spatial-decomposition algorithm described in Section 9.6.2.3 was utilized, and the simulations were carried out at an initial temperature of zero. As such, their system represented a rare-gas solid. The system studied was a slab of atoms with $L_x = L_y = 336$ and $L_z = 896$ atoms for the three orthogonal directions. The notch was a slit beginning midway along $L_x$ for $y = 0$, having a $y$ extension with a length of 120 atoms, extending through the entire thickness $L_z$. The exposed notch faces were in the $y - z$ planes with (110) faces, with the notch pointing in the $\langle 1\bar{1}0 \rangle$ direction. Note that the (110) face does not have the lowest surface energy, thus contradicting the conventional wisdom that would identify the lowest energy surface, i.e., the (111) surface, as the cleavage plane for brittle fracture. The choice of the (110) surface was, however, based on the MD result of Abraham (1996) that indicated that the nonlinear elastic anisotropy of the crystal, and not the anisotropy of surface energy, controls the cleavage behavior. Figure 9.11 shows the system simulated.

Abraham *et al.* (1997b) found that when the speed of fracture propagation approaches one third of the Rayleigh wave speed, the crack tip begins to roughen on *atomic scales*, followed immediately by a dynamic *ductile-to-brittle transition* where plasticity becomes dominant through prolific emission of loop dislocations and the arrest of the crack motion. Figure 9.12 presents magnified off-diagonal views of the cohesive zone during the time period of this transition (times 36, 43, and 54 in the figure from second panel on the left to the right). The atomic roughening is consistent with the onset of the brittle fracture instability suggested by the experiments discussed in Chapters 6 and 7, and hence supporting the notion that this instability is a general feature of rapid fracture of brittle materials. However, although Kelly *et al.* (1967) and Rice and Thomson (1974) had proposed that rare-gas solids are inherently ductile, Abraham *et al.* (1997b) found that their system undergoes brittle fractures along the (110) plane, but fails by ductile plasticity for a notch with (111) or (100) faces. Thus, not only their MD simulations indicated partial break down of the Rice–Thomson and Kelly *et al.* theory, but also provided new insights into fracture behavior of materials.

FIGURE 9.12. Off-diagonal views of the cohesive zone during the time of the brittle-to-ductile transition. Darker and lighter areas indicate, respectively, whether the speed of the atoms is less than or greater than one twelfth of the longitudinal sound speed (after Abraham *et al.*, 1997b; courtesy of Dr. Farid F. Abraham).

## 9.8.3    Comparison with the Experimental Observations

With the enormous increase in the computational power, MD simulations of dynamic fracture have reached that degree of accuracy that their predictions can be directly compared with experimental observations and data. Three important aspects of dynamic fracture that have been reproduced by MD simulations are as follows.

### 9.8.3.1    Fracture Instabilities

One of the first comparisons between the results of MD simulation of dynamic fracture and the experimental observations was taken by Abraham *et al.* (1994) in which the dynamics of a fracture was studied using a $10^6$ atom crystal. The interaction potential between the atoms was a LJ potential. Stress was applied to the system by displacing its opposing boundaries at a strain rate that was approximately equivalent to that obtained in explosive loading applied to fracture faces. These conditions were necessary to achieve sufficient acceleration of the fracture tip so as to achieve high enough crack velocities over the duration of the simulation, and to be able to detect the existence of instabilities in the fracture's motion. The results of these MD simulations were close to the experimental observations in amorphous materials. The fracture was observed to accelerate smoothly until it reached a velocity of $0.32c_R$. The velocity of the fracture experienced large fluctuations when its value exceeded $0.32c_R$, as a result of which the instantaneous velocity of the fracture tip became erratic. These fluctuations were coupled with a "zig-zag" motion of the fracture tip, which formed in its wake a rough fracture surface. These interesting simulations, which are in complete agreement with the experimental

observations described in Chapters 6 and 7, highlighted the robust and general nature of the crack instability.

The robust nature of the crack instability was further highlighted by the work of Zhou *et al.* (1996; see also Gumbsch *et al.*, 1996) in which fracture propagation was investigated in a 400, 000 atom crystal. The atoms interacted with each other via a Morse potential, Eq. (175). By varying the applied strain rates the maximum velocity that a fracture could achieve over the duration of the simulation was varied between $0.18c_R$ and $0.36c_R$, where the strain rates used corresponded, as in the work of Abraham *et al.*, to explosive loading of the system. At a velocity of $0.36c_R$, instability of the crack was observed to occur by its forming several branches. The branching process was observed to be immediately preceded by the nucleation of a dislocation in the crystal together with a build up of the phonon field in the vicinity of the crack tip. These features are again in agreement with the experimental observations described in Chapter 7. In both of the above large-scale MD simulations, the entire fracture process (from initiation to the onset of the instability) occurred over a time of about 1 nanosecond. For this reason, both the strain rates used and the amount of strain in the material at the onset of fracture (approximately an order of magnitude larger than observed in real materials) had to be extremely large. Thus, the close correspondence of the results of the MD simulations with the laboratory results obtained in amorphous materials is rather surprising. The short time scales used in these simulations preclude, of course, examination of steady-state properties of the system. It should therefore be interesting to compare the results obtained in these experiments with steady-state results that can be obtained by smaller MD simulations.

Further insight into the limiting velocity of the crack tip was provided by the MD simulations of Abraham (1996), whose simulations indicated that this velocity can reach 60% of the Rayleigh wave speed by following the highest energy path. Subsequent MD simulations of dynamic fracture in graphite sheets by Omeltchenko *et al.* (1997) confirmed Abraham's results. They found that, for certain crystalline orientations, multiple fracture branches with nearly equal spacing are created as the velocity of the fracture tip reaches $0.6c_R$.

Abraham *et al.* (1997a) carried out MD simulations of dynamic fracture of a 2D material represented by a triangular lattice with more than 2 million atoms. The interatomic forces were treated as central forces, modeled by a LJ spline potential [see Eqs. (33)–(36)] and the analytic embedded-atom model that represents a many-body potential, developed by Holian *et al.* (1991) and described in Section 9.7.1. As discussed by Holian *et al.* (1991), the LJ potential can accurately represent brittle materials, while the embedded-atom model can be used for studying ductile solids. Under these conditions, the reduced melting temperature is $k_B T/\varepsilon = 0.2$ for the EAM and $k_B T/\varepsilon = 0.4$ for the LJ material. Moreover, the longitudinal wave speed $c_l$ [see Eqs. (7.24) and (7.25)] is about $5\sqrt{\varepsilon/m}$ for the EAM material and $9\sqrt{\varepsilon/m}$ for the LJ material, where $m$ is the atomic mass. The transverse speed is $c_t = c_l/\sqrt{3}$, and the Rayleigh wave speed $c_R \simeq c_t$. The simulations were performed at a reduced temperature of $10^{-5}$. Abraham *et al.* (1997) found that, for rapid fracture of brittle (LJ) material, a dynamic instability of fracture growth

FIGURE 9.13. The onset of fracture instability (left), in reduced time intervals of 7 and beginning at 85. The right panels show the fracture zipzag, beginning at reduced time 220 (after Abraham *et al.*, 1997a; courtesy of Dr. Farid F. Abraham).

develops when the crack velocity approaches 1/3 of the Rayleigh wave speed $c_R$. At higher crack velocities, the fracture either follows a wavy path or branches out, with the anisotropy that is due to the large deformation at the tip of the propagating fracture playing the dominant role in determining the path of the fracture. This is of course in contrast with the conventional wisdom (see Chapter 7) which associates the lowest energy surface as the favored cleavage direction. Figure 9.13 nicely demonstrates these results. The simulations also produced dislocations emission from the rapidly moving fracture tip *after* the onset of the crack growth roughening, implying that the dislocations are the *consequence* rather than the *cause* of the dynamic instability. The number of the dislocations emitted was dependent upon the external loading. These results are all in excellent agreement with the experimental observations described in Chapters 6 and 7.

Let us point out that fractures at zero temperature display a clear dynamic instability at a critical energy flux. Up to this point, phonons are able to carry away all excess energy. This instability does not take the form of a simple microbranching instability, partly because atomic bonds can easily rejoin above the main crack line in a single component solid with no environmental impurities available. Instabilities at room temperature have not yet been explored either numerically or analytically.

### 9.8.3.2  Morphology of Fracture Surface

The morphology of fracture surfaces has also been studied by MD simulations. For example, Nakano *et al.* (1994) performed MD simulations of fracture surfaces in porous silica glasses. The surface of the pores was rough, and thus particular attention was paid to root-mean square fluctuations in the height $h$ of surface of the pores, averaged over a length scale $\ell$. The MD results indicated that $h$ scales with $\ell$ as in Eq. (7.130) where, as discussed in Chapters 1, 6 and 7 $\alpha$ is the roughness exponent. The simulations of Nakano *et al.* (1994) yielded $\alpha \simeq 0.87$, in agreement with the prediction of theory of rough surfaces and the experimental data for many fracture surfaces during rapid fracture (see Chapters 6 and 7), thus supporting the claim that in this regime $\alpha$ is universal. Subsequently, Nakano *et al.* (1995) (see also Tsuruta *et al.*, 1996; Omeltchenko *et al.*, 1996) studied, using MD simulation, fracture dynamics in amorphous $Si_3N_4$ films. They showed that the surface rough-

ness exponent $\alpha$ depends on the speed of fracture propagation. At the initial stages of fracture, when the crack tip propagated slowly, $\alpha \simeq 0.44$. However, once the speed of fracture propagation exceeded a certain limit, a crossover was observed to a higher value, $\alpha \simeq 0.8$. These results are in agreement with the experiments of Bouchaud and Navéos (1995) which were described in Chapters 6 and 7. In another effort by this group, Li *et al.* (1996) studied dynamic fracture in $SiSe_2$ nanowires, and found that fracture is initiated in an amorphous region of the surface of the material, while multiple fractures start at the boundaries of the amorphous region. Finally, Kalia *et al.* (1997) investigated dynamic fracture in nanophase $Si_3N_4$, showing that intercluster regions of the material are amorphous, deflecting fracture and hence giving rise to local crack branching. This implies that nanophase $Si_3N_4$ can resist fracture much better than crystalline $Si_3N_4$. The roughness exponent was found to be, $\alpha \simeq 0.84$, in agreement with the experimental data.

### 9.8.3.3    Fracture Propagation Faster Than the Rayleigh Wave Speed

As discussed in Chapter 7, it is generally believed that a brittle crack cannot propagate faster than the Rayleigh wave speed $c_R$. Continuum mechanics predicts that for Mode I tensile loading the forbidden velocity zone for fracture propagation is any speed larger than $c_R$, while for Mode II shear loading the forbidden zone exists only for speeds between $c_R$ and the shear wave speed. However, the limiting speed of a Mode II fracture is also $c_R$ because its forbidden zone acts as an impenetrable barrier for the shear fracture and does not permit it to propagate faster than $c_R$. On the other hand, experimental evidence that was discussed in Section 7.8.15 indicates that, under certain circumstances, a Mode II shear fracture can propagate faster than the shear wave speed.

Abraham and Gao (2000) carried out MD simulation of fracture in a 2D rectangular slab with $L_x = 1424$ and $L_y = 712$ atoms, thus utilizing over $10^6$ atoms. The system used consisted of two crystals joined by a weak interface. The interatomic forces were assumed to be harmonic, excepts for those pairs of atoms with a separation cutting the horizontal center line of the simulation slab, for which a LJ potential was utilized. The cutoff distance was taken to be $2.5\sigma$, where $\sigma$ is the size parameter of the LJ potential. As mentioned above, at zero temperature and pressure, the longitudinal wave speed $c_l$, in reduced units, is 9, the shear (transverse) wave speed is $c_t = c_l/\sqrt{3} \simeq 5.2$, and the Rayleigh wave speed $c_R$ is 4.83, where the LJ parameters $\sigma$ and $\varepsilon$ are used as the basic units of length and energy (thus, for example, $c_l = 9\sqrt{\varepsilon/m}$ in dimensional units). A horizontal crack of 200 atom distance was cut midway along the left-hand vertical slab boundary. The 2D crystal used was a triangular lattice with the initial crack being parallel to the close packed direction. The initial temperature of the system was zero, and the MD simulations were carried out at constant energy (see Section 9.2.4). Both Modes I and II were simulated.

Molecular dynamics simulations of Abraham and Gao (2000) showed that, in Mode I (cracks under tensile loading), the crack quickly approaches a constant velocity of about 4.83, the same as the Rayleigh wave speed, after which it propagated with the same constant speed. In contrast, in Mode II (cracks under shear

loading) the crack tip first reaches the Rayleigh wave speed $c_R$ and for some time propagates with this speed, and then jumps to a higher constant speed with a value of about 8.97, essentially the same as the longitudinal wave speed $c_l$. Similar results were obtained (Gumbsch and Gao, 1999) by MD simulations of propagation of dislocation of indentation. These results are in complete agreement with the experimental observations described in Section 7.8.15.

Let us point out that, while effective potentials, such as the embedded-atom and LJ potentials (with fitted parameters), may be adequate for representing metals, they are poor representatives of non-metallic materials. In this case, one must use a first-principle quantum mechanical description of the materials in order to calculate the potentials (see Sections 9.4 and 9.5). This is a rigorous method since, unlike the case of MD simulations with empirical or semi-empirical potentials, *ab initio* quantum mechanical computations do not use any adjustable parameters. Moreover, as discussed in Sections 9.2 and 9.4, the local density approximation in its plane-wave pseudo-potential formulation can be optimized, so that the computations will be highly efficient. Kaxiras and Duesbery (1993) presented the results of such a study for silicon, and Spence *et al.* (1993) used *ab initio* QMD simulations to investigate the dependence of lattice trapping energies on applied load for fractures propagating in silicon. Pérez and Gumbsch (2000) studied the anisotropy of cleavage fracture in silicon and the effect of sample size on the process. This approach has not, however, been utilized extensively, presumably because its computational cost is much larger than MD simulations using the empirical or semi-empirical potentials discussed above.

## Summary

*Ab initio* quantum mechanical computations based on the density functional theory in the local density approximation, together with plane-wave pseudo-potential formulation, offer an efficient and rigorous method for computing materials properties. Quantum MD simulation method of Car and Parrinello did not change the essentials of such computations, but offered an enormous increase in the efficiency of the method, hence making much larger pieces of materials accessible to such computations.

As a technique for studying materials at the atomic scale, molecular dynamics simulation has been used for several decades. However, development of vector computers and parallel machines, and hence vectorized and parallelized computational algorithms, together with derivation of accurate interatomic potentials, have made MD simulations a powerful tool for studying materials at atomic scales utilizing millions of atoms and molecules.

These two computational strategies have enabled us to investigate and accurately predict various properties of materials. We believe that when one or both of these methods are joined with the multiscale approach that will be described in the next chapter, the possibilities for accurate and efficient *optimal design of materials with specific properties may be limitless*.

# 10
# Multiscale Modeling of Materials: Joining Atomistic Models with Continuum Mechanics

## 10.0 Introduction

Throughout this book we have emphasized that solid materials of industrial importance are highly heterogeneous, with the heterogeneities manifesting themselves at several length scales, ranging from the smallest to macroscopic scales. For many centuries, such materials were discovered, mined, and processed in a serendipitous manner. However, characterization of atoms and the progress made in x-ray diffraction during the early decades of the twentieth century provided the impetus for the search for a theory of materials that can explain their properties in terms of their atomic constituents. It was soon realized that developing such a theory was not practical yet, because

(1) the disparity between the relevant length scales, i.e., those at the atomic and macroscopic scales is huge;
(2) the existing knowledge of the principles of atomic cohesion and the basic properties of materials was totally inadequate, and
(3) the required computational power for solving the problem was (and still is) enormous, while the available computational power was grossly limited. Therefore, the consensus at the time was that any theory of materials could not have predictive power.

Later decades of the twentieth century witnessed development of many qualitative and semi-quantitative models of materials that could explain the principles of atomic cohesion, and basic properties of such fundamental materials as metals and semiconductors. Some of these models, though relatively simple, were surprisingly accurate and helped us make remarkable progress in understanding materials' properties. However, as discussed in Chapter 9, for most materials of industrial importance, the interatomic interactions are complex enough to require very elaborate models. Such models must usually be accompanied by very extensive computations.

In Volume I, as well as in the present Volume, we utilized continuum mechanics and lattice models to describe modeling and simulation of the morphology,

and estimating the effective transport and mechanical properties of heterogeneous materials. These models and approaches are appropriate for describing a material at the microscopic and macroscopic length scales, but cannot provide any insight into its properties at the smallest length scales, namely, the molecular scales. In Chapter 9 we described modeling and simulation of materials and their properties at the atomic and molecular scales. As discussed there, methods for computing the properties of materials at such length scales are divided into two major classes: Those that do not use any empirically- or experimentally-derived quantities, and those that do. The former are the *ab initio* (or first-principles) methods, while the latter are referred to as empirical or semi-empirical techniques. Several important theoretical and computational tools of both classes, developed over the past three decades, such as the density functional theory (DFT) and its variants, the classical molecular dynamics (MD) simulation, the quantum MD (QMD) technique, and the tight-binding (TB) methods, were all described in Chapter 9. The advent of very fast computers and development of efficient strategies for massively-parallel computations, that were also described in Chapter 9, have made it feasible to carry out large-scale computations at the atomic and molecular length scales, making such techniques indispensable tools for investigating and predicting materials' properties on such length scales.

However, we should not "abuse" atomic-scale simulations in the study of macroscopic properties of materials, since atomic-scale mechanisms are in general separated from the macroscopic behavior that they engender by a vast array of intervening *continuum* scales. These *mesoscopic* length scales both filter (i.e., average) and modulate (i.e., set the boundary conditions or driving forces for) atomic-scale phenomena, and are therefore an essential part of the constitution of materials.

As our discussions throughout Volume I and the present Volume should have made it clear by now, continuum models are based on the assumption that the relevant fields that describe the state of a material vary slowly on the atomic scale; otherwise, continuum models that represent macroscopic behavior of a material, and are derived by averaging the material's properties at the smaller length scales, lose their meaning. Therefore, if we are, for example, to describe properties of a material with defects, continuum theories *a fortiori* break down in the vicinity of the defects, or, more generally, any other entity that possesses structure on the atomic scales. It is therefore clear that continuum theories can be "enriched" by incorporation of additional information, and hence avert their breakdown. Thus, *atomistic and continuum models need and reinforce each other*.

Given great advances in both continuum formulation and atomistic simulation of the behavior of materials, we have reached the stage in which we may be able to integrate models for describing materials at the atomic and molecular scales with those at the continuum level, thus developing a methodology for making quantitative predictions for materials' properties that utilizes our knowledge at *all* the relevant length scales. To better motivate this discussion, consider an important phenomenon, namely, fracture of silicon, that was already described in Chapter 9. Suppose that one is to carry out a fracture experiment in samples of silicon

with lengths and widths that are several centimeters each and a thickness of about a millimeter. Such samples would contain about $10^{22}$ atoms. The duration of the experiment is about 50 $\mu$s whereas, as discussed in Chapter 9, the largest atomistic simulations that are currently feasible can follow $10^8$ atoms for around $10^{-9}$ s. Direct atomistic simulation of fracture of such a sample of silicon will therefore require more than *eighteen* orders of magnitude increase in computational power over what we currently have. Such a breakthrough in computational power will not be achieved any time soon, and therefore the critical question is: How can one make a comparison between the results of the atomistic computations and the relevant experimental data, or even the predictions of the continuum theories?

The essential problem facing such simulations is one of length and time scales. However, before answering the above question, we should first remind ourselves that the goal in computer simulations, both at the continuum and atomistic scales, should *not* be performing the largest simulation utilizing the largest possible system, but constructing the smallest one that is capable of answering specific physical questions. In the example of brittle fracture of silicon, as well as in many other phenomena that occur in materials, many important features may profitably be studied by atomistic simulations that are comparatively small, involving only tens of thousands of atoms. Therefore, a fruitful strategy may be based on merging atomistic and continuum simulations. Hence, in the example of fracture of silicon, one may describe the vicinity of the crack tip by atoms, but utilize the continuum elasticity everywhere else.

The combined methodology is called a *multiscale approach* to modeling of materials. In political jargon, multiscale modeling of materials is a *divide-and-conquer* modeling paradigm. As the first step, the entire range of material behavior is divided into a hierarchy of length scales. Next, the relevant physical processes that are irreducible and operate independently at a given scale—sometimes referred to as unit processes—are identified. The unit processes at one scale represent averages of unit processes operating at the immediately lower length scale, and this relation defines a partial ordering of the processes.

Over the past decade, multiscale modeling of materials has evolved from something that was thought of as a distant dream to a research area with intensive methodological development. If its development continues at the current rapid pace, it will, in the relatively near future, develop into a practical tool for industrial applications involving design of materials that possess specific properties. The goal of this chapter is to describe the basic ideas and techniques for multiscale modeling of materials. Because the range of problems and phenomena to which multiscale methods are applicable is very broad, we restrict the discussions to the main theme of this book, namely, modeling of materials and predicting their properties, and describe the essentials of the multiscale modeling approach by discussing a few recent applications. We begin our discussion by briefly describing two main classes of multiscale modeling approach developed so far, after which we describe some of their recent applications. More extensive discussions and overviews pertaining to micromechanics and multiscale modeling of materials are given by, for example, Ortiz and Phillips (1999) and Phillips (2001).

# 10.1    Multiscale Modeling

We should first point out that the multiscale paradigm is more easily stated than carried out in practice. Currently (at the time of writing this book), the analysis of mechanisms and the characterization of the effective properties of materials rely on either extensive numerical computations, or an assortment of analytical tools, ranging from mean-field and effective-medium approximations to variational approaches that have been described throughout this book. Because of this difficulty, multiscale modeling in general, and combined atomistic-continuum modeling in particular, cannot be easily reduced to a self-contained and unified theory. Therefore, at this stage multiscale modeling is an *art* as well as science. Keeping in mind this fact, we now begin to describe the main concepts and ideas of the multiscale approach.

Generally speaking, there are two types of multiscale approach to modeling of various physical phenomena that occur in materials and systems that contain several disparate length as well as time scales. These approaches are either sequential or parallel methods, in the sense described below. What follows is a brief discussion of each approach.

## 10.1.1    Sequential Multiscale Approach: Atomistically-Informed Continuum Models

In this method, which has been developed and used much more extensively than the parallel multiscale modeling (to be described below), beginning with the smallest length/time scales of the problem, the results of one series of computations are used as the input to the next (larger) up the length and/or time scale hierarchy. Hence, the essential idea is to pass (as input) information (output) from finer to coarser scales.

A good example is provided by polycrystalline plasticity, for which the main identifiable length scales are,

(1)  the nanoscale in which unit processes represent the possible behaviors of single-crystal defects, such as individual dislocations or vacancies;
(2)  the mesoscale which is characterized by the collective behavior of large numbers of defects, as in, for example, dislocation dynamics;
(3)  the subgrain scale, characterized by the formation and evolution of subgrain dislocation structures, and
(4)  the polycrystalline scale which is characterized by the collective behavior of a large number of grains.

The atomistic scale and the continua communicate at the nanoscale through *ab initio* computations of material parameters pertaining to continuum theories. Under this scenario, the mesoscopic model sets the functional form of the response functions of a material, while the atomistic models dictate the relevant material-specific parameters of the mesoscopic theory. In effect, one works with continuum models, except that the relevant parameters of the models are provided by atomistic computations, hence the name *atomistically-informed continuum models*.

A beautiful example of implementation of this strategy is the pioneering work of Clementi (1988). He used accurate quantum-mechanical computations, of the type that were described in Chapter 9, in order to evaluate the interaction of several water molecules, from which he developed an accurate empirical interatomic potential that involved two-, three-, and four-body interactions. The potential was then utilized in a MD simulation for evaluating the viscosity of water. Subsequently, motion of water in a channel with or without obstacles was studied, using as input the viscosity that had been computed in the previous step. The resulting understanding was then employed in a fluid dynamics computation for predicting tidal circulations (which is somewhat similar to the classical Bérnard problem; see, for example, Koschmieder, 1993) in Buzzard Bay, Massachusetts.

Another example of sequential multiscale modeling is provided by work in the atmospheric and environmental science (see, for example, Elbern, 1997). In this case, sophisticated computations were used to evaluate reaction barriers of as many as 150 chemical reactions involving over 60 reactants. The results were then used in rate equations for large scales that were coupled to computer codes for spatial grid generation in order to determine and predict chemical meteorology. Elbern (1997) carried out such computations for a domain that covered an area from the eastern North Atlantic to the Black Sea, and from northern Africa to central Scandinavia. The domain was partitioned into many subdomains, each of which was assigned to a processor. Each processor communicated with only its neighboring processors that used its results (output) as the input for their own computations.

A good example of sequential multiscale modeling of materials is provided by the work of Zepeda–Ruiz *et al.* (1999). Their goal was designing experimental protocols toward the development of engineering strategies for strain relaxation of semiconductor films which are grown heteroepitaxially on semiconductor substrates. The strain is caused by the lattice mismatch between the film and substrate, and controlling it, which generates defects in the material, is the key to optimal design of the film's optoelectronic properties. One important design parameter is the thickness of a *compliant* substrate. The thickness and elastic properties of such substrates are comparable to those of the epitaxial film. They behave as if they are unconstrained at their bases, and thus they can aid the film in accommodating the lattice mismatch by either contracting or expanding parallel to the interface between the film and the substrate. The compositional grading of the epitaxial film is another important design parameter.

Matthews and Blakeslee (1974) and Freund and Nix (1996) have already developed a continuum elasticity theory for describing strain relaxation mechanisms in semiconductor heteroepitaxy, and predicting the critical epitaxial film thickness that marks the onset of misfit dislocation formation in heteroepitaxy on both thick and thin substates. Their continuum elasticity theory successfully predicts both the energetics and kinetics of some relaxation phenomena associated with the formation of strained islands grown on semiconductor surfaces. However, the theory does have its limitations. For example, the continuum elasticity theory in conjunction with equilibrium bulk values of the elastic moduli are used commonly for quantitative analyses. However, it is well-known that the elastic stiffness coefficients depend strongly on the stress, which may attain very large values in an

epitaxial thin film. Atomistic simulations can overcome such limitations. Using either an *ab initio* or semi-empirical description of the interatomic interactions between the material's atoms, energy minimization or MD simulations, of the type that were described in Chapter 9, can be utilized for investigating the structure, energetics, and dynamics of a system consisting of an epitaxial film on a substrate.

Using a multiscale approach, Zepeda–Ruiz *et al.* (1999) studied the energetics, strain fields, and semi-coherent interface structures in a layer-by-layer semiconductor heteroepitaxy, such as InAs on GaAs(110) and InAs on GaAs(111)A, as well as the interfacial stability with respect to misfit dislocation formation and the morphology of the surface of the film grown on the substrate. The continuum theory provided a parameterization scheme for the atomistic simulations. A Keating-type potential (see Chapter 9), which contains the contributions of the stretching and bond-bending forces, was utilized for representing the interatomic interactions, and total energy minimization was used for determining the most stable configuration of the system (i.e., the one with the lowest-energy state). The minimization was done based on the conjugate-gradient method (see Section 9.5.2) with respect to the atomic coordinates for a given strain state, which were uniaxial, biaxial, or fully relaxed. Because of misfit dislocations a supercell (see Section 9.1.3) was used that, depending on the state of the system, contained anywhere from 500 to 140,000 atoms. A major conclusion of this work was that the continuum theory of elasticity can be accurate all the way down to the monolayer thickness, which is the finest possible length scale for the theory in the context of layer-by-layer expitaxial growth. In addition, the theory was shown, in conjunction with the atomistic simulations, to provide quantitative predictions for various properties of interest. Indeed, even the linear isotropic elasticity was capable of fitting the results of the atomistic simulations.

## 10.1.2   Parallel Multiscale Approach

This method is not as well-developed as the sequential approach, because it requires very significant computing power which, up until very recently, was not available. In this method, different computational methods, ranging from those for atomic scales to continuum scales, are coupled for a *simultaneous attack* on a given problem. The reason for this coupling is that many physical phenomena are in fact inherently multiscale; that is, one must know what is happening *simultaneously* in different regions and scales of the system in order to understand and predict its macroscopic behavior. A good example is fracture propagation in solid materials which was described in detail in Chapters 6–9. The atoms that constitute the material interact with each other, and nucleation of a crack and its propagation are due to what happens at this length scale, namely, breaking of the bonds between the atoms. As the nucleated crack starts to propagate, complex phenomena, such as plastic deformation, happen at a larger length scale which includes the tip of the propagating fracture. At still larger length scales, which include the region far from the tip of the crack, the material behaves as a continuum which may be described by the classical continuum mechanics. The most important property that

such a multiscale approach must possess, in addition to being efficient, is *accuracy*. However, accuracy in this context is somewhat different from the traditional way that one understands this word, since accuracy in the present context means that the dynamics of a phenomenon under study must be indistinguishable whether they are determined from a multiscale approach or from a system that has the same size but is studied only by a quantum-mechanical approach, i.e., at the most fundamental length scale.

It should therefore be clear to the reader that the multiscale approach can, in principle, solve the problem of how to deal with a heterogeneous material with several distinct length scales. However, in order to be able to make quantitative predictions for materials' properties using the multiscale approach, one must also solve the problem of the disparate time scales, namely, the wide gap between the time scale over which the actual physical phenomena occur, and those that can, with the current computational power, be accessed. This disparity between the two time scales poses severe restrictions to predictive material modeling, regardless of whether one uses sequential or parallel multiscale modeling. For example, in parallel multiscale modeling, the time scale is determined by the finest atomic-scale method that is used in the model, e.g., the QMD or a TB formulation, and this is typically up to one nanosecond. Such time-scale limitations can be partially overcome by using a dynamic Monte Carlo method, or by an accelerated MD technique. In this book, we consider only the question of multiscale modeling of properties of materials that contain several relevant and disparate length scales. modeling of phenomena in which there are several relevant and widely disparate time scales has been discussed by, for example, Voter (1997).

We now begin describing and discussing a few important examples of use of a multiscale approach to modeling physical phenomena in materials.

## 10.2   Defects in Solids: Joining Finite-Element and Atomistic Computations

The first example that we describe is the pioneering work of Tadmor *et al.* (1996a) who studied how a defect, such as a dislocation, crack, or grain boundary, can be embedded within a continuum, but *without* the standard assumptions that are inherent in a simulator based on a continuum model which usually utilizes *ad-hoc* assumptions about failure or fracture of a given region of a material. Their original study was carried out before the current computational power with parallel machines had reached its current state, and as such it had certain limitations. In subsequent papers (see, for example, Shenoy *et al.*, 1998, 1999, 2000; Tadmor *et al.*, 1999), the methodology was refined and extended to various applications. What follows is a summary of this important work.

As pointed out by Tadmor *et al.* (1996a), the analysis of the structure of crystal defects requires consideration of anharmonic effects on the scale of the lattice. Such effects can, of course, be studied by atomistic simulations, such as *ab initio* or MD computations (see, for example, Arias and Joannopoulos, 1994). However,

such computations alone cannot do justice to the problem, nor can they give us a quantitative picture of what is happening during the deformation of a crystal, if we consider the fact that macroscopic deformation of crystals involves dislocation densities as high as $10^{13}$ m$^{-2}$, so that an area as small as 1 mm$^2$ near the tip of a crack can be crossed by over one hundred dislocations. In contrast, the supercell that was used in the *ab initio* computations of Arias and Joannopoulos (1994) contained only 324 atoms. Thus, it would be impossible to consider individual dislocations at the *macroscopic scale*. Even the intermediate length scale, of the order of a few hundred nanometres, is presently hardly within reach of conventional atomistic simulation, as the number of atoms involved in such simulations would be in excess of $10^8$, which makes purely atomistic simulations difficult to carry out. One may develop phenomenological continuum models in which the defects are treated as continuously distributed objects. However, such an approach could not provide deep insight into the role of the defects in the properties of materials.

Tadmore *et al.* (1996a) were interested in studying the deformation processes of interest at the intermediate length scale which, on one hand, involve discrete dislocations in numbers that are too small to be described adequately by a conventional continuum model of crystal plasticity, and, on the other hand, contain too many atoms to be treatable by purely atomistic simulations. Thus, a multiscale approach that joins the two methods, and would seem to be the only practical solution, should have several key attributes, some of which are as follows.

(1) The theory should, at macroscopic length scales, reduce to the continuum crystal elasticity and reproduce its important properties, such as material frame indifference and crystal symmetry (see, for example, Milstein, 1982).

   To achieve this goal, the atoms should be constrained to move in accordance with the continuum displacement field. This would enable one to compute energies and forces from local lattice calculations. Thus, by construction, the resulting continuum would automatically satisfy material frame indifference and exhibit all the symmetries of the crystal. The system would also possess lattice invariance, i.e., its energy would be invariant with respect to distortions of the reference configuration. In particular, the energy density would be periodic under crystallographic slip. An important consequence of this periodicity is the lack of quasi-convexity of the energy functional [see, for example, Fonseca (1988); see also Chapters 2 and 4, as well as Chapters 4, 7, and 10 of Volume I], which would make stable development of lattice defects possible. The relaxation of functionals lacking quasi-convexity would require consideration of minimizing sequences of deformations which exhibit structure on increasingly finer scales, hence necessitating multiscale analysis for simultaneous resolution of macroscopic and microscopic features into the model.

(2) At the atomic scale, the theory should be built upon reliable interatomic interactions, incorporate a lattice constant (or lattice parameter, as referred to by Tadmor *et al.*), and possess all the invariance properties that are expected of a crystal lattice.

