

David Lane, Sander van der Leeuw
Denise Pumain, Geoffrey West
Editors

Methodos Series 7

Complexity Perspectives in Innovation and Social Change



Springer

Complexity Perspectives in Innovation and Social Change

METHODOS SERIES

VOLUME 7

Editor

DANIEL COURGEAU, *Institut National d'Études Démographiques*
ROBERT FRANCK, *Université Catholique de Louvain*

Editorial Advisory Board

PETER ABELL, *London School of Economics*
PATRICK DOREIAN, *University of Pittsburgh*
SANDER GREENLAND, *UCLA School of Public Health*
RAY PAWSON, *Leeds University*
CEES VAN DER EIJK, *University of Amsterdam*
BERNARD WALLISER, *Ecole Nationale des Ponts et Chaussées, Paris*
BJÖRN WITTROCK, *Uppsala University*
GUILLAUME WUNSCH, *Université Catholique de Louvain*

This Book Series is devoted to examining and solving the major methodological problems social sciences are facing. Take for example the gap between empirical and theoretical research, the explanatory power of models, the relevance of multilevel analysis, the weakness of cumulative knowledge, the role of ordinary knowledge in the research process, or the place which should be reserved to “time, change and history” when explaining social facts. These problems are well known and yet they are seldom treated in depth in scientific literature because of their general nature.

So that these problems may be examined and solutions found, the series prompts and fosters the setting-up of international multidisciplinary research teams, and it is work by these teams that appears in the Book Series. The series can also host books produced by a single author which follow the same objectives. Proposals for manuscripts and plans for collective books will be carefully examined.

The epistemological scope of these methodological problems is obvious and resorting to Philosophy of Science becomes a necessity. The main objective of the Series remains however the methodological solutions that can be applied to the problems in hand. Therefore the books of the Series are closely connected to the research practices.

For further volumes:

<http://www.springer.com/series/6279>

David Lane · Denise Pumain ·
Sander Ernst van der Leeuw · Geoffrey West
(Editors)

Complexity Perspectives in Innovation and Social Change

 Springer

Editors

Prof. David Lane
Università Modena e Reggio
Emilia
Fac. Economia
Viale Berengario, 51
41100 Modena
Italy
lane@unimo.it

Prof. Denise Pumain
Universités Paris 1
CNRS
UMR Géographie-Cités
13 rue du Four
75006 Paris
France
pumain@parisgeo.cnrs.fr

Prof. Sander Ernst van der Leeuw
Arizona State University
School of Human Evolution &
Social Change
Tempe AZ 85287-2402
USA
vanderle@asu.edu

Prof. Geoffrey West
Santa Fe Institute
1399 Hyde Park Road
Santa Fe NM 87501
USA
gbw@santafe.edu

ISBN: 978-1-4020-9662-4

e-ISBN: 978-1-4020-9663-1

DOI 10.1007/978-1-4020-9663-1

Library of Congress Control Number: 2008942069

© Springer Science+Business Media B.V. 2009

No part of this work may be reproduced, stored in a retrieval system, or transmitted in any form or by any means, electronic, mechanical, photocopying, microfilming, recording or otherwise, without written permission from the Publisher, with the exception of any material supplied specifically for the purpose of being entered and executed on a computer system, for exclusive use by the purchaser of the work.

Printed on acid-free paper

9 8 7 6 5 4 3 2 1

springer.com

Contents

Introduction	1
David Lane, Denise Pumain and Sander van der Leeuw	
Part I From Biology to Society	
1 From Population to Organization Thinking	11
David Lane, Robert Maxfield, Dwight Read and Sander van der Leeuw	
2 The Innovation Innovation	43
Dwight Read, David Lane and Sander van der Leeuw	
3 The Long-Term Evolution of Social Organization	85
Sander van der Leeuw, David Lane and Dwight Read	
4 Biological Metaphors in Economics: Natural Selection and Competition	117
Andrea Ginzburg	
5 Innovation in the Context of Networks, Hierarchies, and Cohesion ...	153
Douglas R. White	
Part II Innovation and Urban Systems	
6 The Organization of Urban Systems	197
Anne Bretagnolle, Denise Pumain and Céline Vacchiani-Marcuzzo	
7 The Self Similarity of Human Social Organization and Dynamics in Cities	221
Luís M.A. Bettencourt, José Lobo and Geoffrey B. West	
8 Innovation Cycles and Urban Dynamics	237
Denise Pumain, Fabien Paulus and Céline Vacchiani-Marcuzzo	

Part III Innovation and Market Systems

- 9 Building a New Market System: Effective Action, Redirection and Generative Relationships** 263
David Lane and Robert Maxfield
- 10 Incorporating a New Technology into Agent-Artifact Space: The Case of Control System Automation in Europe** 289
Federica Rossi, Paolo Bertossi, Paolo Gurisatti and Luisa Sovieni
- 11 Innovation Policy: Levels and Levers** 311
Federica Rossi and Margherita Russo