Introduction of the lattice parameter is for preventing the displacement field from developing unphysical sublattice-scale structure. Since the continuum equations are solved by a finite-element (FE) method, the lattice parameter can be set so as to impose a lower bound for the FE mesh size. In this manner, the crystal takes on a character similar to the concept of quasi-continua developed by Kunin (1982). The elements in the model are either local or non-local depending on their size, extent of deformation in their region, and energy. The latter type of elements are those that are small, highly deformed and energetic. The domains of the interaction of the two types of element are also different. A local element interacts with deformation only within its own geometrical domain, whereas non-local elements interact also with their neighboring elements. At large scales, the crystal becomes indistinguishable from a nonlinear elastic crystal whereas, in the fine-mesh limit, the theory reverts to lattice statics. This implies, from an atomistic perspective, that the quasi-continua are simply atomistic lattices that are subjected to kinematic constraints, namely, those introduced by the FE interpolation. These constraints eliminate excess atomistic degrees of freedom in regions where the deformation field varies slowly on the scale of the lattices.

(3) In between, at intermediate length scales, the theory should produce a continuous transition from the lattice to the continuum realms.

Note that the distinction between the continuum formulation and atomistic description is not only important from a computational view point, but also from the physical point of view. As discussed throughout this book, at the macroscopic (continuum) level, solid materials are represented as continuous media to which appropriate average material properties are assigned. On the other hand, in atomistic description of a material, a solid is treated as a collection of atoms with interactions that are described by an appropriate energy or potential function. This concept was also the foundation of the lattice models of linear elastic properties of materials that were described and discussed in Chapters 8 and 9 of Volume I, except that in those models the lattice was representative of a coarser piece of materials, and therefore its sites did not represent atoms.

(4) Finally, the theory should also enable one to obtain an accurate treatment of lattice defects, such as dislocations, in case these defects exist or nucleate.

## 10.2.1   The Quasi-continuum Formulation

We first provide a qualitative description of the general ideas and concepts of the approach, after which the details of the model are described. The theory begins from an underlying conventional atomistic model, which is capable of delivering the energy of the crystal as a function of the atomic positions. The configuration of the crystal is then reduced by identifying a subset of representative atoms, which henceforth become the sole independent degrees of freedom of the crystal. The position of the remaining atoms are obtained by piecewise linear interpolation of the representative coordinates, very much similar to the manner by which displace-

ment fields are constructed in the FE method. The selection of the representative atoms may be based on the local variation of the deformation field. For example, one may adapt the mesh in such a way that the variation of the displacement field over each element of the triangulation does not exceed a fraction of the Burgers vectors, hence ensuring that full atomistic resolution is achieved, for example, near dislocation cores and on planes undergoing crystallographic slip. By contrast, far away from any highly-stressed region, the density of the representative atoms decreases rapidly, and the collective motion of a very large number of atoms is governed, without significant loss of accuracy, by a small number of degrees of freedom. In these coarse regions, the behavior of the model is ostensibly indistinguishable from that of a continuum. The effective equilibrium equations are then obtained by minimizing the potential energy of the crystal over the reduced configuration space. Therefore, the number of equilibrium equations that are obtained is commensurate with the number of representative atoms. However, a direct calculation of the effective force field requires, in principle, the evaluation of sums that are extended over the full collection of atoms. Such full sums may be avoided by the introduction of approximate summation rules, whereupon the complexity of the computation of the effective force field becomes of the order of the reduced model.

We now describe the details of the model. In the quasi-continuum (QC) theory of Tadmor *et al.* (1996a) the continuum framework and continuum particle concept are retained, but the macroscopic constitutive law is replaced by one based upon direct atomistic calculations. The continuum particle is represented by a small crystallite of radius $R_c$ which surrounds a representative atom. This crystallite is deformed according to the local continuum displacement field, and its energy is computed based on an appropriate atomistic model. In order to compute the energies, one must make a correspondence between the deformation of the crystallite and the continuum displacement field. A standard approach for doing so is based on the Cauchy–Born rule (see, for example, Ericksen, 1984) according to which the atomic positions are related to the continuum fields through a local deformation gradient $\mathbf{F}$ (which, for infinitesimal deformation, is the usual $\nabla \mathbf{u}$, where $\mathbf{u}$ is the infinitesimal displacement) which is applied to the crystal's undeformed lattice basis. The crystal is then reconstructed from the altered base vectors. In this manner each continuum particle is represented by an infinite crystal undergoing homogeneous deformation. This limit is referred to as the local QC formulation. The key idea is that, since the energy of each point is obtained directly from atomistic simulations, important properties of the crystal, such as its symmetries, are automatically introduced into the description of the material. However, despite being elegant and straightforward to implement, the local QC formulation suffers from several shortcomings that make it necessary to develop a formulation that can deal with non-local effects. Some of such shortcomings, as discussed by Tadmor *et al.*, are as follows.

(1) The most important shortcoming of the QC formulation is that, due to the homogeneous nature of the deformation in this formulation, it is not possible

to model important heterogeneities, such as stacking faults, which are two undeformed crystalline half spaces slipping over each other by a non-lattice translation vector, and are therefore non-uniform. However, within the local QC formulation, such structures can only be modeled via a simple shear deformation which, except for the two atomic layers directly adjacent to the slip plane, results in a structure which is completely different from what is found in real stacking faults.

(2) Another difficulty with a local formulation is that it does not allow for interface defects, such as free surfaces, grain boundaries or other heterogeneous interfaces. Moreover, the lattice parameter, which serves as the crystal's intrinsic length scale, is lost in the Cauchy–Bohr process, hence allowing the energy minimization methods to develop structure on sublattice length scales, which is clearly unphysical. Such deficiencies are critical near the defects, hence emphasizing the need for an approach that can deal with non-local effects.

Such difficulties are circumvented by a non-local QC formulation of the problem. In this approach each atom within the representative crystallite is displaced according to the actual continuum displacement field at its position, implying that the position $\mathbf{r}_n$ of the atom after deformation is given by, $\mathbf{r}_n = \mathbf{R}_n + \mathbf{u}(\mathbf{R}_n)$, where $\mathbf{u}$ is the continuum displacement field, and $\mathbf{R}_n$ is the atom's position before deformation. The local and non-local formulations are equivalent as long as the FE is large enough to entirely contain the representative crystallite centered about its quadrature point; see Figure 10.1. However, as the elements become smaller than $R_c$, members of the representative crystallite will fall inside different elements and experience a non-uniform displacement field; see Figure 10.2. This allows modeling of stacking faults in a straightforward manner. The use of non-local elements near the stacking fault plane captures the true non-uniform deformation. The non-



FIGURE 10.1. Local quasi-continuum/finite-element in which the triangle corners represent the nodes, while the circles are atoms that belong to the crystallite (after Tadmor *et al.*, 1996a).

Representative Atom

FIGURE 10.2. Non-local quasi-continuum/finite-element (solid triangle) surrounded by nearby elements (dashed triangles) (after Tadmor *et al.*, 1996a).



FIGURE 10.3. Interfacial effects represented by a grain boundary (after Tadmor *et al.*, 1996a).

local formulation also allows treatment of problems that involve, for example, grain boundaries. If elements smaller than the representative crystallite radius $R_c$ are placed near such an interface, they will, due to non-locality, contain atoms that are arranged in a different crystal orientation, and thus mimic a grain boundary; see Figure 10.3.

### 10.2.2   Constitutive Models

Tadmor *et al.*'s use of atomistically-derived constitutive relations has the important property that it preserves all the relevant crystal symmetries. The most important of symmetry properties to be preserved is slip-invariance, which expresses the fact that the energy of a solid material is invariant under crystallographic slip, and is described by the one-parameter family of deformation mappings:

$$\mathbf{H}(\epsilon_s) = \mathbf{U} + \epsilon_s \mathbf{s} \otimes \mathbf{n}, \tag{1}$$

where $\mathbf{n}$ is the normal to the slip plane, $\mathbf{s}$ is a vector in the slip direction, $\epsilon_s$ is the slip strain, and $\mathbf{U}$ is the identity matrix. Because of lattice invariance, the energy density $E(\epsilon_s)$ is periodic in $\epsilon_s$ with period $b/d$, where $b = |\mathbf{b}|$ is the magnitude of the translation vector, and $d$ is the distance between adjacent crystallographic planes perpendicular to $\mathbf{n}$. The slip invariance plays a vital role in allowing for the presence of dislocations and other stable defects in the crystal.

A global origin in the undeformed configuration, relating the continuum FE model to the underlying crystal structure, is now set. Then, for every quadrature point in the FE mesh (where the continuum fields are sampled for numerical integration) the nearest atom is selected as a representative atom, and a small neighborhood around it is deformed either according to the local deformation gradient tensor $\mathbf{F}$ (see above), if the element is local, or based on the actual displacement of every atom according to the global continuum displacement fields, if the non-local formulation is used (see below). The total energy of the representative atom is then computed using an appropriate atomistic model. The computed energy and its derivatives at the quadrature point are then supplied to the FE model. By refining the mesh and using the non-local model in highly strained regions, the effect of nonlinear core will be taken into account, while in less strained regions far from the core the local approximation will yield linear elastic behavior, as it should.

### 10.2.3   The Atomistic Model

As the atomistic model, Tadmor *et al.* (1996a) utilized the embedded-atom method (EAM) (see Section 9.7.1) to compute total energies for their system. An important flaw of EAM is that, it grossly underestimates the stacking fault energies, thus hampering its ability for modeling dislocations. However, Ercolessi and Adams (1993) used *ab initio* calculations (see Chapter 9) to fit the parameters of the EAM to a range of material properties, thus enabling the EAM to yield substantially larger intrinsic stacking fault energy. The Ercolessi–Adams potential was utilized by Tadmor *et al.* (1996a), although, in principle, any other atomistic model could have been used.

### 10.2.4   Field Equations and Their Spatial Discretization

The next step is to specify the field equations so that they can be discretized for use in the numerical simulations. One considers a crystal that occupies a reference configuration $B_0$ in $R^3$, represented by a material Cartesian coordinate system,

$\{X_I, \ I = 1, 2, 3\}$. The crystal undergoes a motion described by a deformation mapping $\boldsymbol{\Phi}(\mathbf{X}, t)$. The image of $B_0$ by $\boldsymbol{\Phi}(\cdot, t)$ defines the deformed configuration $B_t$ of the crystal at time $t$, represented by a spatial coordinate system, $\{x_i, \ i = 1, 2, 3\}$. The deformation at time $t$ of an infinitesimal material neighborhood $d\Omega_0$ around a point $\mathbf{X}$ of $B_0$ is completely defined by the linear part of $\boldsymbol{\Phi}(\cdot, t)$ at $\mathbf{X}$. This defines an affine mapping given by

$$dx_i = F_{iJ}(\mathbf{X}, t)dX_J, \tag{2}$$

where $F_{iJ}$ are the components of the deformation gradient tensor $\mathbf{F}$ given by

$$F_{iJ}(\mathbf{X}, t) = \phi_{i,J}(\mathbf{X}, t), \tag{3}$$

where upper-case indices refer to the material frame, lower-case indices to the spatial frame, and $(\cdot), J$ indicates differentiation with respect to $X_J$. In invariant notation, $\mathbf{F} = \boldsymbol{\nabla}_0 \boldsymbol{\Phi}$, where $\boldsymbol{\nabla}_0$ denotes the material gradient operator.

The reference boundary $\partial B_0$ is now partitioned into a Dirchlet (displacement) component $\partial B_{01}$ and a Neumann (traction) component $\partial B_{02}$. A given displacement $\bar{\boldsymbol{\Phi}}$ is imposed on $\partial B_{01}$, and $\partial B_{02}$ is subjected to a given traction $\bar{\mathbf{T}}$. Moreover, body forces per unit volume $\rho_0 \mathbf{B}$ act on the solid material, where $\rho_0$ is the reference mass density and $\mathbf{B}$ is the body force field per unit mass. Stable configurations of the crystal are those that minimize the total potential energy given by

$$E[\boldsymbol{\Phi}] = \inf_{\boldsymbol{\psi}} \left( \int_{B_0} E_s(\boldsymbol{\psi})d\Omega_0 - \int_{B_0} \rho_0 \mathbf{B} \cdot \boldsymbol{\psi} d\Omega_0 - \int_{\partial B_{02}} \bar{\mathbf{T}} \cdot \boldsymbol{\psi} dS_0 \right), \tag{4}$$

where $E_s(\boldsymbol{\psi})$ is the strain energy density computed from the EAM (or any other appropriate atomistic model), and the trial deformation mappings $\boldsymbol{\psi}$ belong to some suitable space of functions over $B_0$ which satisfy the boundary conditions $\boldsymbol{\psi} = \bar{\boldsymbol{\Phi}}$ on $\partial B_{01}$.

If the deformation of the crystal is small, the above formulation reduces to conventional anisotropic elasticity in which energy minimizers are uniquely defined up to a rigid body motion, conditions that are realized in regions of the crystal *far from* the lattice defects. However, as discussed above, the fact that the strain energy density $E_s$ is computed from an atomistic potential implies that it is possible for energy minimizers to develop microstructures on a fine scale, including lattice's defects such as dislocations. On this scale, the periodicity of the lattice, and the resulting periodicity of the energy function with respect to crystallographic slip, are crucial.

Because of the multiscale nature of the problem, the field equations are discretized by an adaptive FE method. First, the reference configuration $B_0$ is partitioned into FEs $\{\Omega_h^e, e = 1, \cdots, M\}$, where $M$ is the number of elements, and $h$ is a measure of the size of the smallest element. The deformation mapping and deformation gradients are discretized by the standard FE method (see, for example, Strang and Fix, 1973):

$$\boldsymbol{\Phi}_h(\mathbf{X}, t) = \sum_{a=1}^N \boldsymbol{\Phi}_a(t)\mathcal{N}_a(\mathbf{X}), \tag{5}$$

$$\mathbf{F}_h(\mathbf{X}, t) = \sum_{a=1}^{N} \mathbf{\Phi}_a(t) \nabla_0 \mathcal{N}_a(\mathbf{X}), \tag{6}$$

where $a = 1, \cdots, N$ denotes the nodes in the mesh, with $N$ being the number of nodes, $\mathbf{\Phi}_a(t)$ are the nodal coordinates at time $t$, and $\mathcal{N}_a(\mathbf{X})$ are the interpolation functions. The main unknowns of the problem are now the nodal coordinates $\mathbf{\Phi}_a(t)$ which are obtained from the constrained minimization problem (4) in which the $\boldsymbol{\psi}$ are replaced by $\boldsymbol{\psi}_h$, with the trial functions $\boldsymbol{\psi}_h$ being of the following form,

$$\boldsymbol{\psi}_h(\mathbf{X}) = \sum_{a=1}^{N} \boldsymbol{\psi}_a \mathcal{N}_a(\mathbf{X}). \tag{7}$$

These trial functions must satisfy the boundary conditions identically on $\partial B_{01}$. All integrals in Eq. (4), when written in terms of $\boldsymbol{\psi}_h$, can be conveniently computed by numerical quadrature at the element level, which reduces all stress-strain calculations to the quadrature points of the elements. Tadmor *et al.* (1996a) used linear three-noded triangular elements with one-point quadrature rule, and constructed all the FE meshes by automatic triangulation based on the Delauney algorithm (see, for example, Sloan, 1987). The constrained minimization problem was solved by a conjugate-gradient approach (see Section 9.5.2), followed by a Newton–Raphson iteration when the initial guess was too far from the solution.

## 10.2.5  *Local Quasi-continuum Formulation*

Given the continuum deformation fields at a quadrature point, one must now compute the energy and its variations at that point. In the local quasi-continuum formulation, each point of a solid material is represented locally by an infinite crystal which is deformed homogeneously, resulting in the loss of the global origin that links the underlying crystal lattice to the continuum. Therefore, the choice of representative atom is unimportant since all atoms are equivalent. Thus, one can assume that the infinite crystal surrounds a representative atom at the origin. Then, the Cauchy–Born approximation (see, for example, Ericksen, 1984) is used so that the infinite crystal is deformed according to the local continuum deformation gradient. Consequently, if $\{\mathbf{A}_I, I = 1, 2, 3\}$ is a crystal basis, then the coordinates of its atoms are given by

$$\mathbf{X}(\mathbf{m}) = m_I \mathbf{A}_I, \quad \mathbf{m} \in Z^3, \tag{8}$$

where $Z$ represents the set of integers. The positions of the atoms in the deformed configuration are then taken to be

$$\mathbf{x}(\mathbf{m}) = \mathbf{F}\mathbf{X}(\mathbf{m}), \quad \mathbf{m} \in Z^3, \tag{9}$$

where $\mathbf{F}$ is the local deformation gradient which is constant within the element. In practice, a region of radius $R_c$ (taken to be about twice the cutoff radius $r_c$ of the atomistic potential) represents the infinite crystal. The applied trial deformation $\mathbf{F}$ should not be so large so as to bring atoms from outside the region $R_c$ to

within the cutoff radius $r_c$. To account for this effect, Tadmor *et al.* introduced the concept of an *influence radius* $R_i$ associated with deformation $\mathbf{F}$, which is the radius that corresponds to the most distant point in the undeformed configuration that is mapped onto the representative atom's cut-off sphere in the deformed configuration. Tadmor *et al.* showed that $R_i$ is given by

$$R_i = r_c \sqrt{\lambda_{\max}}, \tag{10}$$

where $\lambda_{\max}$ is the largest eigenvalue of $(\mathbf{F}^{-1})^{\mathrm{T}} \mathbf{F}^{-1}$. Then, every trial deformation during the minimization process must satisfy the constraint, $R_i \leq R_c$.

Given an acceptable trial deformation, the strain energy density $E_s$, which is a function of $\mathbf{F}$, is computed using the atomistic method. The local contributions to the out-of-balance force residual and global stiffness matrix follow from Eq. (4) (when written in terms of $\boldsymbol{\psi}_h$) as

$$\frac{\partial E}{\partial \boldsymbol{\psi}_a} = \sum_e^{\mathrm{local}} \int_{\Omega_h^e} (\mathbf{P} \cdot \boldsymbol{\nabla}_0 \mathcal{N}_a) d\Omega_0 - \int_{B_0} \rho_0 \mathbf{B} \mathcal{N}_a d\Omega_0 - \int_{\partial B_{02}} \bar{\mathbf{T}} \mathcal{N}_a d S_0, \tag{11}$$

$$\frac{\partial^2 E}{\partial \boldsymbol{\psi}_a \partial \boldsymbol{\psi}_b} = \sum_e^{\mathrm{local}} \int_{\Omega_h^e} [\mathbf{C} : (\boldsymbol{\nabla}_0 \mathcal{N}_a \otimes \boldsymbol{\nabla}_0 \mathcal{N}_b)] d\Omega_0, \tag{12}$$

where $a$ and $b$ are node numbers, $\mathbf{P} = \partial W / \partial \mathbf{F}$ is the first Piola–Kirchhoff stress tensor (see Section 7.11.3), and $\mathbf{C} = \partial^2 W / \partial \mathbf{F}^2$ is the Lagrangian tangent stiffness tensor. These tensors are the finite deformation analogues of the usual Cauchy stress and elastic modulus tensors in linear elasticity described in Chapter 7 of Volume I. Note that the first terms of Eqs. (11) and (12) are sums only over *local* elements $e$ in $B_0$. In terms of the components of $\mathbf{P}$ and $\mathbf{C}$, Eqs. (11) and (12) are rewritten as

$$\frac{\partial E}{\partial \psi_a^i} = \sum_e^{\mathrm{local}} \int_{\Omega_h^e} P_{iJ} \mathcal{N}_{a,J} d\Omega_0 - \int_{B_0} \rho_0 B_i \mathcal{N}_a d\Omega_0 - \int_{\partial B_{02}} \bar{T}_i \mathcal{N}_a d S_0, \tag{13}$$

$$\frac{\partial^2 E}{\partial \psi_a^i \partial \psi_b^k} = \sum_e^{\mathrm{local}} \int_{\Omega_h^e} [C_{iJKL} \mathcal{N}_{a,J} \mathcal{N}_{b,L}] d\Omega_0, \tag{14}$$

where

$$P_{iJ} = \frac{\partial E_s}{\partial F_{iJ}} = \sigma_{ij} F_{Jj}^{-1}, \tag{15}$$

$$C_{iJKL} = \frac{\partial^2 E_s}{\partial F_{iJ} \partial F_{kL}} = (c_{ijkl} + \delta_{ik} \sigma_{jl}) F_{Jj}^{-1} F_{Ll}^{-1}, \tag{16}$$

where the relations between $\mathbf{P}$ and $\mathbf{C}$ on one hand, and their spatial counterparts, the Kirchhoff stress $\boldsymbol{\sigma}$ and the spatial moduli $\mathbf{c}$, on the other hand have been used, with $\delta_{ik}$ being the Kronecker delta. Tadmor *et al.* utilized a three-noded linear element, which means that $\mathbf{F}$ was constant within each element. Thus, the integrals in Eqs. (13) and (14) are computed by evaluating the integrands for the value of $\mathbf{F}$ within each element and multiplying the result by the area of that element.

The treatment of the problem presented so far is completely general and independent of the atomistic model. The appropriate expressions for the components of the stiffness tensor $\mathbf{C}$ were already given in Chapter 9, Eqs. (9.162)–(9.164). The corresponding components of the spatial moduli tensor $\mathbf{c}$ are then computed from Eq. (16). For a purely local formulation, the expressions for $\mathbf{C}$ and $\mathbf{c}$ represent a complete constitutive description of the problem. However, as discussed above, a purely local formulation is unable to capture non-uniform effects, such as stacking faults. Therefore, one must resort to a non-local formulation which is described in the next section. Note, however, that for a pure infinite crystal which is deformed homogeneously, translational invariance reduces the general expression for the EAM to that of a single atom and all of its neighbors that are within a pre-specified cut-off radius $r_c$.

## 10.2.6  Nonlocal Quasi-continuum Formulation

In this formulation each quadrature point is represented by a single atom with neighbors that are displaced in accordance with the continuum displacement fields. This is shown in Figure 10.2. The global lattice is retained, and one must explicitly account for the position of the atoms in the pre-presentative crystallite relative to the continuum mesh. Thus, the global Cartesian coordinates $\mathbf{R}^e$ of the representative atom of element $e$ are introduced. The positions of the atoms belonging to the representative crystallite are written as

$$\mathbf{X}^e(\mathbf{m}) = \mathbf{R}^e + m_I \mathbf{A}_I, \quad \mathbf{m} \in Z^3. \tag{17}$$

Equation (17) connects the Bravais lattice to the FE mesh, and establishes a one-to-one relationship between the atomic sites and the continuum fields. The deformed atomic positions are obtained by interpolation from the FE mesh through,

$$\mathbf{x}^e(\mathbf{m}) = \mathbf{X}^e(\mathbf{m}) + \boldsymbol{\psi}_a \mathcal{N}_a[\mathbf{X}^e(\mathbf{m})]. \tag{18}$$

Unlike the local case, it is not possible in the non-local formulation to define general measures of stress and stiffness, since a uniform strain field does not exist. As a result, the out-of-balance force residual and global stiffness must be written explicitly in terms of the EAM. We write the energy density of the system as [see Eq. (9.159)]

$$E = \frac{1}{\Omega} \sum_i \left[ \frac{1}{2} \sum_{j \neq i} U_{ij}(r_{ij}) + \mathcal{E}_i(\rho_i) \right] = \frac{1}{\Omega} \sum_i [U_i + \mathcal{E}_i(\rho_i)], \tag{19}$$

where we have used a slightly different notation than what was used in Eq. (9.159) in order to avoid confusion. As before, $\rho_i$ is the local embedding density around atom $i$. Then, the non-local contribution to the out-of-balance force residual is given by

$$\frac{\partial E}{\partial \psi_a^i} = \sum_e^{\text{non-local}} \left\{ \frac{1}{\Omega} \sum_{\mathbf{m}} \left[ \left( U'(\rho_e)\rho'(r_e^{\mathbf{m}}) + \frac{1}{2}\mathcal{E}'(r_e^{\mathbf{m}}) \right) \frac{\partial r_e^{\mathbf{m}}}{\partial \psi_a^i} \right] \Omega_h^e \right\}, \tag{20}$$

where the sum over the elements $e$ is over only non-local elements in $B_0$, $\Omega_h^e$ is the area of element $e$, and the total electron density $\rho_e$ at the representative atom of element $e$ is given by [see Eq. (9.160)]

$$\rho_e = \sum_{\mathbf{m}} \rho(r_e^{\mathbf{m}}). \tag{21}$$

Here $\rho$ is the electron density function, and

$$r_e^{\mathbf{m}} = |\mathbf{x}^e(\mathbf{m}) - \mathbf{x}^e(\mathbf{0})| = |m_I \mathbf{A}_I + \boldsymbol{\psi}_a \{\mathcal{N}_a[\mathbf{X}^e(\mathbf{m})] - \mathcal{N}_a(\mathbf{R}^e)\}|. \tag{22}$$

In Eq. (20)

$$\frac{\partial r_e^{\mathbf{m}}}{\partial \psi_a^i} = \{\mathcal{N}_a[\mathbf{X}^e(\mathbf{m})] - \mathcal{N}_a(\mathbf{R}^e)\} \frac{(r_i^{\mathbf{m}})_e}{r_e^{\mathbf{m}}}, \tag{23}$$

where

$$(r_i^{\mathbf{m}})_e = x_i^e(\mathbf{m}) - x_i^e(\mathbf{0}). \tag{24}$$

One should keep in mind that for FE meshes that contain both local and non-local elements, the total out-of-balance force residual $\partial E / \partial \boldsymbol{\psi}_a$ is calculated as a superposition of the two vectors given by Eqs. (13) and (20). The non-local contributions to the global stiffness matrix are similarly calculated:

$$\Omega \frac{\partial^2 E}{\partial \psi_a^i \partial \psi_b^j} = \sum_{e}^{\mathrm{non-local}} U''(\rho_e) \left[ \sum_{\mathbf{m}} \rho'(r_e^{\mathbf{m}}) \frac{\partial r_e^{\mathbf{m}}}{\partial \psi_a^i} \right] \left[ \sum_{\mathbf{n}} \rho'(r_e^{\mathbf{n}}) \frac{\partial r_e^{\mathbf{n}}}{\partial \psi_b^j} \right]$$

$$+ \sum_{\mathbf{m}} \left[ \left( U'(\rho_e) \rho''(r_e^{\mathbf{m}}) + \frac{1}{2} \mathcal{E}''(r_e^{\mathbf{m}}) \right) \frac{\partial r_e^{\mathbf{m}}}{\partial \psi_a^i} \frac{\partial r_e^{\mathbf{m}}}{\partial \psi_b^j} + \left( U'(\rho_e) \rho'(r_e^{\mathbf{m}}) + \frac{1}{2} \mathcal{E}'(r_e^{\mathbf{m}}) \right) \frac{\partial^2 r_e^{\mathbf{m}}}{\partial \psi_a^i \partial \psi_b^j} \right]. \tag{25}$$

Similar to the total out-of-balance force residuals, FE meshes that contain both local and non-local elements yield a stiffness matrix which is obtained by super-position of both Eqs. (14) and (25). The inclusion of non-local elements near highly deformed regions, such as stacking faults, completes the QC formulation.

### 10.2.7   The Criterion for Nonlocality of Elements

The procedure for the computations is completed once one introduces a criterion for determining the status of an element in the FE mesh in terms of it being local or non-local. Although it may seem that a natural criterion would be to consider, and locally compute, elements that are larger than the local crystallite radius $R_c$ and treat the smaller elements non-locally, this criterion would be wasteful from a computational point of view because the non-local QC formulation is needed only close to defect cores and along slip planes where stacking faults develop. Far away from such highly inhomogeneous regions the local formulation, which is less computationally intensive and also more stable, should perform well. Therefore, as pointed out by Tadmor *et al.* (1996a), it is of interest to develop a criterion which

is capable of identifying regions that undergo large inhomogeneous deformation. Based on such considerations, the criterion that they developed is as follows.

Consider the second invariant $II_\epsilon$ of the Lagrangian strain tensor $\epsilon$:

$$II_\epsilon = \frac{1}{2}\left[\epsilon : \epsilon - \mathrm{tr}(\epsilon)^2\right] = \epsilon_{12}^2 + \epsilon_{13}^2 + \epsilon_{23}^2 - (\epsilon_{11}\epsilon_{22} + \epsilon_{22}\epsilon_{33} + \epsilon_{11}\epsilon_{33}), \quad (26)$$

where

$$\epsilon = \frac{1}{2}(\mathbf{F}^T\mathbf{F} - \mathbf{U}). \quad (27)$$

Then, an element is considered as non-local if

$$\sqrt{|II_\epsilon|} > \epsilon_c, \quad (28)$$

where $\epsilon_c$ is a critical strain, the value of which depends on the material and the phenomenon under study. In addition to the elements that satisfy criterion (28), Tadmor *et al.* imposed the additional condition that elements in their immediate vicinity that share atoms with those elements deemed to be non-local by the strain criterion, should be also treated as non-local. Such elements are referred to as non-local *by proximity*.

The task of selecting the non-local elements is not yet complete though. Consider, for example, a material that contains a dislocation, as shown in Figure 10.4, where only elements along the slip plane are shown. In such a case, elements far to the right of the dislocation core will be undeformed, while elements far to the left experience perfect Burgers vector slip. In both cases, these are zero energy modes that correspond to an undeformed crystal which can be treated by local elements. However, with the purely kinematical criterion (28), the elements on the far left of the slip plane will be identified as non-local. Therefore, Tadmor *et al.* (1996a) added a more stringent non-locality criterion to the kinematic criterion (28): In addition to being highly strained, the deformation within the element must produce a non-zero strain energy $E_s$ such that

$$E_s > E_0 \quad (29)$$

where $E_0$ is typically small.

Summarizing, the procedure for determining the status of all elements in the FE mesh is as follows.



FIGURE 10.4. Slip plane elements near a dislocation core (after Tadmor *et al.*, 1996a).

(1) One computes the Lagrangian strain tensor $\boldsymbol{\epsilon}$ for all elements in the FE mesh.

(2) Condition (28) is checked for all elements smaller than $R_c$ in order to identify non-locally strained elements.

(3) The energy of elements identified in (2) is computed by using the non-local formulation of the problem. Only those elements that satisfy criterion (29) are retained and considered as non-local.

(4) The elements that are non-local by proximity are located by identifying all elements that are smaller than $R_c$ and are within a distance $R_c$ of an element that satisfies both criteria (28) and (29).

(5) In addition, all surface and interfacial elements are computed non-locally.

Tadmor *et al.* (1996a) carried out this procedure at the start of each iteration and integrated it into the solution process. As they pointed out, an important point regarding this algorithm is that, to ensure convergence once an element is identified as non-local, it must remain so from that point on even if in future iterations it no longer satisfies the non-locality criterion, or is no longer in proximity to a non-local element.

### 10.2.8  Application to Stacking Faults in FCC Crystals

An important question is whether the above theory can stably support lattice defects, such as dislocations, and, if so, how similar are their core structures to those predicted by a full atomistic simulation. To answer these questions, Tadmor *et al.* (1996a) considered three FCC configurations: Stacking faults within the (111) plane, the Lomer edge dislocation $(001)[1\bar{1}0]$, and the primary FCC edge dislocation $(111)[1\bar{1}0]$. In all cases the predictions of the atomistic-continuum model were shown to be in very good agreement with the lattice statics results, under conditions in which the same potentials were used in addition to equivalent boundary conditions. Here, we describe their simulations for the stacking faults. In what follows the constitutive behavior of the crystal is modeled using the embedded-atom potentials for aluminum due to Ercolessi and Adams (1993) with a cut-off radius, $r_c = 5.56$ Å, the lattice parameter, $a_0 = 4.032$ Å, a representative crystallite radius, $R_c = 9.87$ Å, which for the perfect FCC lattice corresponds to 12 neighbor shells with 249 atoms, a critical non-local strain $\epsilon_c = 10\%$ (which was used in the non-locality criterion), and a zero energy tolerance, $E_0 = 10^{-3}$ eVÅ$^{-3}$.

Stacking faults (SFs) often arise in crystalline deformation processes. In FCC crystals, for example, they are known to form on the dominant $\{111\}\langle 110\rangle$ slip system, as a result of, for example, dissociation of a perfect edge dislocation into two Shockley partial dislocations via the reaction,

$$\frac{1}{2}[\bar{1}10] \to \frac{1}{6}[\bar{2}11] + \frac{1}{6}[\bar{1}2\bar{1}].$$

The two partials can then drift apart and leave a stacking fault ribbon between them. However, due to the imperfect stacking sequence in this ribbon,

there is an energy penalty—the SF energy—associated with its presence which limits the indefinite drift of the partials. The expected out-of-plane displacement jump across the stacking fault plane is $\sqrt{6}a_0/12$, which for aluminum is about 0.823 Å.

As the first step, the SF energy was directly calculated using lattice statics (LS). The unrelaxed SF energy resulting from the use of the Ercolessi–Adams potentials was 7.530 meVÅ$^{-2}$. The LS computation indicated that the atoms on each of the four (111) planes adjacent to the slip plane (two above and two below) contribute equally to the SF energy. This information is useful to the construction of the FE mesh. On relaxation, the SF energy was found to reduce to about 6.5 meVÅ$^{-2}$, which is in reasonable agreement with the observed experimental values for aluminum, which is in the range 7.5–9 meVÅ$^{-2}$.