Part IV Modeling Innovation and Social Change

- 12 The Future of Urban Systems: Exploratory Models** 331
Denise Pumain, Lena Sanders, Anne Bretagnolle, Benoît Glisse and H  l  ne Mathian
- 13 Modeling Innovation** 361
Roberto Serra, Marco Villani and David Lane
- 14 An Agent-Based Model of Information Flows in Social Dynamics** 389
Davide Ferrari, Dwight Read and Sander van der Leeuw
- 15 Exaptive Processes: An Agent Based Model** 413
Marco Villani, Stefano Bonacini, Davide Ferrari and Roberto Serra
- 16 Power Laws in Urban Supply Networks, Social Systems, and Dense Pedestrian Crowds** 433
Dirk Helbing, Christian K  hnert, Stefan L  mmer, Anders Johansson, Bj  rn Gehlsen, Hendrik Ammoser and Geoffrey B. West
- 17 Using Statistical Physics to Understand Relational Space: A Case Study from Mediterranean Prehistory** 451
Tim Evans, Carl Knappett and Ray Rivers
- Conclusion** 481
David Lane, Denise Pumain and Sander van der Leeuw
- Author Index** 489
- Subject Index** 491

Contributors

Hendrik Ammoser Institute for Transport and Economics, Dresden University of Technology, Dresden, Germany

Paolo Bertossi Università di Modena e Reggio Emilia, Reggio Emilia, Italy

Luís M. A. Bettencourt Los Alamos National Laboratory, Santa Fe Institute and Mathematical, Computational & Modeling Sciences Center, Arizona State University, Tempe, AZ, USA

Stefano Bonacini Department of Social, Cognitive and Quantitative Sciences, University of Modena and Reggio Emilia, Reggio Emilia, Italy

Anne Bretagnolle Université Paris I Panthéon Sorbonne, UFR de Géographie, Paris, France

Tim Evans Department of Physics, Imperial College London, London, UK

Davide Ferrari School of Statistics, University of Minnesota, Minneapolis and Saint Paul, Minnesota, USA

Björn Gehlsen Institute for Transport and Economics, Dresden University of Technology, Dresden, Germany

Andrea Ginzburg Università di Modena e Reggio Emilia, Reggio Emilia, Italy, ginzburg@unimore.it

Benoît Glisse Laboratoire LIP6, Université Paris VI, Paris, France

Paolo Gurisatti Università di Trento, Trento, TN, Italy

Dirk Helbing Institute for Transport and Economics, TU Dresden, Dresden, Germany; Collegium Budapest – Institute for Advanced Study, Szentháromság u. 2, 1014 Budapest, Hungary

Anders Johansson Institute for Transport and Economics, Dresden University of Technology, Dresden, Germany

Carl Knappett Department of Art, University of Toronto, Toronto, Ontario, Canada

Christian Kühnert Institute for Transport and Economics, Dresden University of Technology, Dresden, Germany

Stefan Lämmer Institute for Transport and Economics, Dresden University of Technology, Dresden, Germany

David Lane Department of Social, Cognitive and Quantitative Sciences, University of Modena and Reggio Emilia, Via Giglioli Valle, 9 – 42100, Reggio Emilia, Italy

José Lobo School of Human Evolution and Social Change and W.P. Carey School of Business, Arizona State University, Tempe, AZ, USA

Hélène Mathian CNRS, UMR Géographie-cités, Paris, France

Robert Maxfield Stanford University, Stanford, CA, USA

Fabien Paulus Université de Strasbourg, UFR de Géographie, Paris, France

Denise Pumain Université Paris I Panthéon Sorbonne, UFR de Géographie, Paris, France

Dwight Read Department of Anthropology, University of California in Los Angeles, CA, USA

Ray Rivers Department of Physics, Imperial College London, London, UK

Federica Rossi Dipartimento di Economia, University of Torino, Torino, Italy

Margherita Russo Dipartimento di Economia Politica, University of Modena and Reggio Emilia, Modena, Italy

Lena Sanders CNRS, UMR Géographie-cités, Paris, France

Roberto Serra Department of Social, Cognitive and Quantitative Sciences, University of Modena and Reggio Emilia, Reggio Emilia, Italy

Luisa Sovieni Università di Modena e Reggio Emilia, Reggio Emilia, Italy

Céline Vacchiani-Marcuzzo Université de Reims, UFR Lettres et Sciences Humaines, Reims, France

Sander Ernst van der Leeuw School of Human Evolution and Social Change, Arizona State University and Santa fe Institute, AZ, USA

Marco Villani Department of Social, Cognitive and Quantitative Sciences, University of Modena and Reggio Emilia, Reggio Emilia, Italy, mvillani@unimore.it

Geoffrey B. West Los Alamos National Laboratory, Santa Fe Institute and Mathematical, Computational & Modeling Sciences Center, Arizona State University, Tempe, AZ, USA

Douglas R. White School of Social Sciences, University of California, Irvine, CA, USA

Introduction

David Lane, Denise Pumain and Sander van der Leeuw

The project that resulted in this book originates in an encounter of three of the authors (David Lane, Sander van der Leeuw and Geoffrey West) in the summer of 2001 at the Santa Fe Institute around two main themes: (1) a different way of looking at the invention and innovation of artefacts, and (2) the possible impact of innovation on urban dynamics. One of us was a physicist (West), one a statistician (Lane) and one an archaeologist (van der Leeuw). Almost immediately, we asked Denise Pumain (an urban geographer) to join us in this adventure.