Tadmor *et al.* (1996a) carried out two different analyses using the QC-FE method. In the first analysis the FE mesh was generated in the plane that contained the SF. The $x$-direction was taken to coincide with the Shockley partial direction [$\bar{2}11$], while the $y$-direction was set to the slip plane normal [111]. To introduce the initial slip into the model, all nodes above the slip plane to the right by the magnitude of the Shockley partial $a_0\sqrt{6}$ were removed, and then all boundary nodes were constrained to their initial positions. The resulting mesh, shown in Figure 10.5, contained 42 elements and 32 nodes, hence resulting in 36 unconstrained degrees of freedom. As can be seen in this figure, the mesh contains two different types of element. One is the large, nearly equilateral elements, which are



FIGURE 10.5. Finite-element mesh for initially distorted stacking faults (after Tadmor *et al.*, 1996a).

FIGURE 10.6. Underlying crystal structure superimposed on finite elements. Filled points are atoms with zero depth, while open ones represent atoms with 1.426 Å depth (after Tadmor *et al.*, 1996a).

larger than $R_c$ and are thus local. The elements of the second type are long and narrow, and are considered non-local, but because their width is larger than $R_c$, they exhibit non-local effects only in the $y$-direction. Therefore, surface effects play no role. The height of the narrow elements was set equal to the interplanar distance in the (111) direction, $a_0/\sqrt{3}$, and the global origin was selected such that these elements fell between adjacent atomic planes; see Figure 10.6. Hence, the elements that straddle the slip plane were allowed to act as a kinematical mechanism for introducing slip. Thus, from an atomistic point of view, there is a jump in displacement across the slip plane, as expected, while in the continuum there exists a continuous linear variation in slip. However, because the energy is computed atomistically, the manner in which the slip is distributed in the continuum is unimportant. Five layers of the narrow elements were necessary to capture the contributions of the four atomic planes adjacent to the slip plane to the SF energy; see Figure 10.7. In this model, the SF energy is equal to the total energy per unit thickness, divided by the width of the system. The model predicts that the unrelaxed SF energy is identical to that obtained from the LS analysis. Following a Newton–Raphson minimization, Tadmor *et al.* (1996a) found a relaxed SF energy of about 6.2 meVÅ$^{-2}$ which is smaller than, but comparable with, the LS value.

In the second analysis, the FE model was generated in the primary [$\bar{1}$10]-[111] coordinate system, and a deformation corresponding to a 1/2 Burgers vector slip

was used. On relaxation, Tadmor *et al.* (1996a) found the SF energy to be the same as in the previous Shockley partial analysis, and out-of-plane displacements between 0.822 and 0.824 Å were observed in all nodes above the slip plane, which are very close to the expected value.

### 10.2.9 Application to Nanoindentation

Nanoindentation is a process that can be classified as a problem in small-scale contact mechanics, but has major technological applications. Indenters that have a radius of curvature of the order of 50–100 nm are now used routinely in experiments. Recent experiments (see, for example, Zielinski *et al.*, 1995; Gouldstone *et al.*, 2000) measured load-displacement curves and subsurface dislocations. One of the important questions that arises in this setting concerns the conditions attendant to dislocation nucleation. Upon indentation, and after an initial elastic stage, the onset of permanent deformation is mediated by the nucleation and propagation of dislocations. Analysis of the recent experiments indicate that the number of dislocations activated by such processes are typically about 100 or less, and as such they are amenable to effective atomistic simulations (see, for example, Kelchner *et al.*, 1998). However, in such simulations, the indenter sizes that can be considered are much smaller than the experimentally-utilized values, which may in

turn cause premature nucleation of dislocations relative to experimental observations. Likewise, the size of the computational domain is necessarily limited, and therefore the dislocations soon run up against artificial boundaries. In addition, within a strict atomistic simulations, it is difficult to account for the effect of long-range elastic stresses that are, for example, present in a thin film-substrate system. All of these issues indicate that the multiple length scales that are involved must be explicitly accounted for, if a nanoindentation process is to be modeled accurately.

Before any attempts are made for applying a mixed atomistic-continuum model to indentation, a few key issues must be considered. For example, Sharp *et al.* (1993) devised a type of "phase diagram" which divides the indentation response into elastic, elastic-plastic, and brittle regimes. If the phase diagram is truly representative, then the critical question to be asked is, how does this phase diagram depend on the constituent atoms that make the solid? This is particularly important, as one must find a way of describing the plastic deformation that takes place beneath the indenter. Moreover, force-displacement relations that are the result of indentation have been measured frequently. The question is whether such relations can be predicted, and if so, whether they can be used as a way of differentiating between different materials. Finally, indentation of a clean surface is different from one that is coated by another material. Therefore, one must find a way of connecting the two, if this is possible at all.

Tadmore *et al.* (1996b) applied their QC model to the problem of nanoindentation. They considered a pseudo-2D model of indentation with both rectangular punches and rounded indenters. The former case is, however, difficult as it leads to nucleation of dislocations after shallow indents. A crystal orientation such that (111) planes are perpendicular to the face being indented was considered, since this orientation permits the development of conventional FCC edge dislocations discussed above. A challenging problem in such simulations is the ongoing mesh refinement that must accompany such processes as dislocation nucleation and emission. During the indentation process, the large strains under the punch triggers automatic mesh adaption, leading to refinement in areas of interest. Thus, as dislocations nucleate, the mesh refinement follows the paths of the dislocations, hence permitting the dislocations to move away from the punch into the substrate.

Figure 10.8 shows the atomic positions obtained from the mixed atomistic-continuum model using rectangular indenter for an indentation of 5.4 Å. As can be seen, partial dislocations have developed that are separated by stacking faults. Figure 10.9 shows the atomic positions when a rounded indenter was used. In this case the response is elastic for indentation depths that are deeper than those for which the first dislocation was nucleated in the rectangular indenter case, thus raising the possibility of the existence of a connection between the geometry of the indenter and the criterion for emission of dislocations. Knap and Ortiz (2001) extended Tadmore *et al.*'s (1996b) and presented a full 3D QC analysis of the early stages of nanoindentation in gold thin films.

FIGURE 10.8. Atomic positions obtained from simulation of the quasi-continuum theory, using a rectangular indenter for an indentation of 5.4 Å (after Tadmor *et al.*, 1996b).



FIGURE 10.9. Atomic positions obtained from the solution of the quasi-continuum theory, for a rounded indenter for an indentation of 4.6 Å (after Tadmor *et al.*, 1996b).

## 10.3    Fracture Dynamics: Joining Tight-Binding, Molecular Dynamics, and Finite-Element Computations

The second example that we describe is the multiscale modeling of dynamics of fracture propagation in silicon that was developed by Broughton *et al.* (1999; see also Abraham *et al.*, 1998a,b). As they pointed out, in order to obtain meaningful ensemble averages, or access time scales of practical use, one must be able to simulate the system through times of the order of 1 nanosecond. A typical MD time step is about $10^{-15}$ s. Therefore, on a typical parallel machine with about 50–100 processors, one time step is roughly 1 second of real time, implying that $10^6$ time steps may be simulated in approximately 10 days on part of a reasonably powerful parallel machine.

A tight-binding Hamiltonian was used for simulating the quantum-mechanical coupling. The main justification for use of this Hamiltonian was the requirement of computational speed, i.e., the need to simulate a non-trivial number of atoms within 1 second of real time. In addition, a TB Hamiltonian was used for its simplicity and intuitive appeal. At the next coarser level, a molecular dynamics (MD) method was employed for describing the material. The interatomic potential for describing silicon in the MD simulations was the Stillinger–Weber (SW) potential described in Section 9.7.2. As discussed there, one main advantage of the SW potential, in addition to its accuracy, is the ease with which it can be utilized. Finally, at the most coarsened level, FE computations were carried out assuming that linear continuum mechanics is valid. This could be justified by the fact that the FE computations were used only in the far-field (away from the crack tip) region where atoms are perturbed only slightly from equilibrium, and therefore it is unnecessary to employ nonlinear elasticity theory. Moreover, Broughton *et al.* studied a 2D model since the systems of primary interest to them involved plane strain.

Figure 10.10 presents the three primary algorithms, the regions of space that they describe, and the way they are distributed among different processors. Each FE region was assigned to a different processor, and the MD region was domain-decomposed (see Section 9.6.2.3) across several computer nodes. Likewise, the TB region was spread over several processors. The computer program was written in FORTRAN with MPI for the message passing, and hence it was portable to most parallel architecture machines. The advantage of a pseudo-1D topology shown in Figure 10.10 is that much message passing can be performed using the "shift" operator. Moreover, only data within interaction range, defined by the Hamiltonian, must be passed across boundaries, which are defined by the domain decomposition, between processors. This is clearly a parallel multiscale approach.

### 10.3.1    *The Overall Hamiltonian*

Broughton *et al.* (1999) defined a Hamiltonian $\mathcal{H}$ for the entire system. The degrees of freedom of this Hamiltonian are atomic positions **r** and their velocities $\dot{\mathbf{r}}$ for the TB and MD regions, and displacements **u** and their rates of change with time, $\dot{\mathbf{u}}$,

FIGURE 10.10. Domain decomposition of pseudo-1D system that shows the coupling of the length scales. The filled circles show the MD region (after Broughton *et al.*, 1999).



for the FE regions. Conceptually, $\mathcal{H}$ is written as

$$n\mathcal{H} = \mathcal{H}_{\text{TB}}(\{\mathbf{r}, \dot{\mathbf{r}}\} \in \text{TB}) + \mathcal{H}_{\text{MD/TB}}(\{\mathbf{r}, \dot{\mathbf{r}}\} \in \text{MD/TB}) + \mathcal{H}_{\text{MD}}(\{\mathbf{r}, \dot{\mathbf{r}}\} \in \text{MD})$$
$$+ \mathcal{H}_{\text{FE/MD}}(\{\mathbf{u}, \dot{\mathbf{u}}, \mathbf{r}, \dot{\mathbf{r}}\} \in \text{FE/MD}) + \mathcal{H}_{\text{FE}}(\{\mathbf{u}, \dot{\mathbf{u}}\} \in \text{FE}), \tag{30}$$

implying that there are three separate Hamiltonians for each subsystem, as well as Hamiltonians that govern the dynamics of variables at the interface regions between two distinct domains belonging to distinct length scales (Broughton *et al.* referred to such regions as the *handshake regions*). To obtain the equations of motion for all the relevant variables in the system, the appropriate derivatives of $\mathcal{H}$ are taken, and all the variables are updated in time steps using the same integrator. Therefore, the entire time history of the system is obtained numerically, given an appropriate set of initial conditions. Moreover, if one follows a trajectory governed by $\mathcal{H}$, the total energy of the system will be a conserved quantity, ensuring the numerical stability of the simulations.

## 10.3.2   The Tight-Binding Region

Broughton *et al.* (1999) utilized the following TB scheme. As described in Section 9.4.6, semi-empirical TB involves, (a) an *ansatz* for the total energy of the system, and (b) a parameterization of the integrals that occur in the mean-field treatments of the electronic structure of a system. The total energy $E$ of the system is written as

$$E_{\text{TB}} = \sum_{n=1}^{N_o} \epsilon_n + \sum_{i<j} U_r(r_{ij}), \tag{31}$$

where the sum is over all occupied states $N_o$ up to the Fermi level, while the second sum is over all pairs of atoms of the repulsive potential $U_r$. The eigenvalues $\{\epsilon\}$ correspond to the one electron states of a first-principles Hartree–Fock or density functional calculation (see Sections 9.0, 9.1, and 9.4), and are obtained from a

non-orthogonal one-electron Hamiltonian given by

$$\mathcal{H}\Psi_n = \epsilon_n \mathbf{S}\Psi_n. \tag{32}$$

The one-electron wave functions $\{\Psi\}$ were expanded as a linear combination of atomic basis functions $\phi$:

$$\Psi_n = \sum_{i\alpha} c_{i\alpha}^n \phi_{i\alpha}, \tag{33}$$

where the matrix elements of $\mathcal{H}$ and $\mathbf{S}$ were calculated by reducing the equivalent integrals within an extensive database of first-principles calculations to the following parametric forms, expressed in terms of the pairwise functions $h_{\alpha\beta}$ and $s_{\alpha\beta}$,

$$\mathcal{H}_{i\alpha j\beta} \equiv \langle \phi_{i\alpha} | \hat{\mathcal{H}} | \phi_{j\beta} \rangle = h_{\alpha\beta}(\mathbf{r}_{ij}), \tag{34}$$

$$S_{i\alpha j\beta} \equiv \langle \phi_{i\alpha} | \phi_{j\beta} \rangle = s_{\alpha\beta}(\mathbf{r}_{ij}). \tag{35}$$

Here, $n$ denotes the orbital number, while $\alpha$ and $\beta$ label the basis functions which, in the minimal basis of silicon studied by Broughton *et al.* (1999), represent the $s$, $p_x$, $p_y$, and $p_z$ atomic orbitals, and therefore the size of the $\mathcal{H}$ and $\mathbf{S}$ matrices is $4N \times 4N$, where $N$ is the number of atoms in the system. The functions $h_{\alpha\beta}$ and $s_{\alpha\beta}$ smoothly truncate to zero near 5 Å, which is between the third- and fourth-neighbor distances in silicon. The functions $U_r$, $h_{\alpha\beta}$, and $s_{\alpha\beta}$ were obtained by fitting to a database that involved the experimental indirect band gap of the diamond cubic structure and the total energies of crystalline and defective diamond cubic and $\beta$-tin silicon at different densities (Bernstein and Kaxiras, 1997).

The exact form of the basis functions $\phi_{i\alpha}$ are not required, since the one-electron states $\{\Psi\}$ are represented (within this formalism) only by the sets of coefficients $\{c\}$. For given set of atomic coordinates, the coefficients $\{c\}$ were computed by diagonalization. Forces were then computed from the derivative of the TB energy with respect to displacement of the nuclei:

$$\mathbf{f}_i^{\text{TB}} = - \left[ \sum_{n=1}^{N_o} \sum_\alpha c_{i\alpha}^n \sum_{j\beta} c_{j\beta}^n \left( \frac{\partial \mathcal{H}_{i\alpha j\beta}}{\partial \mathbf{r}_i} - \epsilon_n \frac{\partial S_{i\alpha j\beta}}{\partial \mathbf{r}_i} \right) \right] - \sum_{j \neq i} \frac{\partial U_r(r_{ij})}{\partial \mathbf{r}_i}. \tag{36}$$

Orthonormality forces the derivatives of the coefficients with respect to atomic positions to vanish identically. From the knowledge of the forces, atomic coordinates can be computed at any time using exactly the same algorithm as that used for the MD system (see below) with the same time step, since the frequencies in both cases are very similar. In fact, the parameters of the SW potential, used in the MD simulations, are adjusted to ensure equality.

## 10.3.3   Molecular Dynamics Simulation

In the MD region, the interatomic forces were obtained from the SW potential, a full description of which was given in Section 9.7.2. To integrate Newton's

equations of motion, Broughton *et al.* (1999) used the Verlet algorithm (see Section 9.2.2) since it is easily augmented to handle multiple time scale MD simulations. A time step of $\Delta t = 5 \times 10^{-16}$ s was used. The mass $m$ of the silicon atom is $m = 4.6639 \times 10^{-26}$ kg. One can accelerate evaluation of the SW energy and its corresponding forces by taking advantage of atomic neighbor tables, so that computer time scales as $O(N)$, where $N$ is the number of atoms in the system. The two- and three- body terms in the SW potential truncate smoothly to zero just before the second-neighbor distance in zero-pressure diamond cubic structure silicon.

## 10.3.4  Finite-Element Simulation

As mentioned above, Broughton *et al.* (1999) used a FE method for describing the far-field region of the system. Linear elasticity theory was used for developing the FE equations of motion. In this formulation, the total elastic energy of a solid, in the absence of tractions and body forces, is given by

$$E_{\text{FE}} = E_p + E_k,$$

$$E_p = \frac{1}{2} \int \left[ \sum_{i=1}^{3} \sum_{j=1}^{3} \sum_{k=1}^{3} \sum_{l=1}^{3} \epsilon_{ij}(\mathbf{r}) C_{ijkl}(\mathbf{r}) \epsilon_{kl}(\mathbf{r}) \right] d\Omega,$$

$$E_k = \frac{1}{2} \int \rho(\mathbf{r}) \dot{\mathbf{u}}^2(\mathbf{r}) d\Omega, \tag{37}$$

where $\Omega$ is the volume of the system, $E_p$ is the usual Hookian potential energy term which involves the symmetric strain tensor $\epsilon$ and the elastic constant tensor $\mathbf{C}$ (see Chapter 7 of Volume I), $E_k$ is the kinetic energy, and $\rho$ is the mass density. The subscripts $i$, $j$, $k$, and $l$ denote Cartesian directions. The strains and displacements are related through the usual relation,

$$\epsilon_{ij} = \frac{\partial u_i}{\partial x_j} + \frac{\partial u_j}{\partial x_i}. \tag{38}$$

The FE mesh used in this work is described shortly. From the knowledge of the displacements and their time derivatives at the vertices of the FE cells (which were triangular cells) and interpolation functions, one can determine the values of these variables everywhere within each cell. Broughton *et al.* (1999) used linear interpolation inside each cell, and therefore the displacement fields could be represented in piecewise smooth fashion. Equation (37) is now approximated by

$$E_{\text{FE}} = \frac{1}{2} \sum_{m}^{N_c} \sum_{p,q=1}^{6} \left( u_p^m K_{pq}^m u_q^m + \dot{u}_p^m M_{pq}^m \dot{u}_q^m \right), \tag{39}$$

where $\mathbf{K}$ and $\mathbf{M}$ are local stiffness and mass matrices, respectively, $m$ is the cell index, $N_c$ is the total number of FE cells, $\mathbf{u}$ and $\dot{\mathbf{u}}$ defined only at the apices of each triangle, and the sum over $(p, q)$ is over the $(3 \times 2)$ Cartesian directions associated

with the same apices. The stiffness matrix $\mathbf{K}^{(m)}$ associated with the $m$th triangle is given by

$$\mathbf{K}^{(m)} = \frac{L}{4A_m}\mathbf{B}^{\mathrm{T}}\mathbf{C}\mathbf{B}^{(m)}, \tag{40}$$

where $A_m$ is the area of the $m$th triangle, $\mathbf{C}$ is the reduced ($3 \times 3$) elastic constant matrix, $\mathbf{B}^{(m)}$ is the matrix of coordinate differences of the apices of the FE mesh, and $L$ is the thickness of the material in the third dimension. $\mathbf{C}$ depends upon the orientation of the system and is a function of the three basic elastic constants of silicon, namely, $C_{11}$, $C_{12}$, and $C_{44}$, which are given for zero-temperature SW silicon by Ray (1988) and Balamane *et al.* (1992) (these values are listed in Tables 9.4 and 9.6). Broughton *et al.* (1999) used the average of these quantities reported in the literature. $\mathbf{B}^{(m)}$ is given by

$$[\mathbf{B}^{\mathrm{T}}]^{(m)} = \begin{pmatrix} b_1^m & 0 & a_1^m \\ 0 & a_1^m & b_1^m \\ b_2^m & 0 & a_2^m \\ 0 & a_2^m & b_2^m \\ b_3^m & 0 & a_3^m \\ 0 & a_3^m & b_3^m \end{pmatrix}, \tag{41}$$

with

$$b_l^m = y_{l+1}^m - y_{l+2}^m,$$
$$a_l^m = x_{l+2}^m - x_{l+1}^m, \tag{42}$$

where $l = 1, 2$ and 3 denotes the cyclic apex index, and $x$ and $y$ are the 2D FE mesh coordinates with respect to which the displacements $\mathbf{u}$ were defined.

The mass matrix $\mathbf{M}$ requires some care since, in principle, the kinetic energy density varies across any given cell. However, it is necessary to reduce the FE mesh in the interface between FE and MD regions so as to coincide with the perfect atomic lattice. Each atom is apportioned its kinetic energy accordingly. Thus, for the FE mesh, Broughton *et al.* (1999) used the "lumped-mass" approximation which reduces for the smallest mesh size to the atomic limit. In this approximation, one third of the mass in each cell is apportioned to each apex, so that the kinetic energy is given by

$$E_k = \sum_{t=1}^{N_m} M^t (\dot{\mathbf{u}}^t)^2, \tag{43}$$

$$M^t = \frac{\rho L}{3} \sum_{m=1}^{N_c} \sum_{l=1}^{3} \delta_{tm_l} A_m, \tag{44}$$

where $t$ labels the FE mesh points, of which there are $N_m$, and $m_l$ denotes the mesh point index at each of the three apices of cell $m$. The $\dot{\mathbf{u}}$ are vectors of length two since they relate to a mesh point.

Forces that correspond to the displacements in Eq. (39) were computed by taking the spatial derivatives. Displacements and their time derivatives were obtained as functions of time, for given boundary conditions, using the same update algorithm and time step as those used for the MD and TB computations. The force due to the $m$th cell is given by

$$\mathbf{f}_{\text{FE}}^{(m)} = \mathbf{K}^{(m)}\mathbf{u}^{(m)}, \tag{45}$$

where $\mathbf{f}_{\text{FE}}^{(m)}$ and $\mathbf{u}^{(m)}$ are of length six. The total force associated with a mesh point is then the sum of the contributions from each of the cells with apices in common with that point, and

$$\mathbf{f}_{\text{FE}}^{t} = M^{t}\ddot{\mathbf{u}}^{t}. \tag{46}$$

## 10.3.5   Interfacing Finite-Element and Molecular Dynamics Regions

In order to develop the proper model for the interface between the FE and MD regions, two principal issues must be addressed which are, (1) the overlap of the FE mesh with the atoms, and (2) the proper form of the Hamiltonian $\mathcal{H}_{\text{FE/MD}}$.

To address the first issue, Broughton *et al.* (1999) generalized an idea due to Kohlhoff *et al.* (1991). An imaginary surface is drawn between the FE and MD regions. Within the range of the MD interatomic potential (i.e., the SW potential) from this surface, FE mesh points are located at ideal lattice sites. In the absence of diffusion, atoms or mesh points will remain on either side of this interface. However, note that atomic motion may be viewed as displacement around a lattice (mesh) site, and the displacement field may be viewed as motion of an atom away from its perfect site. The same idea can also be used for amorphous materials if a one-to-one mapping of a mesh point to an atom site is made. As one moves away from the interface between the FE and MD regions into the FE region, the mesh spacing may be made larger, with its size being dependent upon the physics of the phenomenon under study.

Since Broughton *et al.* (1999) were interested in examining brittle fracture in Si, they oriented their rectilinear system in such a way that it had (100) faces on all sides. The FE region was represented as a 2D system which, nevertheless, could handle the third dimension in plane strain, since the thickness $L$ of the sample was included in Eq. (44). Thus, a diamond cubic lattice was projected onto a (100) plane. This is shown in Figure 10.11. The FE region can be made periodic in exactly the same way as the MD region. For a periodic system, there are twice as many cells as there are mesh points. Away from the interface region and into the FE region, the mesh was expanded along one dimension (the long axis in Figure 10.10) while the mesh spacing was kept constant in the second dimension. The function chosen for this expansion was of hyperbolic tangent form. Thus, near the interface region there was no expansion of the atomic mesh, while far from the interface region the spacing approached a constant value of ten atomic lattice parameter. The transition region spanned about 200 Å.

FIGURE 10.11. Triangulation of the unit cell shown by dashed lines. The lines on the left and right are continuous, but not those on top and bottom of the unit cell (after Broughton *et al.*, 1999).

To address the issue of the proper Hamiltonian $\mathcal{H}_{\text{FE/MD}}$, Broughton *et al.* (1999) defined a conservative Hamiltonian so as to ensure symplectic time evolution (see Chapter 9) of the atomic and displacement trajectories within the interface region. To conceptualize this Hamiltonian, imagine that two different materials sit on either side of an interface, such that on one side is FE silicon, whereas on the other side one has SW silicon. The cross terms (i.e., the interface Hamiltonian $\mathcal{H}_{\text{FE/MD}}$) can, to first order, be approximated by the average of the two descriptions: All FE triangles that cross the interface contribute half their weight to the Hamiltonian, while the triangles that are fully in the MD region contribute nothing. Similarly, any SW interaction which crosses the interface contributes half its usual weight, while the SW interaction between mesh points, which are fully on the FE side of the interface, contributes nothing to $\mathcal{H}_{\text{FE/MD}}$. Since as discussed above, the atoms and mesh points cannot be distinguished from each other, the SW energy formulation for atomic coordinates **r** and the FE energy formulation for the displacements **u** can be used throughout the interface. The one-to-one mapping of atoms to nodes is no longer needed at distances (from the interface in the FE region) that are greater than twice the SW pair cutoff, which is the distance of greatest three-body range in the SW potential. Figure 10.12 presents such interactions and their ranges. Thus

$$
E_{\text{FE/MD}} = \frac{1}{4} \sum_{m^I=1}^{N_{\text{x}}} \sum_{p,q=1}^{6} u_p^{m^I} K_{pq}^{m^I} u_q^{m^I}
$$
$$
+ \frac{1}{2} \left[ \sum_{(i<j)^I} U_2(r_{(ij)^I}) + \sum_{[i,(j<k)]^I} U_3(\mathbf{r}_{(ij)^I}, \mathbf{r}_{(ik)^I}) \right],
$$

(47)

where $U_2$ and $U_3$ are the two- and three-body terms. Here, the superscript $I$ refers to those interactions that cross the FE/MD boundary [see Eq. (9.176)]. Indeed, as discussed above, $E_{\text{FE/MD}}$ is defined only for interactions that cross the boundary. Equation (47) allows any one atom of the triplet in the three-body terms to be on an opposite side of the interface to the other two.

FIGURE 10.12. The FE/MD interface region. Heavy
lines show the FE triangles that contribute fully
to the interface Hamiltonian, while the thin lines
are those that contribute half weight to the in-
terface Hamiltonian. Dotted lines show two- and
three-body terms of the Stillinger–Weber poten-
tial that cross the boundary and contribute half
weight, while full lines on the right represent full
SW contributions (after Broughton *et al.*, 1999).



Since in the work of Broughton *et al.* (1999), the FE mesh was 2D while the MD
region was 3D, the third dimension of the FE region was treated by mean-field
approximation. Thus, in Eq. (47), $x$- and $y$-displacements of atoms on the MD
side of the FE/MD boundary that contribute to the elastic energy were obtained
by averaging over all equivalent atoms at the depth $z$. In a similar fashion, in
determining the SW energy contribution to $\mathcal{H}_{\text{FE/MD}}$, all $x$- and $y$-displacements in
the third dimension were replicated on the FE side of the boundary by assuming that
atoms are located at ideal lattice sites in that dimension. The overall Hamiltonian
remains conservative.

Two other issues regarding the definition of energy must still be addressed, both
of which involve a reference state. One involves the potential energy, while the
other has to do with the thermal energy. Consider first the potential energy. The
SW potential is referenced to infinitely separated atoms, whereas the FE potential
is referenced to a $T = 0$ unstrained lattice. Therefore, a constant offset energy that
did not affect the dynamics was added to each FE mesh point. The $T = 0$ energy
density for SW silicon at zero pressure is -4.33444 eV/atom. The offset energy was
computed for every FE point using an equation entirely analogous to that used to
compute mass in the "lumped-mass" approximation except that, instead of a mass
density, the SW energy density was used [see Eq. (43)]. This scheme ensured the
correct limiting behavior as the mesh spacing was reduced to atomic dimensions.
For atoms in the interface region, for systems with unusual orientation where
the offset is non-trivial to estimate atom by atom, a $T = 0$ calculation with zero
strain for the coupled FE/MD system can be performed. The offset may thereby
be calculated to maintain the energy/atom constant through the interface.

Rudd and Broughton (1998) showed that in the FE region $(\dot{\mathbf{u}})^2$ is related to the
temperature. However, as one moves away from the interface region and coarsens
the FE grid, atomic degrees of freedom are lost: The FE algorithm involves an
average over these degrees of freedom. Thus, to set the atomic and continuum
thermal energies on an equivalent level, the total FE thermal energy is written
again by using an offset. These corrected energies are denoted by a prime:

$$(E')'_{\text{FE}} = \frac{3}{2}(N_a - N_m)k_B T + (E_k)_{\text{FE}} + \frac{1}{2}N_m k_B T, \qquad (48)$$

$$(E')_{\text{FE}} = \frac{3}{2}(N_a - N_m)k_B T + E_{\text{FE}} + \frac{1}{2}N_m k_B T. \qquad (49)$$

Here, $N_a$ is the number of atoms contained within an equivalent 3D volume, $k_B$ is the Boltzmann's constant, and equipartition has been used. Broughton *et al.* further assumed that the background temperature does not vary during the simulation. Therefore, the first terms of Eqs. (48) and (49) account for the missing degrees of atomic freedom, while the last terms augment the 2D FE plane-strain simulation for the missing third dimension in its degrees of freedom. Similar to the case of potential energy, these offsets do not affect the dynamics of the system and the thermal corrections can be assigned to each mesh point in a manner similar to that described above for the zero-temperature FE potential energy. For finite-temperature simulations, the $\dot{\mathbf{u}}$ were thermalized according to the Maxwell–Boltzmann distribution [see Eq. (9.49)]. Moreover, the appropriate elastic constants for that temperature should be used in the FE equations of motion so as to make the MD and FE regions seamless and compatible. Further, since this method requires a continuation of ideal lattice sites into the FE/MD interface region so as to determine mesh coordinates, the appropriate lattice parameter for a given temperature should be used.

The last issue to be addressed is the dissipation in the FE region. The continuum representation of silicon used by Broughton *et al.* (1999) was based on linear elasticity theory, which is a harmonic theory. Thus, vibrational modes of given $\{k, \omega\}$ relationship, which depend upon the long-wavelength elastic constants of the medium, propagate undamped. In order to thermalize short-wavelength phonons propagating through regions where the mesh spacing changes, and also to allow energy to be dissipated in the FE region, Broughton *et al.* (1999) weakly coupled the FE degrees of freedom to a Brownian heat bath with dynamics that was set to the temperature at which the simulation was being performed, hence coupling the phonon modes of the FE region. The force used in the third step of the velocity Verlet algorithm included random and dissipative terms:

$$\mathbf{f}_{FE}^t = \frac{\partial E_{\text{FE}}}{\partial \mathbf{u}^t} + r(T, y^t) - \eta(y^t) M^t \dot{\mathbf{u}}^t , \qquad (50)$$

where $r$ is a Gaussian random variable, and $\eta$ is a friction coefficient. Through the fluctuation-dissipation theorem, the variance $\sigma$ of the Gaussian is related to $\eta$:

$$\sigma = \sqrt{\frac{2\eta M^t k_B T}{\Delta t}}. \qquad (51)$$

In order to minimally perturb the dynamics of the active zone (i.e., MD and TB), $\eta$ was assumed to be a function of the (time-invariant) FE mesh $y$-coordinate, and was linearly increased from zero in the interface region to a finite value of about 0.1 at the extremal outer edge of the FE regions.

## 10.3.6   *Interfacing Molecular Dynamics and Tight-Binding Regions*

In contrast to the algorithm for the FE/MD interface in which a plane between rows of atoms was defined, the MD/TB interface is conceptually across a plane that consists of atoms. This is necessitated due to the difficulty of assigning (localized) energy in a computationally efficient way to specific bonds in an electronic struc-

FIGURE 10.13. The MD/TB interface region. Full lines represent two- and three-body terms of the Stillinger–Weber potential that contribute fully to the interface Hamiltonian, while the broken ones do not contribute at all (after Broughton *et al.*, 1999).

ture calculation. The MD/TB interface is shown schematically in Figure 10.13. Since in silicon the covalent bonds are local objects, dangling bonds may be tied off with univalent atoms (UAs), and therefore the region chosen for TB description is terminated with such atoms. The **H** and **S** matrices and the repulsive pair potential $U_{rep}$ that couple the UAs to the silicon atoms within the interior of the TB region were chosen to, (a) maintain electro-neutrality on both the UAs and silicons; (b) locate the univalent atom potential energy minimum at a Si-Si distance, not a Si-H distance; (c) provide a bond energy equal to a single Si-Si bond, and (d) provide a longitudinal force constant equal to that of Si. At the perimeter of the TB region, the UAs are constrained to stay at the Si sites of the MD region which, in many cases, implies that more than one UA is found at a given site. Thus, there are no matrix elements, nor $U_{rep}$ terms that couple any of the UAs to one another. In the computations, a circle is drawn around an inner set of atoms and are designated TB silicons. Then, any atom outside this circle, but within range of an inner atom, is designated as a UA. The range that was used as the criterion was the average of the first ($r_0$) and second neighbor [$\sqrt{8/3}r_0$] distances of the equilibrium Si lattice. In Figure 10.13, matrix elements that couple atoms across the light gray region are of Si-Si form. Atoms coupled across the dark gray region use Si-UA matrix elements for which the necessary parameters are given by Bernstein and Kaxiras (1997).

The TB **H** and **S** generalized-eigenvalue problem is solved for the entire Si plus UA system. The only remaining issue is to determine which SW two- and three-body terms are required to couple the UAs to the MD region. Since these atoms are not coupled to one another in the TB region, SW terms that account for such are required. Broughton *et al.* (1999) included all SW pair terms between a UA and Si and either a Si atom in the MD region or another UA-Si, as well as all SW triplets that include at least one MD Si to one UA-Si pair. Thus, the forces that arise on the UA-Si from these terms are added to the forces arising from the TB Hamiltonian on these atoms. Equation (30) should more correctly be written as

$$\mathcal{H} = \mathcal{H}_{FE}(\{\mathbf{u}, \dot{\mathbf{u}}\} \in FE) + \mathcal{H}_{FE/MD}(\{\mathbf{u}, \dot{\mathbf{u}}, \mathbf{r}, \dot{\mathbf{r}}\} \in FE/MD)$$

$$+\mathcal{H}_{MD}(\{\mathbf{r}, \dot{\mathbf{r}}\} \in MD) + \mathcal{H}_{MD/TB}(\{\mathbf{r}, \dot{\mathbf{r}}\} \in MD/TB) + \mathcal{H}_{TB}(\{\mathbf{r}, \dot{\mathbf{r}}\} \in MD/TB),$$

$$(52)$$

FIGURE 10.14. Overlapping TB regions embedded in the MD region. Atomic forces are a function of the overlap (after Broughton *et al.*, 1999).

where the penultimate term involves only SW interactions crossing the boundary, while the last term involves a TB calculation for the combined Si-UA system.

The above prescription produces a conservative Hamiltonian, if there is no dynamic allocation of the TB region to those parts of the material where atomic bonds are breaking. Unfortunately, for many materials and phenomena (such as crack propagation) the 100 or so atoms, whose forces may currently be updated using a non-orthogonal TB Hamiltonian in one second of real time, do not comprise a large region. This problem may be partially addressed by using periodic boundary conditions. One may also address the problem by using more than one processor per TB region to perform the diagonalization, but such algorithms are presently not efficient on coarse-grained scalable architecture computers. Broughton *et al.* (1999) represented the region of breaking bonds by a region of TB segments, an example of which is shown in Figure 10.14 for three overlapping TB regions. In their simulations of fracture propagation (see below), Broughton *et al.* used eight overlapping regions, with each region diagonalized separately, and also handled by a separate processor. After forces on each atom are obtained for each TB region separately, the force to be used in updating the velocity Verlet algorithm is obtained as the average over the different regions; if there is no overlap of TB regions, the same prescription as for a single TB region is used; where a UA of one TB region overlaps a Si of another, the Si value is used. The number of atoms that are propagated using TB forces is therefore less than the total number within all the overlapping regions. These rules are, of course, intuitive.

In the simulation of fracture propagation, the energy and force algorithm, *as implemented* in the MD and TB regions, proceeds by calculating the SW energy for all atoms in the MD processors. The TB processors calculate not only TB energies and forces, but also those SW forces that must be subtracted from those double counted in MD processors. Thus, SW energies, by suitable division of two- and three-body terms, are available for all atoms, which are then used to discriminate different regions. The apex of a crack is found, for example, by locating the atom, with a potential energy greater (more positive) than 60% of the bulk cohesive potential energy, furthest to the left or right of the center of the system. The central TB portion of the overlap region is then placed at that atom.

### 10.3.7  Seamless Simulation

The description of the multiscale method indicates how the simulation is made seamless. The TB region, the region described at the smallest length scale, determines the elastic constants and the atomic force fields used elsewhere in the system. Therefore, the procedure used by Broughton *et al.* (1999) is as follows.

(1)  A pure TB simulation is carried out for a small number of atoms that represent the bulk material (at given temperature and pressure). By appropriate deformation of the computational cell, the elastic constants, and hence the force constant, are computed.
(2)  The SW parameters for Si are then adjusted to reproduce the same estimates.
(3)  The elastic constants from the TB region are also utilized for the stiffness matrix of the FE region.
(4)  Finally, the parameters used for the Si-univalent atom matrix elements are adjusted so that displacement of a univalent atom-Si in the coupled system gives rise to the same force constant of the pure bulk material.

### 10.3.8  Multiscale Simulation of Fracture Propagation in Silicon

The foregoing multiscale algorithm was utilized by Broughton *et al.* (1999) to study rapid brittle fracture of a Si slab "damaged" by a microcrack at its center and deformed under uniaxial tension. The MD region was spatially domain-decomposed onto 24 processors. Each FE region was handled by its own processor. The path of the fracture was monitored and the center of the TB region was placed at the apex of the crack where bond breaking occurs. This region is crucial to determining the kinetics of the crack propagation. A region of eight overlapping TB regions, each being a cylinder of radius 5.3 Å in the $yz$ plane and distributed to a different processor, was used. The exposed notch faces were $x - z$ planes with (100) faces, with the microcrack pointed in the (010) direction. Each FE region contained 258,048 mesh points, there were 1,032,192 atoms in the MD region, and around 280 unique atoms in the TB region. The lengths of the MD region were 10.9 Å (the slab thickness and periodic), 521 Å (before the pull, in the direction of the pull), and 3, 649 Å (the primary direction of propagation and periodic). The full pull length of the FE+MD system was 5, 602 Å. The entire system, including the FE, contained 11,093,376 atoms. The time for a TB force update was 1.5 s, 1.85 s for the MD update, and 0.7 s for the FE.

The rectilinear computational cell comprised (100) faces on all sides. The reduced elastic constant matrix for this geometry was obtained by averaging the results reported by Balamane *et al.* (1992) and Ray (1988) for the zero temperature $C_{11}$, $C_{12}$, and $C_{44}$ elastic constants of SW Si:

$$\mathbf{C} = \begin{pmatrix} 1.578 \times 10^6 & 0.7930 \times 10^6 & 0.0 \\ 0.7930 \times 10^6 & 1.578 \times 10^6 & 0.0 \\ 0.0 & 0.0 & 0.6365 \times 10^6 \end{pmatrix}, \qquad (53)$$

where the units are megabar.

The slab was initialially at zero temperature, and a constant strain rate was imposed on the outermost FE boundaries defining the opposing horizontal faces of the slab. Moreover, a linear velocity gradient was applied within the slab, which resulted in an increasing internal strain with time. The material failed at the notch tip after it had been stretched by about 1.5%. The propagating cracks rapidly achieved a limiting speed (2770 m/s) equal to 85% of the Rayleigh speed $c_R$, the sound speed of the solid Si surface. A very accurate signature of seamless coupling, one that represents a validation of the method, is the fact that stress waves passed from the MD region to the FE regions with no visible reflection at the FE/MD interface, i.e., the coupling of the MD region to the FE region appeared seamless. Moreover, there were no obvious discontinuities at the MD/TB interface.

## 10.4    Other Applications of Multiscale Modeling

In addition to the above examples, over the past few years the multiscale methodology has been utilized for describing several other important sets of phenomena. In this section, we briefly describe two of such applications. Complete details of these works are given in the original papers.

### 10.4.1    Atomistically Induced Stress Distributions in Composite Materials

With the rapid technological advances of the past decade, the semiconductor device feature size has been decreasing. The feature size is now predicted to reach 70 nm by the year 2008, if not sooner. Because materials of such sizes have a high surface-to-volume ratios, the stress inhomogeneities that are caused by surfaces and interfaces play an important role in the performance of the devices, since they affect the donor distribution by trapping donors in tensile stress regions. Thus, it is important to be able to control the growth of the materials in order for them to have the desired properties. This can be aided by accurate models that describe the materials and the stress and strain distributions in them. Although several continuum models have been developed for this purpose (see, for example, Johnson and Freund, 1997), they appear to be incapable of describing atomistically-induced stresses at the interface of two different materials, especially if one or both of the materials are amorphous. A better approach would be based on a multiscale model that can address the atomistic aspects of the problem near the interface, and the macroscopic aspects deep into the substrate.

Several of such multiscale models have been developed recently. Of particular note is the multiscale model that was developed by Lidorikis *et al.* (2001) who used a hybrid FE/MD approach to study stress domains in Si nanopixels covered with amorphous $Si_3N_4$ films. The stress domains originate from the atomic configuration at the amorphous/crystalline interface and extend throughout the substrate, where the semiconductor device operates. Lidorikis *et al.* (2001) simulated a $25 \times 25$ nm Si(111) mesa, covered with a $25 \times 25 \times 5$ nm $Si_3N_4$ film, on

a $50 \times 50 \times 15$ nm Si(111) substrate. The lattice mismatch between Si and $Si_3N_4$ is of the order of 1.2%. The $Si_3N_4$ film, as well as its interface with Si, was treated by the MD simulations, while most of the Si substrate was modeled by FE simulations. Both crystalline (001) and amorphous $Si_3N_4$ films were analyzed. Periodic boundary conditions were applied on the sides of the substrate, while its bottom was held fixed.

To model Si, the SW potential was used, while for $Si_3N_4$ a combination of two- and three-body interaction terms, that included electronic polarizability, charge transfer (taken to be screened Coulomb potential), and covalent bonding effects, was utilized. A variation of the same potential was also used for describing the interactions across the $Si/Si_3N_4$ interface. In the Si substrate, about 20 Å below the interface, the deformations are small and thus were treated by linear continuum elasticity, solved on a computational grid. The displacement field was discretized at the nodes of the grids and, within the grid cells, was interpolated from the nodal values.

The MD and FE regions were then merged seamlessly in the interface region between the two. To do this, the FE mesh was refined down to atomic scale, and was also shifted from the simple-cubic arrangement (used in the continuum region), in order to follow the lattice structure of Si. Within the interface region, the FE and crystalline lattices overlap, hence yielding a one-to-one correspondence between the MD atoms and the FE nodes. For these hybrid atom/node particles, an average Hamiltonian was defined, and the total Hamiltonian was given by an equation similar to (52) (except that, of course, there was no TB region). The Verlet velocity algorithm, described in Chapter 9, was used for integrating the equations of motion in the MD region.

The hybrid simulations for the Si pixels covered with an amorphous $Si_3N_4$ film indicated that, the inhomogeneous stress domains are formed below the interface with a shadow that extends deep into the substrate. These domains are caused by atomistically-induced stresses at the lattice-mismatched amorphous/crystalline interface, and cannot be predicted by a continuum model alone. The chemistry of the interface and the degree of mismatch at the interface control the size of the domains. These simulations were compared with those in which only the MD method with several million atoms were utilized; the agreement was found to be excellent, hence validating the multiscale model. The computational time of the multiscale model was of course much less than that of the MD simulation alone.

## 10.4.2   Chemical Vapor Deposition

Chemical vapor deposition (CVD) is a complex operation that involves several key processes which occur on widely separated length and time scales. At each of the distinct length scales, a particular set of modeling techniques and assumptions is appropriate. One reason for the complexity of CVD is that, frequently, processes of different length scales interact, and thus many problems of interest cannot be cast within the context of a single scale. An appropriate model of CVD should therefore focus on interactions between macroscopic flow and transport, where

the gas phase is regarded as a continuum, and phenomena at the micron length scale, where the discrete nature of the gas phase becomes apparent and important.

For most CVD processes, flow and transport in the reactor at large are described by the usual macroscopic conservation equations for momentum, mass, and energy which give rise to coupled nonlinear partial differential equations that are typically solved with the FE method. However, these continuum models are not valid as length scales approach and shrink below the gas phase mean-free path $\lambda$. This is an important consideration since microelectronic device features are frequently submicron, whereas the mean-free path $\lambda$, under low pressure CVD conditions, can be several hundred $\mu$m. This limitation on the use of continuum models is widely recognized and has led to the development of discrete particle transport models.

Rodgers and Jensen (1998) developed a multiscale model for CVD on length scales ranging from microns to meters. At the macroscopic level the problems of fluid flow and heat and mass transfer in a single wafer, low-pressure CVD reactor were solved using the FE method. Since on the feature scale the continuum models are no longer valid, transport at the feature scale was linked with the continuum model by an effective reactivity function $\mathcal{R}$ that included effects of both the multiscale surface heterogeneity and microscopic transport resistance. A Monte Carlo method was then used for computing $\mathcal{R}$ for any set of reaction pathways that occur over microelectronic device features of any geometry. Feature scale computations were then combined to yield an effective reactivity map over the surface of the substrate, which was subsequently utilized for formulating a flux boundary condition for the continuum model. Iteration between macroscopic and microscopic models was then used to ensure a consistent set of conditions at the micro-macro interface.

For a different application of multiscale modeling methodology to problems involving transport and reaction processes in nano- and microporous materials, see Dadvar and Sahimi (2002, 2003).

## Summary

The goal in multiscale modeling is to predict the performance and behavior of heterogeneous materials in which there are several relevant and widely disparate length scales. Due to the great complexities that are involved, the task of developing a multiscale model is usually referred to as a *grand challenge problem*. The complexities that are involved are clear: At the atomic scale electrons govern the interactions among atoms in a material, and thus quantum-mechanical effects are important and must be taken into account. At the same time, at the macroscopic or engineering scale, forces that arise from macroscopic stresses and/or temperature gradients are the factors that control the performance of materials. At the intermediate length scales, defects such as dislocations control the mechanical properties of materials up to tens of micrometers, while large collections of such defects, including grain boundaries, govern their mesoscopic properties up to length scales

that are of the order of hundreds of micrometers. Coupling of all of these length scales would not have been possible, had it not been for the tremendous advances in the advent of ever more powerful, massively-parallel computers, and the great advances that have been made in the theoretical understanding of materials and their properties. The emerging multiscale methodology demonstrates how coupling of atomistic and continuum approaches results in more predictive power than either approach offers alone.

An interesting and very useful aspect of multiscale modeling is the fact that, it is a multidisciplinary field that brings together scientists from many disciplines. The development of a multiscale model of any phenomenon that deals with one or more aspect of materials' properties should, in principle at least, involve chemists and chemical engineers, applied physicists and mathematicians, and continuum mechanicians.

# References

Abell, G.C., "Empirical chemical pseudopotential theory of molecular and metallic bonding," *Phys. Rev. B* **31**, 6184 (1985).

Abraham, F.F., "Computational statistical mechanics methodology, applications and supercomputing," *Adv. Phys.* **35**, 1 (1986).

Abraham, F.F., "Dynamics of brittle fracture with variable elasticity," *Phys. Rev. Lett.* **77**, 869 (1996).

Abraham, F.F., D. Brodbeck, R.A. Rafey, and W.E. Rudge, "Instability dynamics of fracture: A computer simulation investigation," *Phys. Rev. Lett.* **73**, 272 (1994).

Abraham, F.F., D. Brodbeck, W.E. Rudge, and X. Xu, "A molecular dynamics investigation of rapid fracture mechanics," *J. Mech. Phys. Solids* **45**, 1595 (1997a).

Abraham, F.F., J.Q. Broughton, N. Bernstein, and E. Kaxiras, "Spanning the length scales in dynamic simulation," *Comput. Phys.* **12**, 538 (1998a).

Abraham, F.F., J.Q. Broughton, N. Bernstein, and E. Kaxiras, "Spanning the continuum to quantum length scales in a dynamic simulation of brittle fracture," *Europhys. Lett.* **44**, 783 (1998b).

Abraham, F.F., and H. Gao, "How fast can cracks propagate?" *Phys. Rev. Lett.* **84**, 3113 (2000).

Abraham, F.F., W.E. Rudge, D.J. Auerbach, and S.W. Koch, "Molecular dynamics simulations of the incommensurate phase of Krypton on graphite using more than 100,000 atoms," *Phys. Rev. Lett.* **52**, 445 (1984).

Abraham, F.F., D. Schneider, B. Land, D. Lifka, J. Skovira, J. Gerner, and M. Rosenkrantz, "Instability dynamics in three-dimensional fracture: An atomistic simulation," *J. Mech. Phys. Solids* **45**, 1461 (1997b).

Adda-Bedia, M., and M. Ben Amar, "Stability of Quasiequillibrium cracks under uniaxial loading," *Phys. Rev. Lett.* **76**, 1497 (1996).

Adda-Bedia, M., and Y. Pomeau, "Crack instabilities in a heated glass strip," *Phys. Rev. E* **52**, 4105 (1995).

Aharony, A., "Universal critical amplitude ratios for percolation," *Phys. Rev. B* **22**, 400 (1980).

Aharony, A., "Crossover from linear to nonlinear resistance near percolation," *Phys. Rev. Lett.* **58**, 2726 (1987).

Aktsipetrov, A.A., O. Keller, K. Pedersen, A.A. Nikulin, N.N. Novikova, and A.A. Fedyanin, "Surface-enhanced second-harmonic generation in $C_{60}$-coated silver island films," *Phys. Lett. A* **179**, 149 (1993).

Alder, B.J., and T.E. Wainwright, "Phase transitions for a hard sphere system," *J. Chem. Phys.* **27**, 1208 (1957).

Alder, B.J., and T.E. Wainwright, "Decay of velocity autocorrelation function," *Phys. Rev. A* **1**, 18 (1969).

Allen, M.P., and D.J. Tildesley, *Computer Simulation of Liquids* (Oxford University Press, London, 1987).

Andersen, H.C., "Molecular dynamics simulations at constant pressure and/or temperature," *J. Chem. Phys.* **72**, 2384 (1980).

Andersen, J.V., D. Sornette, and K.-t. Leung, "Tricritical behavior in rupture induced by disorder," *Phys. Rev. Lett.* **78**, 2140 (1997).

Andrews, D.J., "Rupture velocity of plane strain shear cracks," *J. Geophys. Res.* **81**, 5679 (1976).

Anthony, J.B.C.S.R., J.B. Chubb, and J. Congleton, "The crack-branching velocity," *Philos. Mag.* **22**, 1201 (1970).

Arakawa, K., and K. Takahashi, "Branching of a fast crack in polymers," *Int. J. Fract.* **48**, 245 (1991).

Aranson, I.S., V.A. Kalatsky, and V.M. Vinokur, "Continuum field description of crack propagation," *Phys. Rev. Lett.* **85**, 118 (2000).

Arbabi, S., and M. Sahimi, "Elastic properties of three-dimensional percolation networks with stretching and bond-bending forces," *Phys. Rev. B* **38**, 7173 (1988).

Arbabi, S., and M. Sahimi, "On three-dimensional elastic percolation networks with bond-bending forces," *J. Phys. A* **23**, 2211 (1990a).

Arbabi, S., and M. Sahimi, "Test of universality for three-dimensional models of mechanical breakdown in disordered solids," *Phys. Rev. B* **41**, 772 (1990b).

Archuleta, R.J., "Analysis og near-source static and dynamic measurements from the 1979 Imperial Valley earthquake," *Bull. Seismol. Soc. Am.* **72**, 1927 (1982).

Argon, A.S., "Brittle to ductile transition in cleavage fracture," *Acta Mettal.* **35**, 185 (1987).

Arias, T.A., and J.D. Joannopoulos, "*Ab initio* theory of dislocation interactions: From close-range spontaneous annihilation to the long-range continuum limit," *Phys. Rev. Lett.* **73**, 680 (1994).

Ashcroft, N., and N.D. Mermin, *Solid State Physics* (Saunders, London, 1976).

Ashurst, W.T., and W.G. Hoover, "Microscopic fracture studies in the two-dimensional triangular lattice," *Phys. Rev. B* **14**, 1465 (1976).

Åström, J.A., J.P. Mäkinen, M.J. Alva, and J. Timonen, "Elasticity of Poissonian fiber networks," *Phys. Rev. E* **61**, 5550 (2000).

Åström, J.A., S. Saarinen, K.J. Niskanen, and J. Kurkijärvi, "Microscopic mechanics of fiber networks," *J. Appl. Phys.* **75**, 2383 (1994).

Åström, J.A., and J. Timonen, "Crack bifurcations in a strained lattice," *Phys. Rev. B* **54**, R9585 (1996).

Åström, J.A., and J. Timonen, "Fragmentation by crack branching," *Phys. Rev. Lett.* **78**, 3677 (1997a).

Åström, J.A., and J. Timonen, "Fracture of a brittle membrane," *Phys. Rev. Lett.* **79**, 3684 (1997b).

Auerbach, D.J., W. Paul, A.F. Bakker, C. Lutz, W.E. Rudge, and F.F. Abraham, "A special purpose parallel computer for molecular dynamics: Motivation, design, implementation, and application," *J. Phys. Chem.* **91**, 4881 (1987).

Baker, T.C., and F.W. Preston, "The effect of water on the strength of glasses," *J. Appl. Phys.* **17**, 179 (1946).

Bakker, A.F., G.H. Gilmer, M.H. Grabow, and K. Thompson, "A special purpose computer for molecular dynamics calculations," *J. Comput. Phys.* **90**, 313 (1990).

Balamane, H., T. Halicioglu and W. A. Tiller, "Comparative study of silicon empirical interatomic potentials," *Phys. Rev. B* **46**, 2250 (1992).

Ball, R.C., and H. Larralde, "Three-dimensional stability analysis of planar straight cracks propagating quasistatically under type I loading," *Int. J. Fract.* **71**, 365 (1995).

Bao, G., J.W. Hutchinson, and R.M. McMeeking, "Particle reinforcement of ductile matrices against plastic flow and creep," *Acta Metall. Mater.* **39**, 1871 (1991).

Barabási, A.-L., and H.E. Stanley, *Fractal Concepts in Surface Growth* (Cambridge University Press, Cambridge, 1995).

Barabási, A.-L., and T. Vicsek, "Multifractality of self-affine surfaces," *Phys. Rev. A* **44**, 2730 (1990).

Baran, G.R., C. Roques-Carmes, D. Wehbi, and M. Degrange, "Fractal characteristics of fracture surfaces," *J. Am. Ceram. Soc.* **75**, 2687 (1992).

Barber, M., J. Donley, and J.S. Langer, "Steady-state propagation of a crack in a viscoelastic strip," *Phys. Rev. A* **40**, 366 (1989).

Barbosa, F.F., and S.L.A. de Queiroz, "Concentration anisotropy and directionality in the dielectric breakdown problem on a square lattice," *J. Phys.: Condens. Matter* **1**, 2771 (1989).

Barclay, A.L., P.J.J. Sweeney, L.A. Dissado, and G.C. Stevens, "Stochastic modelling of electrical treeing: fractal and statistical characteristics," *J. Phys. D* **23**, 1536 (1990).

Bardhan, K.K., "Nonlinear conduction in composites above percolation threshold-beyond the backbone," *Physica A* **241**, 267 (1997).

Barenblatt, G.I., "The formation of equilibrium cracks during brittle fracture: general ideas and hypothesis, axially symmetric cracks," *Appl. Math. Mech.* (Translation of PMM) **23**, 622 (1959a).

Barenblatt, G.I., "Concerning equilibrium cracks forming during battle fracture: the stability of isolated cracks, relationship with energetic theories," *Appl. Math. Mech.* (Translation of PMM) **23**, 1273 (1959b).

Barton, C.C., "Fractal analysis of the spatial clustering of fractures," in *Fractals and Their Use in the Earth Sciences*, edited by C.C. Barton and P.R. La Pointe (American Geophysical Union, 1992), p. 126.

Barton, C.C., and P.A. Hsieh, *Physical and Hydrological-Flow Properties of Fractures*, Guidebook T385 (American Geophysical Union, Las Vegas, Nevada, 1989).

Barton, C.C., and E. Larsen, "Fractal geometry of two-dimensional fracture networks at Yucca Mountain, southwest Nevada," in *Proceedings of the International Symposium on Fundamentals of Rock Joints*, edited by O. Stephenson (Bjorkliden, Sweden, 1985), p. 77.

Barton, C.C., T.A. Schutter, W.R. Page, and J.K. Samuel, "Computer generation of fracture networks for hydrologic-flow modelling," *Trans. Amer. Geophys. Union* **68**, 1295 (1987).

Barton, J.L., "La relaxation diélectrique de quelques verres ternaires silice oxyde alcalin oxyde alcalin-terreux," *Verres et Refr.* **20**, 328 (1966).

Baskes, M.I., J.S. Nelson, and A.F. Wright, "Semiempirical modified embedded-atom potentials for silicon and germanium," *Phys. Rev. B* **40**, 6085 (1989).

Batrouni, G.G., and A. Hansen, "Fracture in three-dimensional fuse model," *Phys. Rev. Lett.* **80**, 325 (1998).

Bazant, M.Z., E. Kaxiras, and J.F. Justo, "Environment-dependent interatomic potential for bulk silicon," *Phys. Rev B* **56**, 8542, (1997).

Beale, P.D., and P.M. Duxbury, "Theory of dielectric breakdown in metal-loaded dielectrics," *Phys. Rev. B* **37**, 2785 (1988).

Beale, P.D., and D.J. Srolovitz, "Elastic fracture in random materials," *Phys. Rev. B* **37**, 5500 (1988).

Beazley, D.M., and P.S. Lomdahl, "Message-passing multi-cell molecular dynamics on the Connection Machine CM-5," *Parallel Computing* **20**, 173 (1994).

Beazley, D.M., P.S. Lomdahl, N. Gr$\phi$nbech-Jensen, R. Giles, and P. Tamayo, "Parallel algorithms for short-range molecular dynamics," *Annual Rev. Comput. Phys.* **3**, 116 (1995).

Benguigui, L., "Simulation of dielectric failure by means of resistor-diode random lattices," *Phys. Rev. B* **38**, 7211 (1988).

Benguigui, L., and P. Ron, "Laboratory simulation of dielectric breakdown," in *Non-Linearity and Breakdown in Soft Condensed Matter*, edited by K.K. Bardhan, B.K. Chakrabarti, and A. Hansen, *Lecture Notes in Physics* **437** (Springer, Heidelberg, 1994), p. 221.

Benguigui, L., P. Ron, and D.J. Bergman, "Strain and stress at the fracture of percolative media," *J. Physique* **48**, 1547 (1987).

Bennett, C.H., "Efficient estimation of free energy differences from Monte Carlo data," *J. Comput. Phys.* **22**, 245 (1976).

Beran, M.J., "Use of the variational approach to determine bounds for the effective permittivity of random media," *Nuovo Cimento* **38**, 771 (1965).

Beran, M.J., *Statistical Continuum Theories* (Wiley, New York, 1968).

Bergkvist, H., "Some experiments on crack motion and arrest in polyethylmethacrylate," *Eng. Fract. Mech.* **6**, 621 (1974).

Bergman, D.J., in *Fragmentation, Form and Flow in Fractured Media*, *Ann. Israel Phys. Soc.* **8**, 266 (1986).

Bergamn, D.J., "Nonlinear behavior and $1/f$ noise near a conductivity threshold: effects of local microgeometry," *Phys. Rev. B* **39**, 4598 (1989).

Bergman, D.J., and Y. Kantor, "Critical properties of an elastic fractal," *Phys. Rev. Lett.* **53**, 511 (1984).

Bergman, D.J., O. Levy, and D. Stroud, "Theory of optical bistability in a weakly nonlinear composite medium," *Phys. Rev. B* **49**, 129 (1994).

Bernstein, N., and E. Kaxiras, "Nonorthogonal tight-binding Hamiltonians for defects and interfaces in silicon," *Phys. Rev. B* **56**, (1997).

Berry, M.V., "Regular and irregular semiclassical wavefunctions," *J. Phys. A* **10**, 2083 (1977).

Berveiller, M., and A. Zaoui, "An extension of the self-consistent scheme to plastically-flowing polycrystals," *J. Mech. Phys. Solids* **26**, 325 (1979).

Bhattacharya, K., and R.V. Kohn, "Elastic energy minimization and recoverable strains of polycrystalline shape-memory materials," *Arch. Rational Mech. Anal.*, (1997).

Bird, R.B., R.C. Armstrong, and O. Hasseger, *Dynamics of Polymeric Liquids*, 2nd ed. (Wiley, New York, 1987).

Bishop, J.F.W., and R. Hill, "A theory of the plastic distortion of a polycrystalline under combined stresses," *Philos. Mag.* **42**, 414 (1951a).

Bishop, J.F.W., and R. Hill, "A theoretical derivation of the plastic properties of a polycrystalline face-center metal," *Philos. Mag.* **42**, 1298 (1951b).

Blumberg Selinger, R.L., Z.-G. Wang, W.M. Gelbart, and A. Ben-Shaul, "Statistical-theormodynamic approach to fracture," *Phys. Rev. A* **43**, 4396 (1991a).

Blumberg Selinger, R.L., Z.-G. Wang, and W.M. Gelbart, "Effect of temperature and small-scale defects on the strength of solids," *J. Chem. Phys.* **95**, 9128 (1991b).

Blumenfeld, R., and A. Aharony, "Nonlinear resistor fractal networks, topological distances, singly connected bonds and fluctuations," *J. Phys. A* **18**, L443 (1985).

Blumenfeld, R., and D.J. Bergman, "Comment on Nonlinear susceptibilities of granular matter," *Phys. Rev. B* **43**, 13682 (1991a).

Blumenfeld, R., and D.J. Bergman, "Strongly nonlinear composite dielectrics: a perturbation method for finding the potential field and bulk effective properties," *Phys. Rev. B* **44**, 7378 (1991b).

Blumenfeld, R., Y. Meir, A. B. Harris, and A. Aharony, "Infinite set of exponents describing physics on fractal networks," *J. Phys. A* **19**, L791 (1986).

Bolding, B.C., and H.C. Andersen, "Interatomic potential for silicon clusters, crystals and surfaces," *Phys. Rev. B* **41**, 10568 (1989).

Born, M., "Thermodynamics of crystals and melting," *J. Chem. Phys.* **7**, 591 (1939).

Bornert, M. *Morphologie Microstructurale et Comportement Mécanique: Caractérisations Expérimentales, Approches par Bornes et Estimations Autocohérentes Généralisées.*, Ph.D. Thesis, Ecole Nationale des Ponts et Chaussées (1996).

Bornert, M., C. Stolz, and A. Zaoui, "Morphologically representative pattern-based bounding in elasticity," *J. Mech. Phys. Solids* **44**, 307 (1996).

Bouchaud, E., and J.-P. Bouchaud, "Fracture surfaces: apparent roughness, relevant length scales, and fracture toughness," *Phys. Rev. B* **50**, 17752 (1994).

Bouchaud, J.-P., E. Bouchaud, G. Lapasset, and J. Planes, "Models of fractal cracks," *Phys. Rev. Lett.* **71**, 2240 (1993).

Bouchaud, E., J.-P. Bouchaud, J. Planes, and G. Lapasset, "The statistics of crack branching during fast crack propagation," *Fractals* **1**, 1051 (1993a).

Bouchaud, E., L. de Arcangelis, G. Lapasset, and J. Planes, "Fractals breakage mater," *La Recherche* (Paris) **22**, 808 (1991).

Bouchaud, E., G. Lapasset, and J. Planes, "Fractal dimension of fractured surfaces: A universal value?" *Europhys. Lett.* **13**, 73 (1990).

Bouchaud, E., G. Lapasset, J. Planes, and S. Navéos, "Statistics of branched fracture surfaces," *Phys. Rev. B* **48**, 2917 (1993b).

Bouchaud, E., and S. Navéos, "From quasi-static to rapid fracture," *J. Phys. I France* **5**, 547 (1995).

Bouchitte, G., and P. Suquet, "Homogenization, plasticity and yield design," in: *Composite Media and Homogenization Theory*, edited by G. Dal Maso and G. Dell'Antonio (Birkhäuser, Basel, 1991), p. 107.

Boudet, J.F., S. Ciliberto, and V. Steinberg, "Experimental study of the instability of crack propagation in brittle materials," *Europhys. Lett.* **30**, 337 (1995).

Boudet, J.F., S. Ciliberto, and V. Steinberg, "Dynamics of crack propagation in brittle materials," *J. Phys II France* **6**, 1493 (1996).

Bowman, D.R., and D. Stroud, "Model for dielectric breakdown in metal-insulator composites," *Phys. Rev. B* **40**, 4641 (1989).

Boyd, R.W., *Nonlinear Optics* (Academic Press, New York, 1992).

Brace, W.S., and A.S. Orange, "Further studies of the effects of pressure on electrical resistivity of rocks," *J. Geophys. Res.* **73**, 1433 (1968).

Bradley, R.M., and K. Wu, "Studies of titanium dioxide film growth from titanium tetraisopropoxide," *J. Phys. A* **26**, 327 (1993).

Brańka, A.C., and K.W. Wojciechowski, "Generalization of Nosé and Nosé-Hoover isothermal dynamics," *Phys. Rev. E* **62**, 3281 (2000).

Brass, A., H.H. Jensen, and A.J. Berlinsky, "Models of flux pinning in the quasistatic limit", *Phys. Rev. B* **39**,102 (1989).

Brenner, D.W., "Empirical potential for hydrocarbons for use in simulating the chemical vapor deposition of diamond films," *Phys. Rev. B* **42**, 9458 (1990).

Brickstad, B., and F. Nilsson, "Numerical evaluation by FEM of crack propagation experiments," *Int. J. Fract.* **16**, 71 (1980).

Broberg, K.B., "On the behaviour of the process region at a fast running crack tip," in *High Velocity Deformation of Solids*, edited by K. Kawata and J. Shiori (Springer, Berlin, 1979), pp. 182–193.

Broberg, K.B., "The near-tip field at high crack velocities," *Int. J. Fract.* **39**, 1 (1989).

Brockenborough, J., S. Suresh, and H. Wienecke, "Deformation of metal-matrix composites with continuous fibers: Geometrical effects of fiber distribution and shape," *Acta Metall. Mater.* **39**, 735 (1991).

Brouers, F., S. Blacher, A.N. Lagarkov, A.K. Sarychev, P. Gadenne, and V.M. Shalaev, "Theory of giant Raman scattering from semicontinuous metal films," *Phys. Rev. B* **55**, 13234 (1997).

Brouers, F., S. Blacher, and A.K. Sarychev, "Giant field fluctuations and anomalous light scattering from semicontinuous metal films," *Phys. Rev. B* **58**, 15897 (1998).

Brouers, F., J.P. Clerc, G. Giraud, J.M. Laugier and Z.A. Randriamantany, "Dielectric and optical properties close to the percolation threshold. II," *Phys. Rev. B* **47**, 666 (1993).

Broughton, J.Q., F.F. Abraham, N. Bernstein, and E. Kaxiras, "Cocurrent coupling of length scales: Methodology and application," *Phys. Rev. B* **60**, 2391 (1999).

Bruge, F., and S.L. Fornili, "A distributed dynamic load balancer and its implementation on multi-transputer systems for molecular dynamics simulation," *Comput. Phys. Commun.* **60**, 39 (1990).

Budiansky, B., "A reassessment of deformation theories of plasticity," *J. Appl. Mech.* **26**, 259 (1959).

Budiansky, B., "On the elastic moduli of some heterogeneous materials," *J. Mech. Phys. Solids* **13**, 223 (1965).

Bulatov, V.V., S. Yip, and A.S. Argon, "Atomic models of dislocation mobility in silicon," *Philos. Mag.* **72**, 452 (1995).

Buresch, F.E., K. Frye, and T. Muller, in *Fracture Mechanics of Ceramics*, Vol. 5, edited by R.C. Bradt, A.G. Evans, D.P.H. Hasselman, and F.F. Lange (Plenum, New York, 1983), p. 591.

Burns, S.J., and W.W. Webb, "Fracture surface energies and dislocation processes during dynamical cleavage of LiF. II. Experiments," *J. Appl. Phys.* **41**, 2086 (1970).

Burridge, R., G. Conn, and L.B. Freund, "The stability of a rapid mode II shear crack with finite cohesive traction," *J. Geophys. Res.* **84**, 2210 (1979).

Cagin, T., and B.M. Pettitt, "Molecular dynamics with a variable number of molecules," *Mol. Simul.* **6**, 5 (1991).

Caldarelli, G., C. Castellano, and A. Vespignani, "Fractal and topological properties of directed fractures," *Phys. Rev. E* **49**, 2673 (1994).

Car, R., and M. Parrinello, "Unified approach for molecular dynamics and density-functional theory," *Phys. Rev. Lett.* **55**, 2471 (1985).

Carlsson, J., L. Dahlberg, and F. Nilsson, "Experimental studies of the unstable phase of crack propagation in metal and polymers," in *Dynamic Crack Propagation*, edited by G.C. Sih (Noordhoof, Leydon 1972), p. 3.

Ceperley, D.M., "Ground state of fermion one-component plasma—A Monte Carlo study in two and three dimensions," *Phys. Rev. B* **18**, 3126 (1978).

Chakrabarti, B.K., K.K. Bardhan, and R. Ray, "Insulation breakdown in a random nonpercolating network of conductors," *J. Phys. C* **20**, L57 (1987).

Chakrabarti, B.K., and L.G. Benguigui, *Statistical Physics of Fracture and Breakdown in Disordered Systems* (Oxford University Press, London, 1997).

Chakrabarti, B.K., D. Chowdhury, and D. Stauffer, "Molecular dynamic study of fracture in 2D disordered elastic Lennard–Jones solids," *Z. Phys. B* **62**, 343 (1986).

Chakrabarti, B.K., A.K. Roy, and S.S. Manna, "Breakdown exponents in lattice and continuum percolation," *J. Phys. C* **21**, L65 (1988).

Chan, D.Y.C., B.D. Hughes, L. Paterson, and C. Sirakoff, "Resistor networks with distributed breakdown voltages," *Phys. Rev. A* **43**, 2905 (1991).

Chan, K.S., (ed.), *George R. Irwin Symposium on Cleavage Fracture* (The Minerals, Metals & Materials Society, New York, 1997).

Chayes, J.T., L. Chayes, and R. Durret, "Critical behavior of the two-dimensional first passage time," *J. Stat. Phys.* **45**, 933 (1986).

Che, J., T. Cagin, and W.A. Goddard III, "Generalized extended empirical bond-order dependent force fields including nonbond interactions," *Theor. Chem. Acc.* **102**, 346 (1999).

Chelikowsky, J.R., J.C. Phillips, M. Kamal, and M. Strauss, "Surface and thermodynamic interatomic force fields for silicon clusters and bulk phases," *Phys. Rev. Lett.* **62**, 292 (1989).

Chermant, J.L., and M. Coster, "Quantitative fractography," *J. Mater. Sci.* **14**, 509 (1979).

Cheung, K.S., A.S. Argon, and S. Yip, "Activation analysis of dislocation nucleation from crack tip in $\alpha$-Fe," *J. Appl. Phys.* **69**, 2088 (1991).

Cheung, K.S., and S. Yip, "Brittle-ductile transition in intrinsic fracture behavior of crystals," *Phys. Rev. Lett.* **65**, 2804 (1990).

Ching, E.S.C., "Dynamic stresses at a moving crack tip in a model of fracture propagation," *Phys. Rev. E* **49**, 3382 (1994).

Ching, E.S.C., J.S. Langer, and H. Nakanishi, "Linear stability analysis for propagating fracture," *Phys. Rev. E* **53**, 2864 (1996a).

Ching, E.S.C., J.S. Langer, and H. Nakanishi, "Dynamic instabilities in fracture," *Phys. Rev. Lett.* **76**, 1087 (1996b).

Ching, E.S.C., J.S. Langer, and H. Nakanishi, "Model study of fracture propagation-solutions of steady-state propagation and their stability," *Physica A* **221**, 134 (1996c).

Cho, K., and J.D. Joannopoulos, "Ergodicity and dynamical properties of constant-temperature molecular dynamics," *Phys. Rev. A* **45**, 7089 (1992).

Christman, T., A. Needleman, and S. Suresh, "An experimental and numerical study of deformation in metal-ceramic composites," *Acta Metall. Mater.* **37**, 3029 (1989).

Chu, T., and Z. Hashin, "Plastic behavior of composites and porous media under isotropic stress," *Inter. J. Eng. Sci.* **9**, 971 (1971).

Chudnovsky, A., and B. Kunin, "A probabilistic model of brittle crack formation," *J. Appl. Phys.* **62**, 4124 (1987).

Chung, J.W., J.Th.M. De Hosson, and E. van der Giessen, "Failure of a disordered three-dimensional spring network," *Phys. Rev. B* **64**, 064202 (2001).

Ciccotti, G., M. Ferrario, and J.P. Rykaert, "Molecular dynamics of rigid systems in cartesian coordinates. A general formulation," *Mol. Phys.* **47**, 1253 (1982).

Cieplak, M., A. Maritan, and J.R. Banavar, "Optimal paths and domain walls in the strong disorder limit," *Phys. Rev. Lett.* **72**, 2320 (1994).

Cieplak, M., A. Maritan, and J.R. Banavar, "Invasion percolation and Eden growth: geometry and universality," *Phys. Rev. Lett.* **76**, 3754 (1996).

Clarke, L.J., I. Stich, and M.C. Payne, "Large-scale ab initio total energy calculations on parallel computers," *Comp. Phys. Comm.* **72**, 14 (1992).

Clementi, E., "Global scientific and engineering simulations on scalar, vector and parallel LCAP-type supercomputer," *Philos. Trans. R. Soc. Lond. A* **326**, 445 (1988).

Coleman, B.D., "On the strength of classical fibres and fibre bundles," *J. Mech. Phys. Solids* **7**, 60 (1958).

Coppard, R.W., L.A. Dissado, S.M. Rowland, and R. Rakowski, "Dielectric breakdown in metal-loaded polyethylene," *J. Phys.: Condens. Matter* **1**, 3041 (1989).

Cotterell, B., "Velocity effects in fracture propagation," *Appl. Mater. Res.* **4**, 227 (1965).

Cotterell, B., and A.G. Atkins, "A review of the *J* and *I* integrals and their implications for crack growth resistance and toughness in ductile fracture," *Int. J. Fract.* **81**, 357 (1996).

Cotterell, B., and J.R. Rice, "Slightly curved or kinked cracks," *Int. J. Fract.* **14**, 155 (1980).

Cox, H.L., "The elasticity and strength of paper and other fibrous materials," *Br. J. Appl. Phys.* **3**, 72 (1952).

Cracknell, R.F., D. Nicholson, and N. Quirke, "Direct molecular dynamics simulation of flow down a chemical potential gradient in a slit-shaped micropore," *Phys. Rev. Lett.* **74**, 2463 (1995).

Crosby, K.M., and R.M. Bradley, "Fragmentation of thin films bonded to solid substrates: simulation and a mean field theory," *Phys. Rev. E* **55**, 6084 (1997).

Curtin, W.A., "Theory of mechanical properties of ceramic-matrix composites," *J. Am. Ceram. Soc.* **74**, 2837 (1991).

Curtin, W.A., "Toughening in disordered brittle materials," *Phys. Rev. B* **55**, 11270 (1997).

Curtin, W.A., M. Pamel, and H. Scher, "Time-dependent damage evolution and failure in materials. II. Simulations," *Phys. Rev. B* **55**, 12051 (1997).

Curtin, W.A., and H. Scher, "Brittle fracture in disordered materials: A spring network model," *J. Mater. Res.* **5**, 535 (1990a).

Curtin, W.A., and H. Scher, "Mechanics modeling using a spring network," *J. Mater. Res.* **5**, 554 (1990b).

Curtin, W.A., and H. Scher, "Analytic model for scaling of breakdown," *Phys. Rev. Lett.* **67**, 2457 (1991).

Curtin, W.A., and H. Scher, "Algebraic scaling of material strength," *Phys. Rev. B* **45**, 2620 (1992).

Curtin, W.A., and H. Scher, "Time-dependent damage evolution and failure in materials. I. Theory," *Phys. Rev. B* **55**, 12038 (1997).

Dadvar, M., and M. Sahimi, "Pore network model of deactivation of immobilized glucose isomerase in packed-bed reactors II: three-dimensional simulation at the particle level," *Chem. Eng. Sci.* **57**, 939 (2002).

Dadvar, M., and M. Sahimi, "Pore network model of deactivation of immobilized glucose isomerase in packed-bed reactors III: simulation at the reactor level," *Chem. Eng. Sci.* (2003).

Daguier, P., E. Bouchaud, and G. Lapasset, "Roughness of a crack front pinned by microstructural obstacles," *Europhys. Lett.* **31**, 367 (1995).

Daguier, P., S. Henaux, E. Bouchaud, and F. Creuzet, "Quantitative analysis of a fracture surface by atomic force microscopy," *Phys. Rev. E* **53**, 5637 (1996).

Daguier, P., B. Nghiem, E. Bouchaud, and F. Creuzet, "Pinning and depinning of crack fronts in heterogeneous materials," *Phys. Rev. Lett.* **78**, 1062 (1997).

Dally, J.W., "Dynamic photoelastic studies of fracture," *Exp. Mech.* **19**, 349 (1979).

Dally, J.W., "Dynamic photoelasticity and its application to stress wave propagation, fracture mechanics, and fracture control," in *Static and Dynamic Photoelasticity and Caustics*, edited by A. Lagarde (Springer, Berlin, 1987), p. 247.

Dally, J.W., W.L. Fourney, and G.R. Irwin, "On the uniqueness of the stress intensity factor-crack velocity relationship," *Int. F. Fract.* **27**, 169 (1985).

Daniels, H.E., "The statistical theory of the strength of bundles of threads. I.," *Proc. R. Soc. Lond. A* **183**, 404 (1945).

Dauskardt, R., F. Haubensak, and R.O. Ritchie, "On the interpretation of the fractal character of fracture surfaces," *Acta Metall. Mater.* **38**, 143 (1990).

Daw, M.S., and M.I. Baskes, "Semiempirical, quantum mechanical calculation of hydrogen embrittlement in metals," *Phys. Rev. Lett.* **50**, 1285 (1983).

Daw, M.S., and M.I. Baskes, "Embedded-atom method: Derivation and application to impurities, surfaces, and other defects in metals," *Phys. Rev. B* **29**, 6443 (1984).

de Arcangelis, L., A. Hansen, H.J. Herrmann, and S. Roux, "Scaling laws in fracture," *Phys. Rev. B* **40**, 877 (1989).

de Arcangelis, L., and H.J. Herrmann, "Scaling and multiscaling laws in random fuse networks," *Phys. Rev. B* **39**, 2678 (1989).

de Arcangelis, L., S. Redner and H.J. Herrmann, "A random fuse model for breaking processes," *J. Physique* **46**, L585 (1985).

deBotton, G., "The effective yield strength of fiber-reinforced composites," *lnter. J. Solids Structures* **32**, 1743 (1995).

deBotton, G., and P. Ponte Castañeda, "On the ductility of laminated materials," *Inter. J. Solids Structures* **29**, 2329 (1992).

deBotton, G., and P. Ponte Castañeda, "Elastoplastic constitutive relations for fiber-reinforced solids," *Inter. J. Solids Structures* **30**, 1865 (1993).

deBotton, G., and P. Ponte Castañeda, "Variational estimates for the creep behavior of polycrystals," *Proc. R. Soc. Lond. A* **448**, 121 (1995).

de Buhan, P., and A. Taliercio, "A homogenization approach to the yield strength of composites," *Europ. J. Mech. A Solids* **10**, 129 (1991).

Dickinson, J.T., "Fracto-emission," in *Non-Destructive Testing of Fibre-Reinforced Plastics Composites*, edited by J. Summerscales, Vol. 2 (Elsevier, London, 1991).

Dienes, G.J., and A. Paskin, "Molecular dynamic simulations of crack propagation," *J. Phys. Chem. Solids* **48**, 1015 (1987).

Ding, H.-Q., N. Karasawa, and W.A. Goddard III, "Atomic level simulations on a million particles: the cell multipole method for Coulomb and London nonbond interactions," *J. Chem. Phys.* **97**, 4309 (1992).

Dissado, L.A., and J.C. Fothergill, *Electrical Degradation and Breakdown in Polymers* (Peregrinus, London, 1992).

Dissado, L.A., and P.J.J. Sweeney, "Physical model for breakdown structures in solid dielectrics," *Phys. Rev. B* **48**, 16261 (1993).

Dobrodumov, A.M., and A.M. El'yashevich, "Simulation of brittle fracture of polymers by a network model in the Monte Carlo method," *Sov. Phys.-Solid State* **15**, 1259 (1973).

Döll, W., "An experimental study of the heat generated in the plastic region of a running crack in different polymenc materials," *Eng. Fract. Mech.* **5**, 259 (1973).

Döll, W., "A molecular weight dependent transition in polymethylmethacrylate," *J. Mater. Sci.* **10**, 935 (1975).

Doyle, W.T., "The Clausius–Mossotti problem for cubic arrays of spheres," *J. Appl. Phys.* **49**, 795 (1978).

Dreizler, R.M., and E.K.U. Gross, *Density-Functional Theory* (Springer, Berlin, 1990).

Druker, D.C., "On minimum weight design and strength of non-homogeneous plastic bodies," in *Non-homogeneity in Elasticity and Plasticity*, edited by W. Olszag (Pergamon Press, New York, 1959), p. 139.

Drucker, D.C., "The continuum theory of plasticity on the macroscale and the microscale," *J. Mater.* **1**, 873 (1966).

Duarte, J.A.M.S., "The histogram characteristics of perimeter polynomials for directed percolation," *J. Physique* **47**, 383 (1986).

Duarte, J.A.M.S., "Direct lattices: site-to-bond conversion and its uses for the percolation and dilute polymer transitions," *Z. Phys. B* **80**, 299 (1990).

Duarte, J.A.M.S., "Series and Monte Carlo studies of 2 and 3 dimensions for axial hyperscaling in directed percolation," *Physica A* **189**, 43 (1992).

Duarte, J.A.M.S., J.M. Carvalho, and H.J. Ruskin, "The direction of maximum spread in anisotropic forest fires and its critical properties," *Physica A* **183**, 411 (1992).

Dubson, M.A., Y.C. Hui, M.B. Wiessman, and J.C. Garland, "Measurements of the fourth moment of the current distribution in two-dimensional random resistor networks," *Phys. Rev. B* **39**, 6807 (1989).

Dugdale, D.C., "Yielding of steel sheets containing slits," *J. Mech. Phys. Solids* **8**, 100 (1960).

Dulaney, E., and W. Brace, "Velocity behavior of a growing crack," *J. Appl. Phys* **31**, 2233 (1960).

Duva, J.M., and J.W. Hutchinson, "Constitutive potentials for dilutely voided nonlinear materials," *Mech. Mater.* **3**, 41 (1984).

Duxbury, P.M., P.D. Beale, and P.L. Leath, "Size effects of electrical breakdown in quenched random media," *Phys. Rev. Lett.*, **57**, 1052 (1986).

Duxbury, P.M., P.D. Beale, and P.L. Leath, "Breakdown properties of quenched random systems: the random-fuse network," *Phys. Rev. B* **36**, 367 (1987).

Duxbury, P.M., P.D. Beale, and C. Moukarzel, "Breakdown of two-phase random resistor networks," *Phys. Rev. B* **51**, 3476 (1995).

Duxbury, P.M., and P.L. Leath, "The failure distribution in percolation models of breakdown," *J. Phys. A* **20**, L411 (1987).

Duxbury, P.M., and P.L. Leath, "Exactly solvable models of material breakdown," *Phys. Rev. B* **49**, 12676 (1994a).

Duxbury, P.M., and P.L. Leath, "Failure probability and average strength of disordered systems," *Phys. Rev. Lett.* **72**, 2805 (1994b).

Duxbury, P.M., and Y. Li, in *Disorder and Fracture*, edited by C.J. Charmet, S. Roux, and E. Guyon (Plenum, New York, 1990), p. 141.

Dvorak, G.J., "Transformation field analysis of inelastic composite materials," *Proc. R. Soc. Lond. A* **437**, 311 (1992).

Dvorak, G.J., and Y.A. Bahei-El-Din, "A bimodal plasticity theory of fibrous composite material," *Acta Mech.* **69**, 219 (1987).

Dvorak, G.J., Y.A. Bahei-El-Din, Y. Macheret, and C.H. Liu, "An experimental study of elastic-plastic behavior of a fibrous boron-aluminum composite," *J. Mech. Phys. Solids* **36**, 655 (1988).

Edwards, S.F., and D.R. Wilkinson, "The surface statistics of a granular aggregate," *Proc. R. Soc. Lond. A* **381**, 17 (1982).

Efros, A.L., and B.I. Shklovskii, "Critical behaviour of conductivity and dielectric constant near the metal-non-metal transition threshold," *Phys. Status Solidi B* **46**, 475 (1976).

Elbern, H., "Parallelization and load balancing of a comprehensive atmospheric chemistry transport model," *Atmos. Environ.* **31**, 3561 (1997).

Elimes, A., R.A. Romer, and M. Schreiber, *Eur. Phys. J. B* **1**, 29 (1998).

Englman, R., and Z. Jaeger, "Third Bar-Ilan Conference on Frontiers in Condensed Matter Physics," *Physica A* **168**, 655 (1990).

Ercolessi, F., and J. Adams, "Interatomic potentials from first-principles calculations: The force-matching method," *Europhys. Lett.*, **26**, 583 (1993).

Ericksen, J. L., in *Phase Transformations and Material Instabilities in Solids*, edited by M. Gurtin (Academic Press, New York, 1984).

Ertas, D., and M. Kardar, "Dynamic roughening of directed lines," *Phys. Rev. Lett.* **69**, 929 (1992).

Ertas, D., and M. Kardar, "Dynamic relaxation of drifting polymers: a phenomenological approach," *Phys. Rev. E* **48**, 1228 (1993).

Ertas, D., and M. Kardar, "Anisotropic scaling in depinning of a flux line," *Phys. Rev. Lett.* **73**, 1703 (1994).

Ertas, D., and M. Kardar, "Anisotropic scaling in threshold critical dynamics of driven directed lines," *Phys. Rev. B* **53**, 3520 (1996).

Eshelby, J.D., "The elastic field of a crack extending non-uniformly under general anti-plane loading," *J. Mech. Phys. Solids* **17**, 177 (1969).

Eshelby, J.D., "Fracture mechanics," *Sci. Prog.* **59**, 161 (1971).

Español, P., I. Zúniga, and M.A. Rubio, "Effect of boundary conditions in mode I fracture in brittle materials," *Physica D* **96**, 375 (1996).

Essenlink, K., B. Smit, and P.A.J. Hilbers, "Efficient parallel implementation of molecular dynamics on a toroidal machine, Part I: Parallelizing strategy," *J. Comput. Phys.* **106**, 101 (1993).

Ewalds, H.L., and R.J.H. Wanhill, *Fracture Mechanics* (Arnold, London, 1986).

Falk, M.L., "Molecular-dynamics study of ductile and brittle fracture in model noncrystalline solids," *Phys. Rev. B* **60**, 7062 (1999).

Falk, M.L., and J.S. Langer, "Dynamics of viscoplastic deformation in amorphous solids," *Phys. Rev. E* **57**, 7192 (1998).

Family, F., and T. Vicsek, "Scaling of the active zone in the Eden process on percolation networks and the ballistic deposition model," *J. Phys. A* **18**, L75 (1985).

Family, F., and T. Vicsek (eds.), *Dynamics of Fractal Surfaces* (World Scientific, Singapore, 1991).

Family, F., Y.C. Zhang, and T. Vicsek, "Invasion percolation in an external field: dielectric breakdown in random media," *J. Phys. A* **19**, L733 (1986).

Feigel'man, M.V., and V.M. Vinokur, "Thermal fluctuations of vortex lines, pinning, and creep in high-$T_c$ superconductors," *Phys. Rev. B* **41**, 8986 (1990).

Fernandez, L., F. Guinea, and E. Louis, "Random and Dendritic patterns in crack propagation," *J. Phys. A* **21**, L301 (1988).

Ferrante, J., J.R. Smith, and J.H. Rose, "Diatomic molecules and metallic adhesion, cohesion, and chemisorption: a single binding-energy relation," *Phys. Rev. Lett.* **50**, 1385 (1983).

Feynman, R.P., "Forces in molecules," *Phys. Rev.* **56**, 340 (1939).

Field, J.E., "Brittle fracture: Its study and application," *Contemp. Phys.* **12**, 1 (1971).

Fineberg, J., S.P. Gross, M. Marder, and H.L. Swinney, "Instability in dynamic fracture," *Phys. Rev. Lett.* **67**, 457 (1991).

Fineberg, J., S.P. Gross, M. Marder, and H.L. Swinney, "Instability in the propagation of fast cracks," *Phys. Rev. B* **45**, 5146 (1992).

Fineberg, J., S.P. Gross, and E. Sharon, "Micro-branching as an instability in dynamic fracture," in *Nonlinear Analysis of Fracture*, edited by J.R. Willis (Kluwer Academic, Dordrecht, 1997), p. 177.

Fineberg, J., and M. Marder, "Instability in dynamic fracture," *Phys. Rep.* **313**, 1 (1999).

Fisher, D.S., "Sliding charge-density waves as a dynamic critical phenomenon," *Phys. Rev. B* **31**, 1396 (1985).

Fisher, D.S., M.P.A. Fisher, and D.A. Huse, "Thermal fluctuations, quenched disorder, phase transitions, and transport in type-II superconductors," *Phys. Rev. B* **43**, 130 (1991).

Fleck, N.A., and J.W. Hutchinson, "Strain gradient plasticity," *Adv. Appl. Mech.* **33**, 295 (1997).

Fleming, R.M., and C. C. Grimes, "Sliding-mode conductivity in $NbSe_3$: Observation of a threshold electric field and conduction noise," *Phys. Rev. Lett.* **42**, 1423 (1979).

Flytzanis, C., *Prog. Opt.* **29**, 2539 (1992).

Foiles, S.M., M.I. Baskes, and M.S. Daw, "Embedded-atom-method functions for fcc metals Cu, Ag, Au, Ni, Pd, Pt, and their alloys," *Phys. Rev. B* **33**, 7983 (1986).

Fonesca, I., *J. Math. Pure. Appl.* **67**, 175 (1988).

Form, W., N. Ito, and G.A. Kohring, "Vectorized and parallelized algorithms for multi-million particle MD-simulation," *Int. J. Mod. Phys. C* **4**, 1075 (1993).

Fox, G.C., M.A. Johnson, G.A. Lyzenga, S.W. Otto, J.K. Salmon, and D.W. Walker, *Solving Problems on Cocurrent Processors*, Vol. 1 (Prentice Hall, Englewood Cliffs, NJ, 1988).

Frenkel, D., and A.J.L. Ladd, "New Monte Carlo method to compute the free energy of solids. Application to the FCC and HCP phases of hard spheres," *J. Chem. Phys.* **81**, 3188 (1984).

Freund, L.B., "The mechanics of dynamic shear crack propagation," *J. Geophys. Res.* **84**, 2199 (1979).

Freund, L.B., *Dynamic Fracture Mechanics* (Cambridge University Press, Cambridge, 1990).

Freund, L.B., and W.D. Nix, "A critical thickness condition for a strained compliant substrate/epitaxial film system," *Appl. Phys. Lett.* **69**, 173 (1996).

Friel, J.J., and C.S. Pande, "A direct determination of fractal dimension of fracture surfaces using scanning electron microscopy and stereoscopy," *J. Mater. Res.* **8**, 100 (1993).

Fuller, K.N.G., P.G. Fox, and J.E. Field, "The temperature rise at the tip of fast-moving cracks in glassy polymers," *Proc. R. Soc. Lond. A* **341**, 1213 (1975).

Furukawa, H., "Propagation and pattern of crack in two dimensional dynamical lattices," *Prog. Theor. Phys.* **90**, 949 (1993).

Gadenne, P., F. Brouers, V.M. Shalaev, and A.K. Sarychev, "Giant Stokes fields on semicontinuous metal films," *J. Opt. Soc. Am. B* **15**, 68 (1998).

Galambos, J., *The Asymptotic Theory of Extreme Order Statistics* (Wiley, New York, 1978).

Gao, H., "Surface roughening and branching instabilities in dynamic fracture," *J. Mech. Phys. Solids* **41**, 457 (1993).

Gao, L., B.-C. Xie, and Z.-Y. Li, "Crossover exponents in percolating nonlinear normal conductor-insulator mixtures," *Physica A* **271**, 238 (1999).

Gărăjeu, M., and P. Suquet, "Effective properties of porous ideally plastic or viscoplastic materials containing rigid particles," *J. Mech. Phys. Solids* **45**, 873 (1997).

Garboczi, E.J., "Linear dielectric-breakdown electrostatics," *Phys. Rev. B* **38**, 9005 (1988).

Gear, C.W., *Numerical Initial Value Problems in Ordinary Differential Equations* (Prentice-Hall, Englewood Cliffs, NJ, 1971).

Gefen, Y., W.-H. Shih, R.B. Laibowitz, and J.M. Viggiano, "Nonlinear behavior near the percolation metal-insulator transition," *Phys. Rev. Lett.* **57**, 3097 (1986).

George, A., and H. Michot, "Dislocation loops at crack tips: nucleation and growth—an experimental study in silicon," *Mater. Sci. Eng. A* **164**, 118 (1993).

Germain, P., Q.S. Nguyen, and P. Suquet, "Continuum thermodynamics," *J. Appl. Mech.* **50**, 1010 (1983).

Gibiansky, L.V., and S. Torquato, "New approximation for the effective energy of non-linear conducting composites," *J. Appl. Phys.* **84**, 301 (1998a).

Gibiansky, L.V., and S. Torquato, "Effective energy of nonlinear elastic and conducting composites: Approximations and cross-property bounds," *J. Appl. Phys.* **84**, 5969 (1998b).

Gilabert, A., S. Roux, and E. Guyon, "Current-voltage characteristic of a nonlinear resistor network," *J. Physique* **48**, 1609 (1987).

Gilman, J.J., C. Knudsen, and W.P. Walsh, "Cleavage cracks and dislocations in LiF crystals," *J. Appl. Phys.* **6**, 601 (1958).

Gilormini, "Insuffisance de l'extension classique du modéle auto-cohérent au comportement non-linéaire," *C. R. Acad. Sci. Paris Ser. IIb* **320**, 115 (1995).

Goldstein, R.V., and R.K. Salganik, "Brittle fracture of solids with arbitrary cracks," *Int. J. Fract.* **10**, 507 (1974).

Golosovsky, M., M. Tsindlekht, D. Davidov, and A.K. Sarychev, "Effective-medium approach to the microwave properties of the high-$T_c$ superconductor-insulator composite," *Physica C* **209**, 337 (1993).

Gordon, J.E., *Structures or Why Things Don't Fall Down* (Penguin, London, 1978).

Gorkov, L.P., and G. Grüner (eds.), *Charge Density Waves in Solids* (Elsevier, Amsterdam, 1989).

Gouldstone, A., H.J. Koch, K.Y. Zeng, A.E. Ginnakopoulos, and S. Suresh, "Discrete and continuous deformation during nanoindentation of thin films," *Acta Mater.* **48**, 2277 (2000).

Greengard, L., and V.I. Rokhlin, "A fast algorithm for particle simulations," *J. Comput. Phys.* **73**, 325 (1987).

Grest, G.S., B. Dünweg, and K. Kremer, "Vectorized link cell Fortran code for molecular dynamics simulations for a large number of particles," *Comput. Phys. Commun.* **70**, 243 (1989).

Griffith, A.A., "The phenomena of rupture and flow in solids," *Philos. Trans. R. Soc. Lond.* **227**, 163 (1920).

Gross, S., *Dynamics of Fast Fracture*, Ph.D. Thesis, University of Texas, Austin (1995).

Gross, S.P., J. Fineberg, M. Marder, W.D. McCormick, and H.L. Swinney, "Acoustic emissions from rapidly moving cracks," *Phys. Rev. Lett.* **71**, 3162 (1993).

Gu, G.Q., and K.W. Yu, "Effective conductivity of nonlinear composites," *Phys. Rev. B* **46**, 4502 (1992).

Gumbel, E.J., *Statistics of Extremes* (Columbia University Press, New York, 1958).

Gumbsch, P., and H. Gao, "Dislocation faster than the speed of sound," *Science* **283**, 965 (1999).

Gumbsch, P., S.L. Zhou, and B.L. Holian, "Molecular dynamics investigation of dynamic crack stability," *Phys. Rev. B* **55**, 3445 (1997).

Günes, P., S. Simsek, and S. Erkoc, "A comparative study of empirical potential energy functions: Application to clusters," *Int. J. Mod. Phys. C* **11**, 451 (2000).

Gyure, M.F., and P.D. Beale, "Dielectric breakdown of a random array of conducting cylinders," *Phys. Rev. B* **40**, 9433 (1989).

Gyure, M.F., and P.D. Beale, "Dielectric breakdown in continuous models of metal-loaded dielectrics," *Phys. Rev. B* **46**, 3736 (1992).

Haile, J.M., and S. Gupta, "Extensions of molecular dynamics simulation method. II. Isothermal systems," *J. Chem. Phys.* **79**, 3067 (1983).

Hamann, D.R., M. Schlüter, and C. Chiang, "Norm-conserving pseudopotentials," *Phys. Rev. Lett.* **43**, 1494 (1979).

Hammonds, K.D., I.R. McDonald, and D.J. Tildesley, "Computational studies of the structure of carbon dioxide monolayers physisorbed on the basal plane of graphite," *Mol. Phys.* **70**, 175 (1990).

Hansen, A., E.L. Hinrichsen, and S. Roux, "Roughness of crack interfaces," *Phys. Rev. Lett* **66**, 2476 (1991b).

Hansen, A., S. Roux, and E.L. Hinrichsen, "Annealed model for breakdown processes," *Europhys. Lett.* **13**, 517 (1990).

Hansen, A., S. Roux, and H.J. Herrmann, "Rupture of central-force lattices," *J. Phys. France* **50**, 733 (1989).

Hansen, J.P., and McDonald, I.R., *Theory of Simple Liquids* (Academic Press, New York, 1986).

Harlow, D.G., *Proc. R. Soc. Lond. A* **397**, 211 (1985).

Harlow, D.G., and S.L. Phoenix, "The chain-of-bundles probability model for the strength of fibrous materials. I.Analysis and conjectures," *J. Comput. Mater.* **12**, 195 (1978).

Harlow, D.G., and S.L. Phoenix, "Approximations for the strength distribution and size effect in an idealized lattice model of material breakdown," *J. Mech. Phys. Solids* **39**, 173 (1991).

Harris, A.B., "Field-theoretics formulation of the randomly diluted nonlinear resistor network," *Phys. Rev. B* **35**, 5056 (1987).

Harris, A.B., T.C. Lubensky, W.K. Holcomb, and C. Dasgupta, "Renormalization-group approach to percolation," *Phys. Rev. Lett.* **35**, 327 (1974).

Harris, R.A., and L.R. Pratt, "Discretized propagators, Hartree, Hartree–Fock equation, and the Hohenberg–Kohn theorem," *J. Chem. Phys.* **82**, 856 (1985).

Hashin, Z., "The elastic moduli of heterogeneous materials," *J. Appl. Mech.* **29**, 143 (1962).

Hashin, Z. "Failure criteria for unidirectional fiber composites," *J. Appl. Mech.* **47**, 329 (1980).

Hashin, Z., and S. Shtrikman, "On some variational principles in anisotropic and nonhomogeneous elasticity," *J. Mech. Phys. Solids* **10**, 335 (1962a).

Hashin, Z., and S. Shtrikman, "A variational approach to the theory of the elastic behavior of polycrystals," *J. Mech. Phys. Solids* **10**, 343 (1962b).

Hashin, Z., and S. Shtrikman, "A variational approach to the theory of the elastic behavior of multiphase materials," *J. Mech. Phys. Solids* **11**, 127 (1963).

Hassold, G.N., and D.J. Srolovitz, "Brittle fracture in materials with random defects," *Phys. Rev. B* **39**, 9273 (1989).

Hauch, J., and M. Marder, "Energy balance in dynamic fracture, investigated by the potential drop technique," *Int. J. Fract.* (1999).

Heffelfinger, G.S., and F. van Swol, "Diffusion in Lennard–Jones fluids using dual control volume grand canonical molecular dynamics simulation (DCV-GCMD)," *J. Chem. Phys.* **100**, 7548 (1994).

Heiba, A.A., M. Sahimi, L.E. Scriven, and H.T. Davis, "Percolation theory of two-phase flow in porous media," Society of Petroleum Engineers paper 11015, New Orleans, LA (1982).

Heiba, A.A., M. Sahimi, L. E. Scriven, and H. T. Davis, "Percolation theory of two-phase relative permeability," *SPE Reservoir Engineering* **7**, 123 (1992).

Heino, P., and K. Kaski, "Mesoscopic model of crack branching," *Phys. Rev. B* **54**, 6150 (1996).

Heino, P., and K. Kaski, "Mesoscopic Maxwell-dissipative finite element model for crack propagation," *Int. J. Mod. Phys. C* **8**, 383 (1997).

Heinrichs, J., and N. Kumar, "Simple exact treatment of conductance in a random Bethe lattice," *J. Phys. C* **8**, L510 (1975).

Hellmann, H., *Einführung in die Quantumchemie* (Deuticke, Leipzig, 1937).

Hendrickson, B., and R. Leland, "An improved spectral graph partitioning algorithm for mapping parallel computations," *SIAM J. Sci. Stat. Comput.* **16**, 452 (1995).

Herrmann, H.J., A. Hansen, and S. Roux, "Fracture of disordered, elastic lattices in two dimensions," *Phys. Rev. B* **39**, 637 (1989a).

Herrmann, H.J., J. Kertész, and L. de Arcangelis, "Fractal shapes of deterministic cracks," *Europhys. Lett.* **10**, 147 (1989b).

Herrmann, H.J., and S. Roux (eds.), *Statistical Models for the Fracture of Disordered Media* (North-Holland, Amsterdam, 1990).

Herrmann, H.J., and M. Sahimi, "Fluid penetration through a crack in a pressure gradient," *J. Phys. A* **26**, L1145 (1993).

Herrmann, H.J., M. Sahimi, and F. Tzschichholz, "Examples of fractals in soil mechanics," *Fractals* **1**, 795 (1993).

Hervé, E., C. Stolz, and A. Zaoui, "A propos de l'assemblage des sphéres composites de Hashin," *C. R. Acad. Sci. Paris Ser. II* **313**, 857 (1991).

Hesselbo, B., D. Phil. Thesis, Oxford University (1994).

Hill, R., "Elastic properties of reinforced solids: Some theoretical principles," *J. Mech. Phys. Solids* **11**, 357 (1963).

Hill, R., "Continuum micro-mechanics of elastoplastic polycrystals," *J. Mech. Phys. Solids* **13**, 89 (1965a).

Hill, R., "Self-consistent mechanics of composite materials," *J. Mech. Phys. Solids* **13**, 213 (1965b).

Hill, R.M., and L.A. Dissado, "Theoretical basis for the statistics of dielectric breakdown," *J. Phys. C* **16**, 2145 (1983a).

Hill, R.M., and L.A. Dissado, "Examination of the statistics of dielectric breakdown," *J. Phys. C* **16**, 4447 (1983b).

Hinrichsen, E.L., A. Hansen, and S. Roux, "A fracture growth model," *Europhys. Lett.* **8**, 1 (1989).

Hirsch, P.B., R.W. Horne, and M.J. Whelan, "Direct observation of the arrangement and motion of dislocation in aluminum," *Philos. Mag. (series 8)* **1**, 677 (1956).

Hirsch, P.B., J. Samuels, and S.G. Roberts, "The brittle-ductile transition in silicon. II. Interpretation," *Proc. R. Soc. Lond. A* **421**, 25 (1989).

Ho, P.S., and T. Kwok, "Electromagnetism in metals," *Rep. Prog. Phys.* **52**, 301 (1989).

Hoagland, R.G., M.S. Daw, S.M. Foiles, and M.I. Baskes, "An atomic model of crack tip deformation in aluminum using an embedded atom potential," *J. Mater. Res.* **5**, 313 (1990).

Hockney, R.W., S.P. Goel, and J.W. Eastwood, "Quite high-resolution computer models of a plasma," *J. Comput. Phys.* **14**, 48 (1974).

Hodgdon, J.A., and J.P. Sethna, "Derivation of a general three-dimensional crack-propagation law: A generalization of the principle of local symmetry," *Phys. Rev. B* **47**, 4831 (1993).

Hoeing, A., "Some implications of an elastic-electrostatic analogy on certain path-independent integrals," *Int. J. Eng. Sci.* **22**, 87 (1984).

Hohenberg, P., and W. Kohn, "Inhomogeneous electron gas," *Phys. Rev.* **136**, B864 (1964).

Holian, B.L., and R. Ravelo, "Fracture simulations using large-scale molecular dynamics," *Phys. Rev. B* **51**, 11275 (1995).

Holian, B.L., A.F. Voter, and R. Ravelo, "Thermostatted molecular dynamics: How to avoid the Toda demon hidden in Nosé-Hoover dynamics," *Phys. Rev. E* **52**, 2338 (1995).

Holian, B.L., A.F. Voter, N.J. Wagner, R.J. Ravelo, S.P. Chen, W.G. Hoover, C.G. Hoover, J.E. Hammerberg, and T.D. Dontje, "Effects of pairwise versus many-bondy forces on high-stress plastic deformation," *Phys. Rev. A* **43**, 2655 (1991).

Holland, D., and M. Marder, "Ideal brittle fracture of silicon studied with molecular dynamics," *Phys. Rev. Lett.* **80**, 746 (1998).

Hoover, W.G., "Canonical dynamics: Equilibrium phase-space distributions," *Phys. Rev. A* **31**, 1695 (1985).

Hoover, W.G., and F.H. Ree, "Melting transition and communal entropy for hard spheres," *J. Chem. Phys.* **49**, 3609 (1968).

Horowitz, G.E., *Arch. Elektrotech. Berlin* **18**, 555 (1927).

Hough, S.E., "On the use of spectral methods for the determination of fractal dimension," *Geophys. Res. Lett.* **16**, 673 (1989).

Hua, L., H. Rafii-Tabar, and M. Cross, *Philos. Mag. Lett.* **75**, 237 (1997).

Hughes, B.D., *Random Walks and Random Environments*, Vol. 1 (Oxford University Press, London, 1995).

Hui, P.M., "Effective nonlinear response in dilute nonlinear granular materials," *J. Appl. Phys.* **68**, 3009 (1990a).

Hui, P.M., "Enhancement in nonlinear effects in percolating nonlinear resistor networks," *Phys. Rev. B* **41**, 1673 (1990b).

Hui, P.M., "Crossover electric field in percolating perfect-conductor-nonlinear-normal-metal composites," *Phys. Rev. B* **49**, 15344 (1994).

Hui, P.M., P. Cheung, and Y.R. Kwong, "Effective response in nonlinear random composites," *Physica A* **241**, 301 (1997).

Hull, D., "Effect of Crazes on the propagation of cracks in polystyrene," *J. Mater. Sci.* **5**, 357 (1970).

Hull, D., *Fractography* (Cambridge University Press, Cambridge, 1999).

Hull, D., and P. Beardmore, "Velocity of propagation of cleavage cracks in tungsten," *Int. J. Fract. Mech.* **2**, 468 (1966).

Huntington, H.B., "Electromigration in Metals," in *Diffusion in Solids-Recent Developments*, edited by A.S. Novick and J.J. Burton (Academic Press, New York, 1975), p. 120.

Imre, A., T. Pajkossy, and L. Nyikos, "Electrochemical determination of the fractal dimension of fractured surfaces," *Acta Metall. Mater.* **40**, 1819 (1992).

Inglis, C.E., "Stresses in a plate due to the presence of cracks and sharp corners," *Trans. Inst. Naval Archit.* **55**, 219 (1913).

Irwin, G.R., "Analysis of stresses and strains near the end of a crack traversing a plate," *J. Appl. Mech.* **24**, 361 (1956).

Irwin, G.R., "Fracture," in *Handbuch der Physik*, Vol. 6 (Springer, Berlin, 1958), p. 551.

Irwin, G.R., J.W. Dally, T. Kobayashi, W.L. Fourney, M.J. Etheridge, and H.P. Rossmanith, "On the determination of the a-K relationship for birefringent polymers," *Exp. Mech.* **19**, 121 (1979).

Jackson, J.D., *Classical Electrodynamics*, 3rd Ed. (John Wiley & Sons, 1998).

Janak, J.F., "Proof that $\partial E/\partial n_i = \epsilon_i$ in density-functional theory," *Phys. Rev. B* **18**, 7165 (1978).

Jarvinen, R., T. Mantyla, and P. Kettunen, "Improved adhesion between a sputtered alumina coating and a copper substrate," *Thin Solid Films* **114**, 311 (1984).

Jellinek, J., and R.S. Berry, "Generalization of Nosé's isothermal molecular dynamics," *Phys. Rev. A* **38**, 3069 (1988).

Joannopoulos, J.D., P. Bash, and A. Rappe, *Chemical Design Automation News* **6** (No. 8) (1991).

Johnson, E., "Process region changes for rapidly propagating cracks," *Int. J. Fract.* **55**, 47 (1992).

Johnson, E., "Process region influence on energy release rate and crack tip velocity during rapid crack propagation," *Int. J. Fract.* **61**, 183 (1993).

Johnson, H.T., and L.B. Freund, "Mechanics of coherent and dislocation island morphologies in strained epitaxial material systems," *J. Appl. Phys.* **81**, 6081 (1997).

Johnson, J.W., and D.G. Holloway, "On the shape and size of the fracture zones on glass fracture surfaces," *Philos. Mag.* **14**, 731 (1966).

Johnson, J.W., and D.G. Holloway, "Microstructure of the mist zone on glass fracture surfaces," *Philos. Mag.* **17**, 899 (1968).

Johnson, R.A., "Interstitials and vacancies in $\alpha$-iron", *Phys. Rev.* **134A**, 1329 (1964).

Johnson, R.A., "Analytic nearest-neighbor model for fcc metals," *Phys. Rev. B* **37**, 3924 (1988).

Jones, R.O., and O. Gunnarsson, "The density functional formalism, its applications and prospects," *Rev. Mod. Phys.* **61**, 689 (1989).

Kahng, B., G.G. Batrouni, S. Redner, L. de Arcangelis, and H.J. Herrmann, "Electrical breakdown in a fuse network with random, continuously distributed breaking strengths," *Phys. Rev. B* **37**, 7625 (1988).

Kalia, R.K., A. Nakano, A. Omeltchenko, K. Tsuruta, and P. Vashishta, "Role of ultrafine microstructures in dynamic fracture in nanophase silicon nitride," *Phys. Rev. Lett.* **78**, 2144 (1997).

Kallivayalil, J., and A.T. Zehnder, "A method for thermo-mechanical analysis of steady state dynamic crack growth," *Int. J. Fract.* **66** 99 (1994).

Kalthoff, J.F., "On the propagation direction of bifurcated cracks," in *Dynamic Crack Propagation*, edited by G.C. Sih (Noordhoof, Leydon, 1972), p. 449.

Kalthoff, J.F., "The shadow optical method of caustics," in *Static and Dynamic Photoelasticity and Caustics*, edited by A. Lagarde (Springer, Berlin, 1987), p. 401.

Kane, E.O., "Phonon spectra of diamond and zinc-blende semiconductors," *Phys. Rev. B* **31**, 7865 (1985).

Kantor, Y., and I. Webman, "Elastic properties of random percolating systems," *Phys. Rev. Lett.* **52**, 1891 (1984).

Kardar, M., "Fluctuations of interfaces and fronts," in *Disorder and Fracture*, edited by J.C. Charmet, S. Roux, and E. Guyon (Plenum, New York, 1990), p. 3.

Kardar, M., G. Parisi, and Y.C. Zhang, "Dynamic scaling of growing interfaces," *Phys. Rev. Lett.* **56**, 889 (1986).

Karma, A., D.A. Kessler, and H. Levine, "Phase-field model of Mode III dynamic fracture," *Phys. Rev. Lett.* **87**, 045501-1 (2001).

Kausch, H.H., *Polymer Fracture* (Springer, Berlin, 1987).

Kaveh, M., and N.F. Mott, "The conductivity of disordered systems and the scaling theory," *J. Phys. C* **14**, L659 (1981).

Kawarabayashi, T., B. Kramer, and T. Ohtsuki, "Anderson transitions in three-dimensional disordered systems with randomly varying magnetic flux," *Phys. Rev. B* **57**, 11842 (1998).

Kaxiras, E., and M.S. Duesbery, "Free energies of generalized stacking faults in Si and implications for the brittle-ductile transition," *Phys. Rev. Lett.* **70**, 3752 (1993).

Kaxiras, E., and K.C. Pandey, "New classical potential for accurate simulation of atomic processes in Si," *Phys. Rev. B* **38**, 12736 (1988).

Kelchner, C.J., S.J. Plimpton, and J.C. Hamilton, "Dislocation nucleation and defect structure during surface indentation," *Phys. Rev. B* **58**, 11085 (1998).

Kellomäki, M., J. Åström, and J. Timonen, "Rigidity and dynamics of random spring networks," *Phys. Rev. Lett.* **77**, 2730 (1996).

Kelly, A., W.R. Tyson, and A.H. Cotterell, "Ductile and brittle crystals," *Philos. Mag.* **15**, 567 (1967).

Kenkel, S.W., and J.P. Straley, "Percolation theory of nonlinear circuit elements," *Phys. Rev. Lett.* **49**, 767 (1982).

Kent, R., and R. Rat, "Static electricity phenomena in the manufacture and handling of solid propellants," *J. Electrostat.* **17**, 299 (1985).

Kerkhof, F., "Wave fractographic investigations of brittle fracture dynamics," in *Dynamic Crack Propagation,* edited by G.C. Sih (Noordhoff International Publishing, Leyden, 1973), p. 3.

Khanta, M., D.P. Pope, and V. Vitek, "Dislocation screening and the brittle-to-ductile transition: A Kosterlitz-Thouless type instability," *Phys. Rev. Lett.* **73**, 684 (1994).

Kim, J.M., and J.M. Kosterlitz, "Growth in a restricted solid-on-solid model," *Phys. REv. Lett.* **62**, 2289 (1989).

Kim, K.S., "Dynamic fracture under normal impact loading of the crack faces," *J. Appl. Mech.* **52**, 585 (1985).

Kim, S., and S. Karrila, *Microhydrodynamics. Principles and Selected Applications*, (Butterworth-Heinemann, Boston, 1991).

Kinloch, A.J., and R.J. Young, *Fracture Behavior of Polymers* (Applied Sciences, London, 1983).

Kinzel, W., "Directed percolation," in *Percolation Structures and Processes*, edited by G. Deutscher, R. Zallen, and J. Adler (Adam Hilger, Bristol, 1983), p. 425.

Knap, J., and M. Ortiz, "An analysis of the quasicontinuum method," *J. Mech. Phys. Solids* **49**, 1899 (2001).

Knaur, J.A., and P.P. Budenstein, *IEEE Trans.* **EI-15**, 313 (1980).

Kneipp, K., Y. Wang, H. Kneipp, L.T. Perelman, I. Itzkan, R. Dasari, and M.S. Feld, "Single molecule detection using surface-enhanced Raman scattering (SERS)," *Phys. Rev. Lett.* **78**, 1667 (1997).

Kobayashi, A., N. Ohtani, and T. Sato, "Phenomenological aspects of viscoelastic crack propagation," *J. Appl. Poly. Sci.* **18**, 1625 (1974).

Kobayashi, A.S., in *Handbook on Experimental Mechanics*, edited by A.S. Kobayashi (American Society for Testing and Materials, Prentice-Hall, Englewood Cliffs, NJ, 1987), p. 231.

Kohlhoff, S., P. Gumbsch, and H. F. Fischmeister, "Crack propagation in BCC crystals studied with a combined finite-element and atomistic model," *Philos. Mag.* **64A**, 851 (1991).

Kohn, R.V., and T.D. Little, "Some model problems of polycrystal plasticity with deficient basic crystals," (1997).

Kohn, R.V., and G.W. Milton, "On bounding the effective conductivity of anisotropic composites," in *Homogenization and Effective Moduli of Materials and Media*, edited by J.L. Ericksen, D. Kinderlehrer, R.V. Kohn, and J.-L. Lions (Springer, New York, 1986), p. 97.

Kohn, W., and L.J. Sham, "Self-consistent equations including exchange and correlation effects," *Phys. Rev.* **140**, A1133 (1965).

Kohr, K.E., and S. Das Sarma, "Model-potential study of $(2n + 1) \times (2n + 1)$ reconstructions on Si(111) surface," *Phys. Rev. B* **40**, 1319 (1989).

Kolsky, H., *Stress Waves in Solids* (Oxford University Press, London, 1953).

Koplik, J., and H. Levine, "Interface moving through a random background," *Phys. Rev. B* **32**, 280 (1985).

Koschmieder, E.L., *Bérnard Cells and Taylor Vortices* (Cambridge University Press, London, 1993).

Kosterlitz, J.M., and D.J. Thouless, "Ordering, metastability and phase transitions in two-dimensional systems," *J. Phys. C* **6**, 1181 (1973).

Kramer, B., and A. MacKinnon, "Localization: theory and experiment," *Rep. Prog. Phys.* **56**, 1469 (1993).

Kuang, L., and H.J. Simon, "Diffusely scattered second harmonic generation from a silver film due to surface plasmons," *Phys. Lett. A* **197**, 257 (1995).

Kulakhmetova, S.A., V.A. Saraikin, and L.I. Slepyan, "Plane problem of a crack in a lattice," *Mech. Solids* **19**, 102 (1984).

Kulawansa, D.M., L.C. Jensen, S.C. Langford, and J.T. Dickinson, "Scanning electron microscopy of the mirror region of silicate glass fracture surfaces," *J. Mater. Res.* **9**, 476 (1993).

Kun, F., A. Zapperi, and H.J. Herrmann, "Damage in fiber bundle models," *Eur. Phys. J. B* **17**, 269 (2000).

Kunin, I.A., *Elastic Media with Microstructure* (Springer, Berlin, 1982).

Kusy, R.P., and D.T. Turner, "Influence of the molecular weight of PMMA on fracture surface energy in notched tension," *Polymer* **17**, 161 (1975).

Kusy, R.P., and D.T. Turner, "Influence of the molecular weight of poly(methyl methacrylate) on fracture morphology in notched tension," *Polymer* **18**, 391 (1977).

Lagarkov, A.N., L.V. Panina, and A.K. Sarychev, "Effective magnetic permeability of composite materials near the percolation threshold," *Sov. Phys. JETP* **66**, 123 (1987) [*Zh. Eskp. Teor. Fiz.* **93**, 215 (1987)].

Lamaignere, L., F. Carmona, and D. Sornette, "Experimental realization of critical thermal fuse rupture," *Phys. Rev. Lett.* **77**, 2738 (1996).

Landau, L.D., and E.M. Lifshitz, *Statistical Physics* (Pergamon Press, London, 1980).

Langer, J.S., "Models of crack propagation," *Phys. Rev. A* **46**, 3123 (1992).

Langer, J.S., "Dynamic model of onset and propagation of fracture," *Phys. Rev. Lett.* **70**, 3592 (1993).

Langer, J.S., and A.E. Lobkovsky, "Critical examination of cohesive-zone models in the theory of dynamic fracture," *J. Mech. Phys. Solids* **46**, 1521 (1998).

Langer, J.S., and H. Nakanishi, "Models of crack propagation. II. Two-dimensional model with dissipation on the fracture surface," *Phys. Rev. E* **48**, 439 (1993).

Langer, J.S., and C. Tang, "Rupture propagation in a model of an earthquake fault," *Phys. Rev. Lett.* **67**, 1043 (1991).

Langford, S.C., M. Zhenyi, and J.T. Dickinson, "Photon emission as a probe of chaotic processes accompanying fracture," *J. Mater. Res.* **4**, 1272 (1989).

Larkin, A.I., and Yu.N. Ovchinnikov, "Pinning in type II superconductors," *J. Low Temp. Phys.* **34**, 409 (1979).

Larralde, H., and R.C. Ball, "The shape of slowly growing cracks," *Europhys. Lett.* **30**, 87 (1995).

Larson, R.G., "Derivation of generalized Darcy equations for creeping flow in porous media," *Ind. Eng. Chem. Fund.* **20**, 132 (1981).

Lawn, B., *Fracture of Brittle Solids*, 2nd ed. (Cambridge University Press, Cambridge, 1993).

Leath, P.L., and P.M. Duxbury, "Fracture of heterogeneous materials with continuous distributions of local breaking strengths," *Phys. Rev. B* **49**, 14905 (1994).

Lee, H.-C., and M.E. Mear, "Effective properties of power-law solids containing elliptical inhomogenities. I. Rigid inclusions," *Mech. Mater.* **13**, 313 (1992).

Lee, H.-C., and K.W. Yu, "Effective medium theory for strongly nonlinear composites: comparison with numerical simulations," *Phys. Lett. A* **197**, 341 (1995).

Lemaire, E., Y. Ould Mohamed Abdelhaye, J. Larue, R. Benoit, P. Levitz, and H. Van Damme, "Pattern formation in noncohesive and cohesive granular media," *Fractals* **1**, 968 (1993).

Leung, K.-t., and Z. Néda, "Pattern formation and selection in quasistatic fracture," *Phys. Rev. Lett.* **85**, 662 (2000).

Levin, V.M., "Thermal expansion coefficients of heterogeneous materials," *Mekh. Tverd. Tela* **2**, 83 (1967).

Levy, O., and D.J. Bergman, "The bulk effective response of non-linear random resistor networks: numerical study and analytic approximations," *J. Phys.: Condens. Matter* **5**, 7095 (1993).

Levy, O., and D.J. Bergman, "Intrinsic optical bistability and resonances in nonlinear composites," *Physica A* **207** 157 (1994a).

Levy, O., and D.J. Bergman, "Critical behavior of the weakly nonlinear conductivity and flicker noise of two-component composites," *Phys. Rev. B* **50**, 3652 (1994b).

Levy, O., and D.J. Bergman, and D.G. Stroud, "Harmonic generation, induced nonlinearity, and optical bistability in nonlinear composites," *Phys. Rev. E* **52**, 3184 (1995).

Levy, O., and R.V. Kohn, "Duality relations for non-Ohmic composites, with application to behavior near percolation," *J. Stat. Phys.* **90**, 159 (1998).

Levy-Nathansohn, R., and D.J. Bergman, "Decoupling and testing of the generalized Ohm's law," *Phys. Rev. B* **55** 5425 (1997).

Li, G., P. Ponte Castañeda, and A.S. Douglas, "Constitutive models for ductile solids reinforced by rigid spheroidal inclusions," *Mech. Mater* **15**, 279 (1993).

Li, R., and K. Sieradzki, "Ductile-brittle transition in random porous Au," *Phys. Rev. Lett.* **68**, 1168 (1992).

Li, W., R.K. Kalia, and P. Vashishta, "Amorphization and fracture in silicon diselenide nanowires: a molecular dynamics study," *Phys. Rev. Lett.* **77**, 2241 (1996).

Li, X.-P., R.W. Nunes, and D. Vanderbilt, "Density-matrix electronic-structure method with linear system-size scaling," *Phys. Rev. B* **47**, 10891 (1993).

Li, Y.S., and P.M. Duxbury, "Size and location of the largest current in a random resistor network," *Phys. Rev. B* **36**, 5411 (1987).

Li, Y.S., and P.M. Duxbury, "Crack arrest by residual bonding in resistor and spring networks," *Phys. Rev. B* **38**, 9257 (1988).

Liao, H.B., R.F. Xiao, H. Wang, K.S. Wong, and G.K.L. Wong, "Large third-order optical nonlinearity in $Au:TiO_2$ composite films measured on a femtosecond time scale," *Appl. Phys. Lett.* **72**, 1817 (1998).

Lidorikis, E., M.E. Bachlechner, R.K. Kalia, A. Nakano, P. Vashishta, and G.Z. Voyiadjis, "Coupling length scales for multiscale atomistic-continuum simulations: Atomistically induced stress distribution in $Si/Si_3N_4$ nanopixels," *Phys. Rev. Lett.* **87**, 086104-1 (2001)

Lim, S., Tsotsis, T.T., and Sahimi, M., "Molecular dynamics simulation of porous amorphous carbon and transport of fluid mixtures therein," *J. Chem. Phys.* (2003).

Lin, J.J., "Nonlinear $I - V$ characteristics of the granular $PrBa_2Cu_3O_{7-\delta}$ compound," *J. Phys. Soc. Japan* **61**, 4125 (1992).

Lin, S.L., J. Mellor-Crumney, B.M. Pettitt, and G.N. Phillips, Jr., "Molecular dynamics on a distributed-memory multiprocessor," *J. Comput. Phys.* **13**, 1022 (1992).

Lobb, C.J., P.M. Hui, and D. Stroud, "Nonuniversal breakdown behavior in superconducting and dielectric composites," *Phys. Rev. B* **36**, 1956 (1987).

Lomdahl, P.S., D.M. Beazley, P. Tamayo, and N. Grφnbech-Jensen, "Multi-million particle molecular dynamics on the CM-5," *Int. J. Mod. Phys. C* **4**, 1074 (1993).

López, J.M., and J. Schmittbuhl, "Anomalous scaling of fracture surfaces," *Phys. Rev. E* **57**, 6405 (1998).

Louis, E., and F. Guinea, "The fractal nature of fracture," *Europhys. Lett.* **3**, 871 (1987).

Lupkowski, M., and F. van Swol, "Ultrathin films under shear," *J. Chem. Phys.* **95**, 1995 (1995).

Ma, H., Xiao, R., and P. Sheng, "Third-order optical nonlinearity enhancement through composite microstructures," *J. Opt. Soc. Am. B* **15**, 1022 (1998).

MacElroy, J.M.D., "Nonequilibrium molecular dynamics simulation of diffusion and flow in thin microporous membranes," *J. Chem. Phys.* **101**, 5274 (1994).

Machová, A., "Dynamic microcrack initiation in alpha -iron," *Mater. Sci. Eng.* **A206**, 279 (1996).

Machta, J., and R.A. Guyer, "Largest current in a random resistor network," *Phys. Rev. B* **36**, 2142 (1987).

Mahadevan, M., R.M. Bradley, and J.-M. Debierre, "Simulations of an electromigration-induced edge instability in single-crystal metal lines," *Europhys. Lett.* **45**, 680 (1999).

Måløy, K.J., A. Hansen, E.L. Hinrichsen and S. Roux, "Experimental measurements of the roughness of brittle cracks," *Phys. Rev. Lett.* **68**, 213 (1992).

Mandal, P., A. Neumann, A.G. Jansen, P. Wyder, and R. Deltour, "Temperature and magnetic-field dependence of the resistivity of carbon-black polymer composites," *Phys. Rev. B* **55**, 452 (1997).

Mandel, J., *Plasticité Classique et Viscoplasticité* CISM Udine Courses and Lectures, **97** (Springer, Berlin, 1972).

Mandelbrot, B.B., *The Fractal Geometry of Nature* (W.H. Freeman, San Francisco, 1982).

Mandelbrot, B.B., "Self-affine fractals and fractal dimension," *Physica Scripta* **32**, 257 (1985).

Mandelbrot, B.B., D.E. Passoja, and A.J. Paullay, "Fractal character of fracture surfaces of metals," *Nature* **308**, 721 (1984).

Mandelbrot, B.B., and J.W. Van Ness, "Fractional Brownian motions, fractional noise and applications," *SIAM Rev.* **10**, 422 (1968).

Manna, S.S., and B.K. Chakrabarti, "Dielectric breakdown in the presence of random conductors," *Phys. Rev. B* **36**, 4078 (1987).

Manogg, P., "Investigation of the rupture of a Plexiglas plate by means of an optical method involving high speed filming of the shadows originating around holes drilled in the plate," *Int. J. Fract.* **2**, 604 (1966).

Mantese, J.W., W.I. Goldberg, D.H. Darling, H.G. Craighead, U.J. Gibson, R.A. Buhrman, and W.W. Webb, "Excess low frequency conduction noise in a granular composite," *Solid State Commun.* **37**, 353 (1981).

Marcellini, P., "Periodic Solutions and homogenization of nonlinear variational problems," *Ann. Mat. Pura. Appl.* **4**, 139 (1978).

Marder, M., "New dynamical equation for cracks," *Phys. Rev. Lett.* **66**, 2484 (1991).

Marder, M., "Statistical mechanics of cracks," *Phys. Rev. E* **54**, 3442 (1996).

Marder, M., and J. Fineberg, "How things break," *Phys. Today* **49** (No. 9), 24 (1996).

Marder, M., and S.P. Gross, "Origin of crack tip instabilities," *J. Mech. Phys. Solids* **43**, 1 (1995).

Marder, M., and X. Liu, "Instability in lattice fracture," *Phys. Rev. Lett.* **71**, 2417 (1993).

Markel, V.A., V.M. Shalaev, E.B. Stechel, W. Kim, and R.L. Armstrong, "Small-particle composites. I. Linear optical properties," *Phys. Rev. B* **53**, 2425 (1996).

Markel, V.A., V.M. Shalaev, P. Zhang, W. Huynh, L. Tay, T.L. Haslett, and M. Moskovits, "Near-field optical spectroscopy of individual surface-plasmon modes in colloid clusters," *Phys. Rev. B* **59**, 10903 (1999).

Marks, N.A., D.R. McKenzie, B.A. Paithorpe, M. Bernasconi, and M. Parrinello, "*Ab initio* simulations of tetrahedral amorphous carbon," *Phys. Rev. B* **54**, 9703 (1996).

Martin, J.E., and M.B. Heaney, "Reversible thermal fusing model of carbon black current-limiting thermistors," *Phys. Rev. B* **62**, 9390 (2000).

Martín, T., P. Español, M.A. Rubio, and I. Zúñiga, "Dynamic fracture in a discrete model of a brittle elastic solid," *Phys. Rev. E* **61**, 6120 (2000).

Martins, J.L., and A. Zunger, "Bond lengths around isovalent impurities and in semiconductor solid solutions," *Phys. Rev. B* **30**, 6217 (1984).

Martyna, G.J., M.L. Klein, and M. Tuckerman, "Nose-Hoover chains: the canonical ensemble via continuous dynamics," *J. Chem. Phys.* **97**, 2635 (1992).

Martyna, G.J., D.T. Tobias, and M.L. Klein, "Constant pressure molecular dynamics algorithms," *J. Chem. Phys.* **101**, 4177 (1994).

Matthews, J.W., and A.E. Blakeslee, "Defects in epitaxial multilayers. I. Misfit dislocations," *J. Cryst. Growth* **32**, 265 (1974).

Mauri, F., G. Galli, and R. Car, "Orbital formulation for electronic-structure calculations with linear system-size scaling," *Phys. Rev. B*, **47**, 9973 (1993).

McAnulty, P., L.V. Meisel, and P.J. Cote, "Hyperbolic distributions and fractal character of fracture surfaces," *Phys. Rev. A* **46**, 3523 (1992).

Meakin, P., "A simple model for elastic fracture in thin films," *Thin Solid Films* **151**, 165 (1987).

Meakin, P., "Simple kinetic models for material failure and deformation," Herrmann and Roux (1990), p. 291.

Meakin, P., *Fractals, Scaling and Growth far from Equilibrium* (Combridge University Press, Cambridge, 1998).

Meakin, P., G. Li, L.M. Sander, E. Louis, and F. Guinea, "A simple two-dimensional model for crack propagation", *J. Phys. A* **22**, 1393 D (1989).

Mecholsky, J.J., in *Strength of Inorganic Glass*, edited by C.R. Kurkjian (Plenum, New York, 1985).

Mecholsky, J.J., T.J. Mackin, and D.E. Passoja, *Adv. Ceramics* **22**, 127 (1988).

Mecholsky, J.J., D.E. Passoja, and K.S. Feinberg-Ringel, "Quantitative analysis of brittle fracture surfaces using fractal geometry," *J. Am. Ceram. Soc.* **72**, 60 (1989).

Meek, J.M., and J.D. Craggs, *Electrical Breakdown of Gases* (Wiley, New York, 1978).

Mehrabi, A.R., H., Rassamdana, and M., Sahimi, "Characterization of long-range correlations in complex distributions and profiles," *Phys. Rev. E* **56**, 712 (1997).

Mehrabi, A.R., and M. Sahimi, "Diffusion of ionic particles in charged disordered media," *Phys. Rev. Lett.* **82**, 735 (1999).

Meir, Y., R. Blumenfeld, A. Aharony, and A.B. Harris, "Series analysis of randomly diluted nonlinear resistor networks," *Phys. Rev. B* **34**, 3424 (1986).

Mel'cuk, A.I., R.C. Giles, and H. Gould, "Molecular dynamics simulation of liquids on the Connection Machine," *Computers in Physics*, 311 (May/June 1991).

Mercer, J.L., "Tight-binding models for compounds: Application to SiC," *Phys. Rev. B* **54**, 4650 (1996).

Mermin, N.D., "Thermal properties of the inhomogeneous electron gas," *Phys. Rev.* **137**, A1441 (1965).

Michel, J.C., "A self-consistent estimate of the potential of a composite made of two power-law phases with the same exponent," *C. R. Acad. Sci. Paris Ser. II* **322**, 447 (1996).

Mikitishin, S.I., Y.-N. Skhonitskii, and A.N. Tynnyi, *Fiz. Khim. Mekh. Mater.* **5**, 69 (1969).

Miksis, M.J., "Effective dielectric constant of a nonlinear composite material," *SIAM J. Appl. Math.* **43**, 1140 (1983).

Miller, S., and R. Reifenberger, "Improved method for fractal analysis using scanning probe microscopy," *J. Vac. Sci. Technol.* **B10**, 1203 (1992).

Milman, V.Y., "Fracture surfaces: A critical review of fractal studied and a novel morphological analysis of scanning tunneling microscopy measurements," *Prog. Mater. Sci.* **38**, 425 (1994).

Milman, V.Y., R. Blumenfeld, N.A. Stelmashenko, and R.C. Ball, "Comment on 'Experimental measurements of the roughness of brittle cracks'," *Phys. Rev. Lett.* **71**, 204 (1993).

Milman, V.Y., N.A. Stelmashenko, and R. Blumenfeld, "Fracture surfaces: A critical review of fractal studies and a novel morphological analysis of scanning tunneling microscopy measurements," *Prog. Mater. Sci.* **38**, 425 (1994).

Milstein, F., in *Mechanics of Solids*, edited by H.G. Hopkins and M.J. Sewell (Pergamon Press, Oxford, 1982).

Milton, G.W., "Bounds on the electromagnetic, elastic, and other properties of two-component composites," *Phys. Rev. Lett.* **46**, 542 (1981a).

Milton, G.W., "Bounds on the elastic and transport properties of two-component composites," *J. Mech. Phys. Solids* **30**, 177 (1981b).

Milton, G.W., "Bounds on the elastic and transport properties of two-component composites," *J. Mech. Phys. Solids* **30**, 177 (1982).

Milton, G., "Modeling the properties of composites by laminates," in *Homogenization and Effective Moduli of Materials and Media*, edited by J.L. Ericksen, D. Linderlehrer, R. Kohn, and J.-L. Lions (Springer, Berlin, 1986), p. 150.

Mistriotis, A.D., N. Flytzanis, and S.C. Farantos, "Potential model for silicon clusters," *Phys. Rev. B* **39**, 1212 (1989).

Mitchell, M.W., and D.A. Bonnell, "Quantitative topographic analysis of fractal surfaces by scanning tunneling microscopy," *J. Mater. Res.* **5**, 2244 (1990).

Molinari, A., G.R. Canova, and S. Ahzi, "A self-consistent approach of the large deformation polycrystal viscoplasticity," *Acta Metall. Mater.* **35**, 2983 (1987).

Morales, J.J., and M.J. Nuevo, "Comparison of link-cell and neighborhood tables on a range of computers," *Comput. Phys. Commun.* **69**, 223 (1992).

Morel, S., J. Schmittbuhl, E. Bouchaud, and G. Valentin, "Scaling of crack surfaces and implications for fracture mechanics," *Phys. Rev. Lett.* **85**, 1678 (2000).

Morel, S., J. Schmittbuhl, J.M. López, and G. Valentin, *Phys. Rev. E* **58**, 6999 (1998).

Moreno, Y., J.B. Gómez, and A.F. Pacheco, "Fracture and second-order phase transitions," *Phys. Rev. Lett.* **85**, 2865 (2000).

Mori, Y., K. Kaneko, and M. Wadati, "Fracture dynamics by quenching. I. Crack patterns," *J. Phys. Soc. Japan* **60**, 1591 (1991).

Morrissey, J.W., and J.R. Rice, "Crack front waves," *J. Mech. Phys. Solids* **46**, 467 (1998).

Moskovitz, M., "Surface-enhanced spectroscopy," *Rev. Mod. Phys.* **57**, 783 (1985).

Mosolov, A.B., "Mechanics of fractal cracks in brittle solids," *Europhys. Lett.* **24**, 673 (1993).

Mott, N.F., "Brittle fracture in mild steel plates," *Engineering* **165**, 16 (1948).

Moulinec, H., and P. Suquet, "A FTT-based numerical method for computing the mechanical properties of composites from images of their microstructure," in *Microstructure-Properly Interactions in Composite Materials*, edited by R. Pyrz (Kluwer, The Netherlands, 1995), p. 235.

Movchan, A.B., and J.R. Willis, "Dynamic weight functions for a moving crack. II. Shear loading," *J. Mech. Phys. Solids* **43**, 1369 (1995).

Mu, Z.Q., and C.W. Lung, "Studies on the fractal dimension and fracture toughness of steel," *J. Phys. D* **21**, 848 (1988).

Müller, K., B. Mehlig, F. Milde, and M. Schreiber, "Statistics of wave functions in disordered and in classically chaotic systems," *Phys. Rev. Lett.* **78**, 215 (1997).

Mura, D., L. Colombo, R. Bertoncini, and G. Mula, "Structure and chemical order of bulk $Si_{1-x}C_x$ amorphous alloys," *Phys. Rev. B* **58**, 10357 (1998).

Murray, R.T., C.J. Keily, and M. Hopkinson, "Crack formation in III-V epilayers grown under tensile strain on InP(001) substrates," *Philos. Mag.* **74**, 383 (1996).

Murthy, C.S., S.F. O'Shea, and I.R. McDonald, "Electrostatic interactions in molecular crystals. Lattice dynamics of solid nitrogen and carbon dioxide," *Mol. Phys.* **50**, 531 (1983).

Muskhelishvili, N.I., *Some Basic Problems of the Mathematical Theory of Elasticity* (P. Noordhoff, Groningen, Holland, 1953).

Nakamura, T., and D.M. Parks, "Three-dimensional stress field near the crack front of a thin elastic plate," *J. Appl. Mech.* **55**, 805 (1988).

Nakano, A., R.K. Kalia, and P. Vashishta, "Growth of pore interfaces and roughness of fracture surfaces in porous silica: million particle molecular-dynamics simulations," *Phys. Rev. Lett.* **73**, 2336 (1994).

Nakano, A., R.K. Kalia, and P. Vashishta, "Dynamics and morphology of brittle cracks: a molecular-dynamics study of silicon nitride," *Phys. Rev. Lett.* **75**, 3138 (1995).

Namgoong, E., and J.S. Chun, "The effect of ultrasonic vibration on hard chromium plating in a modified selfregulating high speed bath," *Thin Solid Films* **120**, 153 (1984).

Narayan, O., and D.S. Fisher, "Dynamics of sliding charge-density waves in $4 - \epsilon$ dimensions," *Phys. Rev. Lett.* **68**, 3615 (1992).

Narayan, O., and D.S. Fisher, "Nonlinear fluid flow in random media: critical phenomena near threshold," *Phys. Rev. B* **49**, 9469 (1994).

Nelson, D.R., and B.I. Halperin, "Dislocation-mediated melting in two dimensions," *Phys. Rev. B* **19**, 2457 (1979).

Nie, S., and S.R. Emory, "Probing single molecules and single nanoparticles by surface-enhanced Raman scattering," *Science* **275**, 1102 (1997).

Niemeyer, L., L. Pietronero, and H.J. Wiesmann, "Fractal dimension of dielectric breakdown," *Phys. Rev. Lett.* **52**, 1033 (1984).

Nishioka, K., and J.K. Lee, "Temperature dependence of the ideal fracture strength of a b.c.c. crystal," *Philos. Mag. A* **44**, 779 (1981).

Nishioka, K., S. Nakamura, T. Shimamoto, and H. Fujiwara, "Lattice instability theory of fracture," *Scripta Metall.* **14**, 497 (1980).

Noble, B., *Methods Based on the Wiener–Hopf Technique for the Solution of Partial Differential Equations* (Pergamon, New York, 1958).

Noble, B., and J.W. Daniel, *Applied Linear Algebra*, 2nd ed. (Prentice-Hall, Englewood Cliffs, NJ, 1977).

Nosé, S., "A unified formulation of the constant temperature molecular dynamics methods," *J. Chem. Phys.* **81**, 511 (1984).

Noskov, M.D., V.R. Kukhta, and V.V. Lopatin, "Simulation of the electrical discharge development in inhomogeneous insulators," *J. Phys. D* **28**, 1187 (1995).

Nyikos, L., and T. Pajkossy, "Short communication—Diffusion to fractal surfaces," *Electrochim. Acta* **31**, 1347 (1985).

Obukhov, S.P., "The problem of directed percolation," *Physica A* **101**, 145 (1980).

O'Dwyer, J.J., *The Theory of Electrical Conduction and Breakdown in Solid Dielectrics* (Clarendon Press, Oxford, 1973).

Ohring, M., *Reliability and Failure of Electronic Materials and Devices* (Academic Press, San Diego, 1998).

Olsen, K.B., R. Madariaga, and R.J. Archuleta, "Three-dimensional dynamic simulation of the 1992 Landers earthquake," *Science* **278**, 834 (1997).

Olson, T., "Improvements on Taylor's upper bound for rigid-plastic composites," *Mater. Sci. Eng.* **A 175**, 15 (1994).

Omeltchenko, A., A. Nakano, R.K. Kalia, and P. Vashishta, "Structure, mechanical properties, and thermal transport in microporous silicon nitride-molecular-dynamics simulations on a parallel machine," *Europhys. Lett.* **33**, 667 (1996).

Omeltchenko, A., J. Yu, R.K. Kalia, and P. Vashishta, "Crack front propagation and fracture in a graphite sheet: a molecular dynamics study on parallel computers," *Phys. Rev. Lett.* **78**, 2148 (1997).

Ordejón, P., D.A. Drabold, R.M. Martin, and M.P. Grumbach, "Linear system-size scaling methods for electronic-structure calculations," *Phys. Rev. B* **51**, 1456 (1995).

Orowan, E., "Crystal plasticity—III. On the mechanism of the glide process," *Z. Phys.* **89**, 605 (1934).

Orowan, E., "Energy criteria of fracture," *Weld. Res. Supp.* **34**, 157 (1955).

Ortiz, M., and R. Phillips, "Nanomechanics of defects in solids," *Adv. Appl. Mech.* **36**, 1 (1999).

Pande, C.S., L.E. Richards, N. Louat, B.D. Dempsey, and A.J. Schwoeble, "Fractal characterization of fractured surfaces," *Acta Metall.* **35**, 1633 (1987).

Pannetta, C., L. Reggiani, and Gy. Trefán, "Scaling and universality in electrical failure of thin films," *Phys. Rev. Lett.* **84**, 5006 (2000).

Parleton, L.G., "Determination of the growth of branched cracks by numerical methods," *Eng. Fract. Mech.* **11**, 343 (1979).

Parr, R.G., and W. Yang, *Density Functional Theory of Atoms and Molecules* (Oxford University Press, New York, 1989).

Parrinello, M., and A. Rahman, "Polymorphic transitions in single crystals: A new molecular dynamics method," *J. Appl. Phys.* **52**, 7182 (1981).

Paskin, P., A. Gohar, and G.J. Dienes, "Computer simulation of crack propagation," *Phys. Rev. Lett.* **44**, 940 (1980).

Paskin, P., D.K. Som, and G.J. Dienes, "Computer simulation of crack propagation: lattice trapping," *J. Phys. C* **14**, L171 (1981).

Passoja, D.E., *Adv. Ceramics* **22**, 101 (1988).

Passoja, D.E., and D.J. Amborski, *Microstruct. Sci.* **6**, 143 (1978).

Payne, M.C., M.P. Teter, D. Allen, T.A. Arias, and J.D. Joannopoulos, "Iterative minimization techniques for *ab initio* total-energy calculations: molecular dynamics and conjugate gradients," *Rev. Mod. Phys.* **64**, 1045 (1992).

Pederson, M.R., and K.A. Jackson, "Pseudoenergies for simulations on metallic systems," *Phys. Rev. B* **43**, 7312 (1991).

Perdew, J.P., K. Burke, and M. Ernzerhof, "Generalized gradient approximation made simple," *Phys. Rev. Lett.* **77**, 3865 (1996).

Perdew, J.P., and A. Zunger, "Self-interaction correction to density-functional approximations for many-electron systems," *Phys. Rev. B* **23**, 5048 (1981).

Pérez, R., and P. Gumbsch, "Directional Anisotropy in the cleavage fracture of silicon," *Phys. Rev. Lett.* **84**, 5347 (2000).

Perrin, G., and J.R. Rice, "Disordering of a dynamic planar crack in a model elastic medium of randomly variable toughness," *J. Mech. Phys. Solids* **42**, 1047 (1994).

Phillips, R., *Crystals, Defects, and Microstructures: Modeling across Scales* (Cambridge University Press, Cambridge, 2001).

Phoenix, S.L., and I.J. Beyerlein, "Distributions and size scalings for strength in a one-dimensional random lattice with load distribution to nearest and next-nearest neighbors," *Phys. Rev. E* **62**, 1622 (2000).

Phoenix, S.L., and R.L. Smith, "A comparison of probabilistic techniques for the strength of fibrous materials under local load-sharing among fibers," *Int. J. Solids Struct.* **19**, 479 (1983).

Phoenix, S.L., and R. Raj, "Scalings in fracture probabilities for a brittle matrix fiber composite," *Acta Metallurg. Mater.* **40**, 2813 (1992).

Pickard, C.J., B. Winkler, R.K. Chen, M.C. Payne, M.H. Lee, J.S. Lin, J.A. White, V. Milman, and D. Vanderbilt, "Structural properties of Lanthanide and actinide compounds within the plane wave pseudopotential approach," *Phys. Rev. Lett.* **85**, 5122 (2000).

Pietronero, L., and H.J. Wiesmann, "From physical dielectric breakdown to the stochastic fractal model," *Z. Phys. B* **70**, 87 (1988).

Plimpton, S., "Fast parallel algorithms for short-range molecular dynamics," *J. Comput. Phys.* **117**, 1 (1995).

Poirier, C., M. Ammi, D. Bideau, and J.-P. Troadec, "Experimental study of the geometrical effects in the localization of deformation," *Phys. Rev. Lett.* **68**, 216 (1992).

Polanyi, M., *Z. Phys.* **89**, 85 (1934).

Ponte Castañeda, P., "The effective mechanical properties of nonlinear isotropic composites," *J. Mech. Phys. Solids*, **39**, 45 (1991a).

Ponte Castañeda, P., "The effective properties of brittle/ductile incompressible composites," in *Inelastic deformation of composite materials*, edited by G.J. Dvorak (Springer, New York, 1991b), p. 215.

Ponte Castañeda, P., "New variational principles in plasticity and their application to composite materials," *J. Mech. Phys. Solids* **40**, 1757 (1992a).

Ponte Castañeda, P., "Bounds and estimates for the properties of nonlinear heterogeneous systems," *Philos. Trans. R. Soc. London A* **340**, 531 (1992b).

Ponte Castañeda, P., "A new vanational principle and its application to nonlinear heterogeneous systems," *SIAM J. Appl. Math.* **52**, 1321 (1992c).

Ponte Castañeda, P., "Exact second-order estimates for the effective mechanical properties of nonlinear composite materials," *J. Mech. Phys. Solids* **44**, 827 (1996a).

Ponte Castañeda, P., "Variational methods for estimating behaviour of nonlinear composite materials," in *Continuum Models and Discrete systems* (CMDS 8), edited by K.Z. Markov, (World Scientific, Singapore, 1996b), p. 268.

Ponte Castañeda, P., "Nonlinear composite materials: Effective constitutive behavior and microstructure evolution," in *Continuum Micromechanics*, edited by P. Suquet (CISM Courses and Lecture Notes No. 377, Springer, Wien, 1997), p. 131.

Ponte Castañeda, P., "Three-point bounds and other estimates for strongly nonlinear composites," *Phys. Rev. B* **57**, 12077 (1998).

Ponte Castañeda, P., and G. deBotton, "On the homogenized yield strength of two-phase composites," *Proc. R. Soc. Lond. A* **438**, 419 (1992).

Ponte Castañeda, P., G. deBotton, and G. Li, "Effective properties of nonlinear inhomogeneous dielectrics," *Phys. Rev. B* **46**, 4387 (1992).

Ponte Castañeda, P., and M. Kailasam, "Nonlinear electrical conductivity in heterogeneous media," *Proc. R. Soc. Lond. A* **453**, 793 (1997); **453**, 1791(E).

Ponte Castañeda, P., and M.V. Nebozhyn, "Exact second-order estimates of the self-consistent type for nonlinear composite materials," *Mech. Mater.* (1997).

Ponte Castañeda, P., and P. Suquet, "On the effective mechanical behavior of weakly inhomogeneous nonlinear materials," *Europ. J. Mech. A Solids* **14**, 205 (1995).

Ponte Castañeda, P., and P. Suquet, "Nonlinear composites," *Adv. Appl. Mech.* **34** (Academic Press, San Diego, 1998), p. 172.

Ponte Castañeda, P., and J.R. Willis, "On the overall properties of nonlinearly viscous composites," *Proc. R. Soc. Lond. A* **416**, 217 (1988).

Ponte Castañeda, P., and J.R. Willis, "The effective behavior of nonlinear composites: A comparison between two methods," in *Continuum Models and Discrete Systems (CMDS 7)*, edited by K.H. Anthony and H.-H. Wagner (Trans Tech, Aedermannsdorf, Switzerland, 1993), p. 351.

Ponte Castañeda, P., and J.R. Willis, "The effect of spatial distribution on the effective behavior of composite materials and cracked media," *J. Mech. Phys. Solids* **43**, 1919 (1995).

Poon, C.Y., R.S. Sayles, and T.A. Jones, "Surface measurement and fractal characterization of naturally fractured rocks," *J. Appl. Phys.* **25**, 1269 (1992).

Porto, M., S. Havlin, S. Schwarzer, and A. Bunde, "Optimal path in strong disorder and shortest path in invasion percolation with trapping," *Phys. Rev. Lett.* **79**, 4060 (1997).

Prakash, V., and R.J. Clifton, "Experimental and analytical investigation of dynamic fracture under conditions of plane strain," in *Fracture Mechanics: 22nd Symposium*, vol. 1, edited by H.A. Ernst, A. Saxena, and D.L. McDowell (American Society for Testing and Materials, Philadelphia, 1992), p. 412.

Pratt, A.K., and P.L. Green, "Measurement of the dynamic fracture toughness of polymethylmethacrylate by high-speed photography," *Eng. Fract. Mech.* **6**, 71 (1974).

Pulay, P., "Ab initio calculation of force constants and equilibrium geometries in polyatomic molecules. I. Theory," *Mol. Phys.* **17**, 197 (1969).

Qiu, S.-Y., C.Z. Wang, K.M. Ho, and C.T. Chan, "Tight-binding molecular dynamics with linear system-size scaling," *J. Phys.: Condens. Matter* **6**, 9153 (1994).

Qiu, Y.P., and G.J. Weng, "A theory of plasticity for porous materials and particle-reinforced composites," *J. Appl. Mech.* **59**, 261 (1992).

Rahman, A., "Correlation in the motion of atoms in liquid argon," *Phys. Rev.* **136**, A405 (1964).

Rahman, A., and F.H. Stillinger, "Molecular dynamics study of liquid water," *J. Chem. Phys.* **55**, 3336 (1971).

Raine, A.R.C., D. Fincham, and W. Smith, "Systolic loop methods for molecular dynamics simulation using multiple transputers," *Comput. Phys. Commun.* **55**, 13 (1989).

Räisänen, V.I., M.J. Alava, K.J. Niskanen, and R.M. Nieminen, "Does the shear-lag model apply to random fiber networks?" *J. Mater. Res.* **12**, 2725 (1997).

Räisänen, V.I., E.T. Seppala, M.J. Alava, and P.M. Duxbury, "Quasistatic cracks and minimal energy surfaces," *Phys. Rev. Lett.* **80**, 329 (1998).

Rambaldi, S., and O. Pinazza, "An accurate fractional Brownian motion generator," *Physica A* **208**, 21 (1994).

Ramanathan, S., and D.S. Fisher, "Dynamics and instabilities of planar tensile cracks in heterogeneous media," *Phys. Rev. Lett.* **79**, 877 (1997).

Ramanathan, S., and D.S. Fisher, "Onset of propagation of planar cracks in heterogeneous media," *Phys. Rev. B* **58**, 6026 (1998).

Rammal, R., and A.-M.S. Tremblay, "Resistance noise in nonlinear resistor networks," *Phys. Rev. Lett.* **58**, 415 (1987).

Ramulu, M., and A.S. Kobayashi, "Mechanics of crack curving and branching—a dynamic fracture analysis," *Int. J. Fract.* **27**, 187 (1985).

Ramulu, M., and A.S. Kobayashi, "Strain energy density criteria for dynamic fracture and dynamic crack branching," *Theore. Appl. Fract. Mech.* **5**, 117 (1986).

Rapaport, D.C., "Multi-million particle molecular dynamics. III. Design consideration for data-parallel processing," *Comput. Phys. Commun.* **76**, 301 (1993).

Rapaport, D.C., *The Art of Molecular Dynamics* (Cambridge University Press, London, 1995).

Rappe, A., and J.D. Joannopoulos, in *Computer Simulation in Materials Science*, edited by M. Meyer and V. Pontikis, NATO ASI Vol. 205, p. 409 (1991).

Rautiainen, T.T., M.J. Alava, and K. Kaski, "Dynamics of fracture in dissipative systems," *Phys. Rev. E* **51**, R2727 (1998).

Ravi-Chandar, K., and W.G. Knauss, "Dynamic crack-tip stress under stress wave loading—a comparison of theory and experiment," *Int. J. Fract.* **20**, 209 (1982).

Ravi-Chandar, K., and W.G. Knauss, "An experimental investigation into dynamic fracture. I. Crack initiation and arrest," *Int. J. Fract.* **25**, 247 (1984a).

Ravi-Chandar, K., and W.G. Knauss, "An experimental investigation into dynamic fracture: II. Microstructural aspects," *Int. J. Fract.* **26**, 65 (1984b).

Ravi-Chandar, K., and W.G. Knauss, "An experimental investigation into dynamic fracture: III. On steady-state crack propagation and crack branching," *Int. J. Fract.* **26**, 141 (1984c).

Ravi-Chandar, K., and B. Yang, "On the role of microcracks in the dynamic fracture of brittle materials," *J. Mech. Phys. Solids* **45**, 535 (1997).

Ray J.R., "Elastic constants and statistical ensembles in molecular dynamics," *Comput. Phys. Rep.* **8**, 109 (1988).

Ray, P., and B.K. Chakrabarti, "The critical behaviour of fracture properties of dilute brittle solids near the percolation threshold," *J. Phys. C* **18**, L185 (1985a).

Ray, P., and B.K. Chakrabarti, "A microscopic approach to the statistical fracture analysis of disordered brittle solids," *Solid State Commun.* **53**, 477 (1985b).

Ray, P., and B.K. Chakrabarti, "Strength of disordered solids," *Phys. Rev. B* **38**, 715 (1988).

Reidle, J., P. Gumbsch, H.F. Fischmeister, V.G. Glebovsky, and V.N. Semenov, "Fracture studies of tungsten single crystals," *Mater. Lett.* **20**, 311 (1994).

Reuss, A., "Calculation of the flow limits of mixed crystals on the basis of the plasticity of the monocrystals," *Z. Angew. Math. Mech.* **9**, 49 (1929).

Rice, J.R., "Mathematical analysis in the mechanics of fracture," in *Fracture: An Advanced Treatise*, Vol. II, edited by H. Liebowitz (Academic Press, New York, 1968), p. 191.

Rice, J.R., "On the structure of stress-strain relations for time-dependent plastic deformation in metals," *J. Appl. Mech.* **37**, 728 (1970).

Rice, J.R., and G.E. Beltz, "The activation energy for dislocation nucleation at a crack," *J. Mech. Phys. Solids* **42**, 333 (1994).

Rice, J.R., Y. Ben-Zion, and K.S. Kim, "Three dimensional perturbation solution for a dynamic planar crack moving unsteadily in a model elastic solid," *J. Mech. Phys. Solids* **42**, 813 (1994).

Rice, J.R., and R. Thomson, "Ductile versus brittle behaviour of crystals," *Philos. Mag.* **29**, 73 (1974).

Robertson, D.H., D.W. Brenner, and J.W. Mintmire, "Energetics of nanoscale graphitic tubules," *Phys. Rev. B* **45**, 12592 (1992).

Robertson, M.C., C.G. Sammis, M. Sahimi, and A.J. Martin, "Fractal analysis of three-dimensional spatial distributions of earthquakes with a percolation interpretation," *J. Geophys. Res.* **100B**, 609 (1995).

Rodbell, K.P., M.V. Rodriguez, and P.J. Ficalora, "The kinetics of electromigration," *J. Appl. Phys.* **61**, 2844 (1987).

Rodgers, S.T., and K.F. Jensen, "Multiscale modeling of chemical vapor deposition," *J. Appl. Phys.* **83**, 524 (1998).

Ronsin, O., F. Heslot, and B. Perrin, "Experimental study of quasistatic brittle crack propagation," *Phys. Rev. Lett.* **75**, 2252 (1995).

Rosakis, A.J., J. Duffy, and L.B. Freund, "The determination of dynamic fracture toughness of AISI 4340 steel by the shadow spot method," *J. Mech. Phys. Solids* **32**, 443 (1984).

Rosakis, A.J., and L.B. Freund, "The effect of crack tip plasticity on the determination of dynamic stress intensity factors by the optical method of caustics," *J. Appl. Mech.* **48**, 302 (1981).

Rosakis, A.J., O. Samudrala, and D. Coker, "Cracks faster than the shear wave speed," *Science* **284**, 1337 (1999).

Rosakis, A.J., and A.T. Zehnder, "On the dynamic fracture of structural metals," *Int. J. Fract.* **27**, 169 (1985).

Roth, J., F. Gähler, and H.-R. Trebin, "A molecular dynamics run with 5180116000 particles," *Int. J. Mod. Phys. C* **11**, 317 (2000).

Roux, S., and D. Francois, "A simple model for ductile fracture of porous materials," *Scripta Metall.* **25**, 1087 (1991).

Roux, S., A. Hansen, and E. Guyon, "Criticality in non-linear transport properties of heterogeneous materials," *J. Physique* **48**, 2125 (1987).

Roux, S., A. Hansen, H.J. Herrmann, and E. Guyon, "Rupture of heterogeneous media in the limit of infinite disorder," *J. Stat. Phys.* **52**, 251 (1988).

Roux, S., and H.J. Herrmann, "Disorder-induced nonlinear conductivity," *Europhys. Lett.* **4**, 1227 (1987).

Rudd, E., and J.Q. Broughton, "Coarse-grained molecular dynamics and the atomic limit of finite elements," *Phys. Rev. B* **58**, 5893 (1998).

Rundle, J.B., and W. Klein, "Nonclassical nucleation and growth of cohesive tensile cracks," *Phys. Rev. Lett.* **63**, 171 (1989).

Runge, E., and E.K.U. Gross, "Density-functional theory for time-dependent systems," *Phys. Rev. Lett.* **52**, 997 (1984).

Ryckaert, J.P., G. Ciccotti, and H.J.C. Berendsen, "Numerical integration of the cartesian equations of motion of a system with constraints: Molecular dynamics of n-alkane," *J. Comput. Phys.* **23**, 327 (1977).

Ryckaert, J.P., "Special geometrical constraints in the molecular dynamics of chain molecules," *Mol. Phys.* **55**, 549 (1985).

Sachs, G., "Zur Ableitung einer Fleissbedingun," *Z. Ver. Dtsch. Ing.* **72**, 734 (1928).

Safonov, V.P., V.M. Shalaev, V.A. Markel, Y.E. Danilova, N.N. Lepeshkin, W. Kim, S.G. Rautian, and R.L. Armstrong, "Spectral dependence of selective photomodification in fractal aggregates of colloidal particles," *Phys. Rev. Lett.* **80**, 1102 (1998).

Sahimi, M., "Nonlinear transport processes in disordered media," *AIChE J.* **39**, 369 (1993a).

Sahimi, M., "Flow phenomena in rocks: From continuum models to fractals, percolation, cellular automata and simulated annealing," *Rev. Mod. Phys.* **65**, 1395 (1993b).

Sahimi, M., *Applications of Percolation Theory* (Taylor and Francis, London, 1994a).

Sahimi, M., "Long-range correlated percolation and flow and transport in heterogeneous porous media," *J. Physique I France* **4**, 1263 (1994b).

Sahimi, M., "Effect of long-range correlations on transport phenomena in disordered media," *AIChE J.* **41**, 229 (1995a).

Sahimi, M., *Flow and Transport in Porous Media and Fractured Rock* (VCH, Weinheim, Germany, 1995b).

Sahimi, M., "Non-linear and non-local transport processes in heterogeneous media: From long-range correlated percolation to fracture and materials breakdown," *Phys. Rep.* **306**, 295 (1998).

Sahimi, M., and S. Arbabi, "Force distribution, multiscaling, and fluctuations in disordered elastic media," *Phys. Rev. B* **40**, 4975 (1989).

Sahimi, M., and S. Arbabi, "On correction to scaling for two- and three-dimensional scalar and vector percolation," *J. Stat. Phys.* **62**, 453 (1991).

Sahimi, M., and S. Arbabi, "Percolation and fracture in disordered solids and granular media: Approach to a fixed point," *Phys. Rev. Lett.* **68**, 608 (1992).

Sahimi, M., and S. Arbabi, "Mechanics of disordered solids. III. Fracture properties," *Phys. Rev. B* **47**, 713 (1993).

Sahimi, M., and S. Arbabi, "Scaling laws for fracture of heterogeneous materials and rock," *Phys. Rev. Lett.* **77**, 3689 (1996).

Sahimi, M., and J.D. Goddard, "Elastic percolation models for cohesive mechanical failure in heterogeneous systems," *Phys. Rev. B* **33**, 7848 (1986).

Sahimi, M., M.C. Robertson, and C.G. Sammis, "Fractal distribution of earthquake hypocenters and its relation with fault patterns and percolation," *Phys. Rev. Lett.* **70**, 2186 (1993).

Sammonds, P.R., P.G. Meredith, M.R. Ayling, and S.A. Murell, in *Fracture of Concrete and Rock*, edited by S.P. Shah, S.E. Swartz, and B. Barr (Elsevier, London, 1989), p. 101.

Sánchez, A., F. Guinea, E. Louis, and V. Hakim, "On the fractal characteristics of the eta model," *Physica A* **191**, 123 (1992).

Sanchez-Palancia, E., *Non-homogeneous Media and Vibration Theory* (Lecture notes in Physics) **127** (Springer, Heidelberg, 1980), p. 12.

Sanz-Serna, J.M., "Symplectic integration for Hamiltonian systems: An overview," *Acta Numerica* **1**, 243 (1992).

Sarychev, A.K., D.J. Bergman, and Y. Yagil, "Optical and microwave properties of metal-insulator thin films: possibility of light localization," *Physica A* **207**, 372 (1994).

Sarychev, A.K., D.J. Bergman, and Y. Yagil, "Theory of the optical and microwave properties of metal-dielectric films," *Phys. Rev. B* **51**, 5366 (1995).

Sarychev, A.K., and V.M. Shalaev, "Electromagnetic field fluctuations and optical nonlinearities in metal-dielectric composites," *Phys. Rep.* **335**, 275 (2000).

Sarychev, A.K., and V.M. Shalaev, "Giant high-order field moments in metal-dielectric composites," *Physica A* **266**, 115 (1999).

Sarychev, A.K., V.A. Shubin, and V.M. Shalaev, "Anderson localization of surface plasmons and nonlinear optics of metal-dielectric composites," *Phys. Rev. B* **60**, 16389 (1999).

Sawada, Y., S. Ohta, M. Yamazaki, and H. Honjo, "Self-similarity and a phase-transition-like behavior of a random growing structure governed by a nonequilibrium parameter," *Phys. Rev. A* **26**, 3557 (1982).

Schardin, H., "Velocity effects in fracture," in *Fracture,* edited by B.L. Averbach (MIT Press, Cambridge, MA, 1959), p. 297.

Schimschak, M., and J. Krug, "Electromigration-induced breakup of two-dimensional voids," *Phys. Rev. Lett.* **80**, 1674 (1998).

Schmittbuhl, J., S., Gentier, and S. Roux, S., "Field measurements of the roughness of fault surfaces," *Geophys. Res. Lett.* **20**, 639 (1993).

Schmittbuhl, J., S. Roux, J.-P. Vilotte, and K.J. Måløy, "Interfacial crack pinning: effect of nonlocal interactions," *Phys. Rev. Lett.* **74**, 1787 (1995).

Schöen, M., "Structure of a simple molecular dynamics. II. Design considerations for distributed processing," *Comput. Phys. Commun.* **52**, 175 (1989).

Schwartz, L.M., S. Feng, M.F. Thorpe, and P.N. Sen, "Behavior of depleted elastic networks: Comparison of effective-medium and numerical calculations," *Phys. Rev. B* **32**, 4607 (1985).

Scott, I.G., "Basic Acoustic Emission," in *Nondestructive Testing Monographs and Tracts*, Vol. 6 (Gordon and Breach New York, 1991), p. 124.

Shalaev, V.M., R. Botet, and A.V. Butenko, "Localization of collective dipole excitations on fractals," *Phys. Rev. B* **48**, 6662 (1993).

Shalaev, V.M., R. Botet, J. Mercer, and E.B. Stechel, "Optical properties of self-affine thin films," *Phys. Rev. B* **54**, 8235 (1996a).

Shalaev, V.M., C. Douketis, T. Haslett, T. Stuckless, and M. Moskovits, "Two-photon electron emission from smooth and rough metal films in the threshold region," *Phys. Rev. B* **53**, 11193 (1996b).

Shalaev, V.M., V.A. Markel, E.Y. Poliakov, R.L. Armstrong, V.P. Safonov, A.K. Sarychev, "Nonlinear optical phenomena in nanostructured fractal materials," *J. Nonlin. Optic. Phys. & Mat.* **7**, 131 (1998).

Shalaev, V.M., and A.K. Sarychev, "Nonlinear optics of random metal-dielectric films," *Phys. Rev. B* **57**, 13265 (1998).

Sharon, E., and J. Fineberg, "Microbranching instability and the dynamic fracture of brittle materials," *Phys. Rev. B* **54**, 7128 (1996).

Sharon, E., and J. Fineberg, "Universal features of the micro-branching instability in dynamic fracture," *Philos. Mag. B* **78**, 243 (1998).

Sharon, E., and J. Fineberg, "On massless cracks and the continuum theory of fracture," *Nature* **397**, 333 (1999).

Sharon, E., S.P. Gross, and J. Fineberg, "Local crack branching as a mechanism for instability in dynamic fracture," *Phys. Rev. Lett.* **74**, 5096 (1995).

Sharon, E., S.P. Gross, and J. Fineberg, "Energy dissipation in dynamic fracture," *Phys. Rev. Lett.* **76**, 2117 (1996).

Sharp, S.J., M.F. Ashby, and N.A. Fleck, "Material response under static and sliding indentation loads," *Acta Met.* **41**, 685 (1993).

Shenoy, V.B., R. Miller, E.B. Tadmor, R. Phillips, and M. Ortiz, "Quasicontinuum models of interfacial structure and deformation," *Phys. Rev. Lett.* **80**, 742 (1998).

Shenoy, V.B., R. Miller, E.B. Tadmor, D. Rodney, and R. Phillips, "An adaptive finite element approach to atomic-scale mechanics—the quasicontinuum method," *J. Mech. Phys. Solids* **47**, 611 (1999).

Shenoy, V.B., R. Phillips, and E.B. Tadmor, "Nucleation of dislocation beneath a plane strain indenter," *J. Mech. Phys. Solids* **48**, 649 (2000).

Shimojo, F., I. Ebbsjö, R.K. Kalia, A. Nakano, J.P. Rino, and P. Vashishta, "Molecular dynamics simulation of structural transformation in silicon carbide under pressure," *Phys. Rev. Lett.* **84**, 3338 (2000).

Shinkai, N., "Fracture and fractography of flat glass," in *Fractography of Glass*, edited by R.C. Bradt and R.E. Tressler (Plenum, New York, 1994), p. 253.

Shioya, T. and R. Ishida, "Microscopic fracture modes of brittle polymers in dynamic crack propagation," in *Dynamic Failure of Materials*, edited by H.P. Rossmanith and A.J. Rosakis (Elsevier Applied Science, Essex, 1991), p. 351.

Shirley, E.L., D.C. Allan, R.M. Martin, and J.D. Joannopoulos, "Extended norm-conserving pseudopotentials," *Phys. Rev. B* **40**, 3652 (1989).

Sieradzki, K., G.J. Dienes, A. Paskin, and B. Massoumzadeh, "Atomistics of crack propagation," *Acta Mettal.* **36**, 651 (1988).

Sieradzki, K., and R. Li, "Fracture behavior of a solid with random porosity," *Phys. Rev. Lett.* **56**, 2509 (1986).

Sieradzki, K., and R.C. Newman, "Brittle behaviour of ductile metals during stress-corrosion cracking," *Philos. Mag. A* **51**, 95 (1985).

Sig, G.C., "Some basic problems in fracture mechanics and new concepts," *Eng. Fract. Mech.* **5**, 365 (1973).

Sinclair, J.E., and B.R. Lawn, "An atomistic study of cracks in diamond-structure crystals," *Proc. R. Soc. Lond. A* **329**, 83 (1972).

Skjeltorp, A.T., and P. Meakin, "Fracture in microsphere monolayers studied by experiment and computer simulation," *Nature* **335**, 424 (1988).

Slepyan, L., "Dynamics of a crack in a lattice," *Sov. Phys. Dokl.* **26**, 538 (1981).

Sloan, S.W., "A fast algorithm for constructing Delaunay triangulations in the plane," *Adv. Engng. Software*, **9**, 34 (1987).

Smekal, E., "Zum Bruchvorgang bei sprodem stoffrerhalten unter ein und mehrachsinen beanspruchungen," *Osterr. Ing. Arch.* **7**, 49 (1953).

Smith, R.L., S.L. Phoenix, M.R. Greenfield, R.B. Henstenburg, and R.E. Pitt, "Lower-tail approximations for the probability of failure of three-dimensional fibrous composites with hexagonal geometry," *Proc. R. Soc. Lond. A* **388**, 353 (1983).

Smith, W., and T.R. Forester, "Parallel macromolecular simulations and the replicated data strategy II: The RD-SHAKE algorithm," *Comput. Phys. Commun.* **79**, 63 (1994).

Smyshlyaev, V.P., and N.A. Fleck, "Bounds and estimates for the overall plastic behavior of composites with strain gradients effects," *Proc. R. Soc. Lond. A* **451**, 795 (1995).

Söderberg, M., "Resistive breakdown of inhomogeneous media," *Phys. Rev. B* **35**, 352 (1987).

Sornette, D., "Weibull-like failure distribution induced by fluctuations in percolation," *J. Physique* **49**, 889 (1988).

Sornette, D., "Elasticity and failure of a set of elements loaded in parallel," *J. Phys. A* **22**, L243 (1989).

Sornette, D., and C. Vanneste, "Dynamics and memory effects in rupture of thermal fuse networks," *Phys. Rev. Lett.* **68**, 612 (1992).

Sornette, D., and C. Vanneste, "Dendrites and fronts in a model of dynamical rupture with damage," *Phys. Rev. E* **50**, 4327 (1994).

Soules, T.F., and R.F. Busbey, "The rheological properties and fracture of a molecular dynamic simulation of sodium silicate glass," *J. Chem. Phys.* **78**, 6307 (1983).

Spence, J.C.H., Y.M. Huang, and O. Sankey, "Lattice trapping and surface reconstruction for silicon cleavage on (111). *Ab-initio* quantum molecular dynamics calculations," *Acta Mettal. Mater.* **41**, 2815 (1993).

Spitzig, W.A., R.E. Smelser, and O. Richmond, "The evolution of damage and fracture in iron compacts with various initial porosities," *Acta Metall. Mater.* **36**, 1201 (1998).

Sridhar, N., W. Yang, D.J. Srolovitz, and E.R. Fuller, Jr., "Microstructral mechanics model of anisotropic-thermal-expansion-induced microcracking," *J. Am. Ceramic Soc.* **77**, 1123 (1994).

Srolovitz, D.J., and P.D. Beale, "Computer simulation of failure in an elastic model with randomly distributed defects," *J. Amer. Cer. Soc.* **71**, 362 (1988).

Stadler, J., R. Mikulla, and H.-R. Trebin, "IMD: a software package for molecular dynamics studies on parallel computers," *Int. J. Mod. Phys. C* **8**, 1131 (1997).

Stanley, H.E., and P. Meakin, "Multifractal phenomena in physics and chemistry," *Nature* **335**, 405 (1988).

Starkloff, T., and J.D. Joannopoulos, "Local pseudopotential theory for transition metals," *Phys. Rev. B* **16**, 5212 (1977).

Stauffer, D., and A. Aharony, *Introduction to Percolation Theory*, 2nd ed. (Taylor and Francis, London, 1992).

Steele, W.A., "The physical interaction of gases with crystalline solids. I. Gas-solid energies and properties of isolated adsorbed atoms," *Surf. Sci.* **36**, 317 (1973).

Stephens, M.D., and M. Sahimi, "Distribution of fracture strengths in disordered continua," *Phys. Rev. B* **36**, 8656 (1987).

Stinchcombe, R.B., "Conductivity and spin-wave stiffness in disordered systems-an exactly soluble model," *J. Phys. C* **7**, 179 (1974).

Stinchcombe, R.B., P.M. Duxbury, and P. Shukla, "The minimum gap on diluted Cayley trees," *J. Phys. A* **19**, 3903 (1986).

Stillinger, F.H., and T.A. Weber, "Computer simulation of local order in condensed phases of silicon," *Phys. Rev. B* **31**, 5262 (1985).

St. John, C., "The brittle-to-ductile transition in pre-cleaved silicon single crystals," *Philos. Mag.* **32**, 1193 (1975).

Stockman, M.I., "Inhomogeneous eigenmode localization, chaos, and correlations in large disordered clusters," *Phys. Rev. E* **56**, 6494 (1997).

Stockman, M.I., L.N. Pandey, and T.F. George, "Inhomogeneous localization of polar eigenmodes in fractals," *Phys. Rev. B* **53**, 2183 (1996).

Stockman, M.I., L.N. Pandey, L.S. Muratov, and T.F. George, "Giant fluctuations of local optical fields in fractal clusters," *Phys. Rev. Lett.* **72**, 2486 (1994).

Stockman, M.I., L.N. Pandey, L.S. Muratov, and T.F. George, "Optical absorption and localization of eigenmodes in disordered clusters," *Phys. Rev. B* **51** 185 (1995)

Stoddard, S.D., and J. Ford, "Numerical experiments on the stochastic behavior of a Lennard–Jones gas system," *Phys. Rev. A* **8**, 1504 (1973).

Straley, J.P., "Random resistor tree in an applied field," *J. Phys. C* **10**, 3009 (1977).

Straley, J.P., and S.W. Kenkel, "Percolation theory for nonlinear conductors," *Phys. Rev. B* **29**, 6299 (1984).

Strang, G., and G.J. Fix, *An Analysis of the Finite Element Method* (Prentice-Hall, Englewood Cliffs, 1973).

Stroh, A.N., "A theory of the fracture of metals," *Adv. Phys.* **6**, 418 (1957).

Stroud, D., and P.M. Hui, "Nonlinear susceptibilities of granular matter," *Phys. Rev. B* **37**, 8719 (1988).

Stroud, D., and V.E. Wood, "Decoupling approximation for the nonlinear-optical response of composite media," *J. Opt. Soc. Am.* **B6**, 778 (1989).

Stroud, D., and X. Zhang, "Cubic nonlinearities in small-particle composites: local-field induced giant enhancements," *Physica A* **207**, 55 (1994).

Suquet, P., "Analyse limite et homogénéisation," *C. R. Acad. Sci. Paris Ser. II* **296**, 1355 (1983).

Suquet, P., "Elements of homogenization for inelastic solid mechanics," in *Homogenization Techniques for Composite Media*, edited by E. Sanchez-Palencia and A. Zaoui, Lecture Notes in Physics 272. (Springer, Berlin, 1987), p. 193.

Suquet, P., "On bounds for the overall potential of power-law materials containing voids with arbitrary shape," *Mech. Res. Comm.* **19**, 51 (1992).

Suquet, P., "Overall potentials and extremal surfaces of power law or ideally plastic materials," *J. Mech. Phys. Solids* **41**, 981 (1993a).

Suquet, P., "Bounds and estimates for the overall properties of nonlinear composites," in *Micromechanics of Materials*, edited by J.J. Marigo and G. Rousselier (Eyrolles, Paris, 1993b), p. 361.

Suquet, P., "Overall properties of nonlinear composites: A modified secant moduli theory and its link with Ponte Castañeda's nonlinear variational procedure," *C. R. Acad. Sci. Paris Ser. II* **320**, 563 (1995).

Suquet, P., "Effective properties of nonlinear composites," in *Continuum Micromechanics*, edited by P. Suquet, CISM Courses and Lecture Notes No. 377 (Springer, Wien, 1997), p. 197.

Suquet, P., and P. Ponte Castañeda, "Small-contrast perturbation expansions for the effective properties of nonlinear composites," *C. R. Acad Sci. Paris Ser. II* **317**, 1515 (1993).

Suwanmethanond, V., E. Goo, P.K.T. Liu, G. Johnson, M. Sahimi, and T.T. Tsotsis, "Porous silicon carbide sintered substrates for high-temperature membranes," *Ind. Eng. Chem. Res.* **39**, 3264 (2000).

Swope, W., and H.C. Andersen, "$10^6$-Particle molecular dynamics study of homogeneous nucleation of crystals in a supercooled atomic liquid," *Phys. Rev. B* **41**, 7042 (1990).

Tadmor, E.B., M. Ortiz, and R. Phillips, "Quasicontinuum analysis of defects in solids," *Philos. Mag. A* **73**, 1529 (1996a).

Tadmor, E.B., R. Phillips, and M. Ortiz, "Mixed atomistic and continuum models of deformation in solids," *Langmuir* **12**, 4529 (1996b).

Tadmor, E.B., R. Miller, R. Phillips, and M. Ortiz, "Nanoindentation and incipient plasticity," *J. Mater. Res.* **14**, 2233 (1999).

Takayasu, H., "Simulation of electric breakdown and resulting variant of percolation fractals," *Phys. Rev. Lett.* **54**, 1099 (1985).

Talbot, D.R.S., and J.R. Willis, "Variational principles for inhomogeneous nonlinear media," *IMA J. Appl. Math.* **35**, 39 (1985).

Talbot, D.R.S., and J.R. Willis, "Bounds and self-consistent estimates for the overall properties of nonlinear composites," *IMA J. Appl. Math.* **39**, 215 (1987).

Talbot, D.R.S., and J.R. Willis, "Some explicit bounds for the overall behaviour of nonlinear composites," *Int. J. Solids Struct.* **29**, 1981 (1992).

Talbot, D.R.S. and J.R. Willis, "Upper and lower bounds for the overall properties of a nonlinear composite dielectric I. Random microgeometry," *Proc. R. Soc. Lond. A* **447**, 365 (1994).

Talbot, D.R.S., and J.R. Willis, "Three-point bounds for the overall properties of a nonlinear composite dielectric," *IMA J. Appl. Math.* **57**, 41 (1996).

Tandon, G.P., and G.J. Weng, "A theory of particle-reinforced plasticity," *J. Appl. Mech.* **55**, 126 (1988).

Taylor, G.I., "The mechanism of plastic deformation of crystals. Part I.—Theoretical," *Proc. R. Soc. Lond. A* **145**, 362 (1934).

Taylor, G.I., "Plastic strains in metals," *J. Inst. Metals* **62**, 307 (1938).

Termonia, Y., and P. Smith, "Theoretical study of the ultimate mechanical properties of poly(*p*-phenylene-terephthalamide) fibres," *Polymer* **27**, 1845 (1986).

Termonia, Y., and P. Meakin, "Formation of fractal cracks in a kinetic fracture model," *Nature* **320**, 429 (1986).

Termonia, Y., P. Meakin, and P. Smith, "Theoretical study of the molecular weight on the maximum tensile strength of polymer fibre," *Macromolecules* **18**, 2246 (1985).

Termonia, Y., P. Meakin, and P. Smith, "Theoretical study of the influence of strain rate and temperature on the maximum strength of perfectly ordered and oriented polyethylene," *Macromolecules* **19**, 154 (1986).

Termonia, Y., and P. Smith, "Kinetic model for tensile deformation of polymers. 1. Effect of Molecular Weight," *Macromolecules* **20**, 835 (1987).

Termonia, Y., and P. Smith, "Kinetic Model for tensile deformation of polymers. 2. Effect of entanglement spacing," *Macromolecules* **21**, 2184 (1988).

Tersoff, J., "New empirical approach for the structure and energy of covalent systems," *Phys. Rev. B* **37**, 6991 (1988).

Tersoff, J., "Modelling solid-state chemistry: Interatomic potentials for multicomponent systems," *Phys. Rev. B* **39**, 5566 (1989).

Theocaris, P.S., and E.E. Gdoutos, "An optical method for determining opening-mode and edge sliding-mode stress intensity factors," *J. Appl. Mech.* **39**, 91 (1972).

Theocaris, P.S., and H.G. Georgiadis, "Bifurcation predictions for moving cracks by the T-criterion," *Int. J. Fract.* **29**, 181 (1985).

Thijssen, J.M., *Computational Physics* (Cambridge University Press, Cambridge, 1999).

Thomson, R., "The physics of fracture," *Solid State Phys.* **39**, 1 (1986).

Thomson, R., C. Hsieh, and V. Rana, "Lattice trapping of fracture cracks," *J. Appl. Phys.* **42**, 3154 (1971).

Tokatly, I.V., and O. Pankratov, "Many-body diagrammatic expansion in a Kohn–Sham basis: Implications for time-dependent density functional theory of excited states," *Phys. Rev. Lett.* **86**, 2078 (2001).

Torquato, S., "Electrical conductivity of two-phase disordered composite media," *J. Appl. Phys.* **58**, 3790 (1985a).

Torquato, S., "Bulk properties of two-phase disordered media. II. Effective conductivity of a dilute suspension of penetrable spheres," *J. Chem. Phys.* **83**, 4776 (1985b).

Torrie, G.M., and J.P. Valleau, "Nonphysical sampling distributions in Monte Carlo free-energy estimation: Umbrella sampling," *J. Comput. Phys.* **23**, 157 (1977).

Tsuruta, K., A. Omeltchenko, R.K. Kalia, and P. Vashishta, "Early stages of sintering of silicon nitride nanoclusters: a molecular-dynamics study on parallel machines," *Europhys. Lett.* **33**, 441 (1996).

Tua, P.F., and J. Bernasconi, "Monte Carlo simulations of two-dimensional randomly diluted networks of nonlinear resistors," *Phys. Rev. B* **37**, 1986 (1988).

Tvergaard, V., "Analysis of tensile properties for a whisker-reinforced metal-matrix composite," *Acta Metall. Mater.* **38**, 185 (1990).

Tzschichholz, F., "Peeling instability in Cosserat-like media," *Phys. Rev. B* **45**, 12691 (1992).

Tzschichholz, F., "Fracturing of brittle homogeneous solids: finite-size scalings," *Phys. Rev. B* **52**, 9270 (1995).

Tzschichholz, F., and H.J. Herrmann, "Reaction-diffusion model for the hydration and setting of cement," *Phys. Rev. E* **53**, 2629 (1996).

Tzschichholz, F., H.J. Herrmann, H.E. Roman, and M. Puff, "Beam model for hydraulic fracturing," *Phys. Rev. B* **49**, 7056 (1994).

Uenoyama, T., L. Esaki, and H. Kotera, "Theory of stability in a nonlinear resistive network," *Appl. Phys. Lett.* **61**, 363 (1992).

Underwood, E.E., and K. Banerji, "Fractals in fractography," *Mater. Sci. Eng.* **80**, 1 (1986).

Van Damme, H., F. Obrecht, P. Levitz, L. Gatineau, and C. Laroche, "Fractal viscous fingering in clay slurries," *Nature* **320**, 731 (1986).

Van Damme, H., C. Laroche, and L. Gatineau, "Radial fingering in viscoelastic media, an experimental study," *Revue Phys. Appl.* **22**, 241 (1987a).

Van Damme, H., C. Laroche, L. Gatineau, and P. Levitz, "Viscoelastic effects in fingering between miscible fluids," *J. Physique* **48**, 1121 (1987b).

van den Born, I.C., A. Santen, H.D. Hoekstra, and J.Th.M. De Hosson, "Mechanical strength of highly porous ceramics," *Phys. Rev. B* **43**, 3794 (1991).

Vanderbilt, D., "Soft self-consistent pseudopotentials in a generalized eigenvalue formalism," *Phys. Rev. B* **41**, 7892 (1990).

Vanneste, C., and D. Sornette, "The dynamical thermal fuse model," *J. Phys. France I* **2**, 1621 (1992).

Verges, J.A., "Localization length in a random magnetic field," *Phys. Rev. B* **57**, 870 (1998).

Verlet, L., ""Computer "experiments" on classical fluids. I. Thermodynamical properties of Lennard–Jones molecules," *Phys. Rev.* **159**, 98 (1967).

Voigt, W., "Ueber die Bezienhung zwischen den beiden Elasticitäts-constanten isotroper," *Ann. Physik* **38**, 573 (1889).

Vold, M.J., "Computer simulation of floc formation in a colloidal suspension," *J. Colloid Inter. Sci.* **18**, 684 (1963).

Voss, R.F., "Random fractal forgeries," in *Fundamental Algorithms for Computer Graphics*, edited by R.A. Earnshaw, NATO ASI Series, Vol. 17 (Springer-Verlag, Heidelberg, 1985), p. 805.

Vu, B.Q., and V.K. Kinra, "Brittle fracture of plates in tension: static field radiated by a suddenly stopping crack," *Eng. Fract. Mech.* **15** 107 (1981).

Wagner, N.J., B.L. Holian, and A.F. Voter, "Molecular-dynamics simulations of two-dimensional materials at high strain rates," *Phys. Rev. A* **45**, 8457 (1992).

Wan, W.M.V., H.C. Lee, P.M. Hui, and K.W. Yu, "Mean-field theory of strongly nonlinear random composites: Strong power-law nonlinearity and scaling behavior," *Phys. Rev. B* **54**, 3946 (1996).

Wang, Z.G., D.L. Chen, X.X. Jiang, S.H. Ai, and C.H. Shih, "Relationship between fractal dimension and fatigue threshold value in dual-phase steels," *Scripta Metall.* **22**, 827 (1988).

Wang, Z.-G., U. Landman, R.L. Blumberg Selinger, and W.A. Gelbart, "Molecular-dynamics study of elasticity and failure of ideal solids," *Phys. Rev. B* **44**, 378 (1991).

Washabaugh, P.D., and W.G. Knauss, "A reconciliation of dynamic crack velocity and Rayleigh wave speed in isotropic brittle solids," *Int. J. Fract.* **65**, 97 (1994).

Webb, III, E.B., and G.S. Grest, "Liquid/vapor surface tension of metals: Embedded atom method with charge gradient corrections," *Phys. Rev. Lett.* **86**, 2066 (2001).

Weichert, R., and K. Schonert, "On the temperature at the tip of a fast running crack," *J. Mech. Phys. Solids* **22**, 127 (1974).

Weiner, J.H., and M. Pear, "Crack and dislocation propagation in an idealized crystal model," *J. Appl. Phys.* **46**, 2398 (1975).

Whitehead, S., *Dielectric Breakdown of Solids* (Clarendon, Oxford, 1951).

Wiener, O., "Die theorie des mischkʻorpers fȑ das feld des stationären strömung," *Math.-Physichen Klasse der Königl. Sächsischen Gesellschaft der Wissenschaften* **32**, 509 (1912).

Wiesmann, H.J., and H.R. Zeller, "A fractal model of dielectric breakdown and prebreakdown in solid dielectrics," *J. Appl. Phys.* **60**, 1770 (1986).

Williford, R.E., "Scaling similarities between fracture surfaces, energies, and a structure parameter," *Scripta Metall.* **22**, 197 (1988).

Willis, J.R., "Bounds and self-consistent estimates for the overall moduli of anisotropic composites," *J. Mech. Phys. Solids* **25**, 185 (1977).

Willis, J.R., "Variational principles and bounds for the overall properties of composites," in *Continuum Models and Discrete Systems (CMDS 2)*, edited by J. Provan (University of Waterloo Press, Waterloo, Canada, 1978), p. 185.

Willis, J.R., "Variational and related methods for the overall properties of composites," *Adv. Appl. Mech.* **21**, 1 (1981).

Willis, J.R., "The overall response of composite materials," *ASME J. Appl. Mech.* **50**, 1202 (1983).

Willis, J.R., "Variational estimates for the overall response of an inhomogeneous nonlinear dielectric," in *Homogenization and Effective Moduli of Materials and Media*, edited by J.L. Ericksen, D. Kinderlehrer, R. Kohn, and J.-L. Lions (Springer, New York, 1986), p. 247.

Willis, J.R., "The structure of overall constitutive relations for a class of nonlinear composites," *IMA J. Appl. Math.* **43**, 231 (1989a).

Willis, J.R., in *Micromechanics and Inhomogeniety, the Toshio Mura Anniversary Volume*, edited by G.J. Weng, M. Taya, and H. Abe (Springer, New York, 1989b), p. 581.

Willis, J.R., *Elasticity: Mathematical Methods and Applications* (Halston Press, New York, 1990), p. 397.

Willis, J.R., "On methods for bounding the overall properties of nonlinear composites," *Phys. Solids* **39**, 73 (1991).

Willis, J.R., "On method for bounding the overall properties of nonlinear composites: Correction and addition," *J. Mech. Phys. Solids* **40**, 441 (1992).

Willis, J.R., and A.B. Movchan, "Dynamic weight functions for a moving crack. I. Mode I loading," *J. Mech. Phys. Solids* **43**, 319 (1995).

Willis, J.R., and A.B. Movchan, "Three dimensional dynamic perturbation of a propagating crack," *J. Mech. Phys. Solids* **45**, 591 (1997).

Winkler, S., D.A. Shockey, and D.R. Curran, "Crack propagation at supersonic velocities. I," *Int. J. Fract.* **6**, 151 (1970).

Witten, T.A., and L.M. Sander, "Diffusion-limited aggregation, a kinetic critical phenomenon," *Phys. Rev. Lett.* **47**, 1400 (1981).

Wood, W.W., and F.R. Parker, "Monte Carlo equation of state of molecules interacting with the Lennard–Jones Potential. I. A supercritical isotherm at about twice the critical temperature," *J. Chem. Phys.* **27**, 720 (1957).

Wright, D.C., D.J. Bergman, and Y. Kantor, "Resistance fluctuations in random resistor networks above and below the percolation threshold," *Phys. Rev. B* **33**, 396 (1986).

Wu, B.Q., and P.L. Leath, "Failure probabilities and tough-brittle crossover of heterogeneous materials with continuous disorder," *Phys. Rev. B* **59**, 4002 (1999).

Wu, B.Q., and P.L. Leath, "Fracture strength of one-dimensional systems with continuous disorder: A single-crack approximation," *Phys. Rev. B* **61**, 15028 (2000).

Wu, K., and R.M. Bradley, "Theory of electromigration failure in polycrystalline metal films," *Phys. Rev. B* **50**, 12468 (1994).

Xu, G., A.S. Argon, and M. Ortiz, "Nucleation of dislocation from crack tips under mixed modes of loading: Implications for brittle against ductile behaviour of crystals," *Philos. Mag.* **72**, 415 (1995).

Xu, L., M. Sahimi, and T.T. Tsotsis, "Nonequilibrium molecular dynamics simulations of transport and separation of gas mixtures in nanoporous materials," *Phys. Rev. E* **62**, 6942 (2000b).

Xu, L., M.G. Sedigh, M. Sahimi, and T.T. Tsotsis, "Nonequilibrium molecular dynamics simulation of transport of gas mixtures in nanopores," *Phys. Rev. Lett.* **80**, 3511 (1998).

Xu, L., T.T. Tsotsis, and M. Sahimi, "Nonequilibrium molecular dynamics simulation of transport and separation of gases in carbon nanopores. I. Basic results," *J. Chem. Phys.* **111**, 3252 (1999).

Xu, L., M.G. Sedigh, T.T. Tsotsis, and M. Sahimi, "Nonequilibrium molecular dynamics simulation of transport and separation of gases in carbon nanopores. II. Binary and ternary mixtures and comparison with the experimental data," *J. Chem. Phys.* **112**, 910 (2000a).

Xu, X.P., and A. Needleman, "Numerical simulations of fast crack growth in brittle solids," *J. Mech. Phys. Solids* **42** 1397(1994).

Yagil, Y., G. Deutscher, and D.J. Bergman, "Electrical breakdown measurements of semicontinuous metal films," *Phys. Rev. Lett.* **69**, 1423 (1992).

Yagil, Y., G. Deutscher, and D.J. Bergman, "The role of microgeometry in the electrical breakdown of metal-insulator mixtures," *Int. J. Mod. Phys. B* **7**, 3353 (1993).

Yagil, Y., G. Deutscher, and D.J. Bergman, "Nonlinear electrical response and breakdown of semicontinuous metal films," *Physica A* **207**, 323 (1994).

Yagil, Y., P. Gadanne, C. Julien, and G. Deutscher, "Optical properties of thin semicontinuous gold films over a wavelength range of 2.5 to 500 $\mu$m," *Phys. Rev. B* **46**, 2503 (1992).

Yan, Y., G. Li, and L.M. Sander, "Fracture growth in 2d elastic networks with Born model," *Europhys. Lett.* **10**, 7 (1989).

Yang, C.S., and P.M. Hui, "Effective nonlinear response in random nonlinear resistor networks: numerical studies," *Phys. Rev. B* **44**, 12559 (1991).

Yoffe, E.H., "The moving Griffith crack," *Philos. Mag. (series 7)* **42**, 739 (1951).

Yonenaga, I., and K. Sumino, "Mechanical properties and dislocation dynamics of GaP," *J. Mater. Res.* **4**, 355 (1989).

Young, A.P., "Melting and the vector Coulomb gas in two dimensions," *Phys. Rev. B* **19**, 1855 (1979).

Yoon, Y.-G., M.S.C. Mazzoni, H.J. Choi, J. Ihm, and S.G. Louie, "Structural deformation and intertube conductance of crossed carbon nanotube junctions," *Phys. Rev. Lett.* **86**, 688 (2001).

Yu, K.W., and G.Q. Gu, "Electrostatic boundary-value problems of nonlinear media: a perturbation approach," *Phys. Lett. A* **168**, 313 (1992).

Yu, K.W., and G.Q. Gu, "Effective conductivity of nonlinear composites. II. Effective-medium approximation," *Phys. Rev. B* **47**, 7568 (1993).

Yu, K.W., and G.Q. Gu, "Variational calculation of strongly nonlinear composites," *Phys. Lett. A* **193**, 311 (1994).

Yu, K.W., and G.Q. Gu, "Effective conductivity of strongly nonlinear composites: Variational approach," *Phys. Lett. A* **205**, 295 (1995).

Yu, K.W., and P.M. Hui, "Percolation effects in two-component nonlinear composites: Crossover from linear to nonlinear behavior," *Phys. Rev. B* **50**, 13327 (1994).

Yu, K.W., P.M. Hui, and D. Stroud, "Effective dielectric response of nonlinear composites," *Phys. Rev. B* **47**, 14150 (1993).

Yuse, A., and M. Sano, "Transition between crack patterns in quenched glass plates," *Nature* **362** 329 (1993).

Zapperi, S., P. Ray, H.E. Stanley, and A. Vespignani, "First-order transition in the breakdown of disordered media," *Phys. Rev. Lett.* **78**, 1408 (1997).

Zehnder, A.T., and A.J. Rosakis, "On the temperature distribution at the vicinity of dynamically propagating cracks in 4340 steel," *J. Mech. Phys. Solids* **29**, 385 (1991).

Zeng, X.C., D.J. Bergman, P.M. Hui, and D. Stroud, "Effective-medium theory for weakly nonlinear composites," *Phys. Rev. B* **38**, 10970 (1988).

Zeng, X.C., P.M. Hui, D.J. Bergman, and D. Stroud, "Mean field theory for weakly nonlinear composites," *Physica A* **157**, 192 (1989).

Zepeda-Ruiz, L.A., D. Maroudas, and W.H. Weinberg, "Theoretical study of the energetics, strain fields, and semicoherent interface structures in layer-by-layer semiconductor heteroepitaxy," *J. Appl. Phys.* **85**, 3677 (1999).

Zhang, G.M., "Cross-over components in percolation superconductor-nonlinear-conductor mixtures," *Phys. Rev. B* **53**, 20 (1996a).

Zhang, G.M., "Higher order nonlinear response in random resistor networks: numerical studies for arbitrary nonlinearity," *Z. Phys. B* **99**, 599 (1996b).

Zhang, X., and D. Stroud, "Numerical studies of the nonlinear properties of composites," *Phys. Rev. B* **49**, 944 (1994).

Zhang, Y.M., and T.C. Wang, "Lattice Instability at a fast moving crack tip," *J. Appl. Phys.* **80**, 4332 (1996).

Zhenyi, M., S.C. Langford, J.T. Dickinson, M.H. Eengelhard, and D.R. Baer, "Scanning tunneling microscope observations of MgO fracture surfaces," *J. Mater. Res.* **6**, 183 (1990).

Zhou, S.J., D.M. Beazley, P.S. Lomdahl, and B.L. Holian, "Large-scale molecular dynamics simulations of three-dimensional ductile failure," *Phys. Rev. Lett.* **78**, 479 (1997).

Zhou, S.J., A.E. Carlsson, and R. Thomson, "Dislocation nucleation and crack stability: lattice Green's-function treatment of cracks in a model hexagonal lattice," *Phys. Rev. B* **47**, 7710 (1993).

Zhou, S.J., A.E. Carlsson, and R. Thomson, "Crack blunting effects on dislocation emission from cracks," *Phys. Rev. Lett.* **72**, 852 (1994).

Zhou, S.J., and W.A. Curtin, "Failure of fiber composites: a lattice Green function model," *Acta Metall. Matter* **43**, 3093 (1995).

Zhou, S.J., P.S. Lomdahl, R. Thomson, B.L. Holian, "Dynamic crack processes via molecular dynamics," *Phys. Rev. Lett.* **76**, 2318 (1996).

Zielinski, W., H. Huang, S. Venkataraman, and W.W. Gerberich, "Dislocation distribution under a microindentation into an iron-silicon single crystal," *Philos. Mag. A* **72**, 1221 (1995).

Zimmerman, C., W. Klemm, and K. Schonert, "Dynamic energy release rate and fracture heat in polymethylmethacrylate (PMMA) and a high strength steel," *Eng. Fract. Mech.* **20**, 777 (1984).

# Index