Fortunately for us, our meeting coincided with an initiative of the head of the newly formed unit for a project officer of the Future and Emerging Technologies Program of the European Union's Directorate for Information Science and Technology, Dr. Ralph Dum, to stimulate a wide spectrum of research into the potential uses of Complex Systems approaches. Hence, we brought four teams together, at three European universities: the University of Modena and Reggio Emilia (Lane), the University of Paris I (Panthéon-Sorbonne) (Pumain and van der Leeuw), and Imperial College, London (West).

Ralph encouraged us to apply, and after the usual vetting procedure our proposal for a project that considered the Information Society as a Complex System (ISCOM) was accepted and funded. We began work in July 2002, and the project lasted for four years, until the end of June 2006. Those years turned out to be a very exciting intellectual adventure for all of us, as well as the members of our teams and the colleagues whom we invited to our workshops in London, Santa Fe, Venice, Modena, and Paris. Some of the results lie in front of you. But we do not think that we are exaggerating if we say that the collaboration influenced the thinking of all of us to such a point that other results will follow in due time, whether under our name or under that of the many other members of the team (see the list that follows this introduction).

Two conclusions stand out from the project. Firstly that innovation and invention have, in a sense, been among the stepchildren of modern research, whether in the

D. Lane (✉)

Department of Social, Cognitive and Quantitative Sciences, University of Modena and Reggio Emilia, Reggio Emilia, Italy

social sciences or in the humanities, and secondly that the role of innovation in urban dynamics is much more important than is generally acknowledged.

We live in a world that is driven by invention and innovation, but that has not always been so. In the XVIIth century, innovation was a 'dirty word'. The world order was deemed to be immutable; people behaved as their ancestors had done (or at least they believed they did, and they often strived hard to meet that ideal) (Girard 1990).

Little by little, though, over the last three centuries, 'history' and 'tradition' ceded place to 'nature' as the concept invoked to explain the world order. We still speak in many instances of 'it is natural' when we wish to express the fact that we think that a kind of behavior is in harmony with the world order. In the process, in the first part of the XIXth century, History has become a discipline, rather than the omnipresent way to explain what happened or what happens.

Simultaneously, we observe a growing *emphasis on the new rather than the old* – particularly during, and as a result of, the Enlightenment and the Industrial Revolution (Girard, 1990). As science and technology gained in importance, the conceptual and instrumental toolkit of the (western) world grew exponentially, and in doing so enabled humanity to identify and tackle more and more challenges. As a result, the number of inventions and innovations around us is increasing dramatically. This is clearly visible if one looks at the number of inventions that are patented in the industrial countries.

After the industrial and nuclear revolutions, we are now witnessing the silicon, information technology and communications revolutions, and the nanotechnology, biotechnology, and cognitive revolutions are on the horizon, each of which is opening another whole new domain of knowledge, know-how and innovation. For the moment at least, there does not seem an end in sight to this acceleration of change in our world.

In that context it is in our opinion surprising that the scientific community has generated so little understanding of the process of invention and innovation itself. Generally, the world reacts *a posteriori* to innovations once they have been introduced. Could we not attempt to shift our stance from a re-active to a pro-active one, and come to understand and guide the process of invention and innovation itself? That would put us in control rather than dealing with things after they have gotten out of hand, and it would potentially allow us to accelerate the innovative process in those domains in which that is most needed, and maybe slow it in others

What has thus far held back our understanding of the process of invention and innovation? Our tentative working hypothesis is that this lack of understanding is directly related to the fact that the majority of the scientific community has looked at invention and innovation using a positivist, scientific perspective. In essence, invention and innovation have mainly been studied 'a posteriori'. From such a perspective, creation cannot be described or understood. Hence, we have left 'invention' completely to one side in innovation studies, relegating it to the domain of 'personal creativity', and we have focused uniquely on innovation, i.e. on the ways in which an invention is adopted and spreads throughout a population.

The first towns in the world were founded about six or seven thousand years ago. After a slow start, urbanization is now everywhere around us. Currently, about 50% of the world population lives in towns, and this number is growing so rapidly, that the 80% threshold may be reached within twenty to forty years. In effect, urbanization seems to be an unstoppable ‘explosion’ that is only equalled by the explosion in inventions we have seen over the last century or so. Hence, the idea that there might be a relationship between the ‘innovation explosion’ we have just referred to and the ‘urban explosion’ seems worth looking into.

Our team has convincingly demonstrated that innovation (as represented by the number of people involved in research, the number of research organizations, the number of patents submitted, etc.) scales super-linearly with the size of urban agglomerations, while energy scales sub-linearly and services linearly (cf. Bettencourt, Lobo, & Strumsky, 2007; Strumsky, Lobo, & Fleming, 2005, Bettencourt et al., 2006; Pumain et al., 2006 and Chapters 7 and 8 in this book). This seems to point to the fact that, whereas economies of scale in energy use are an important phenomenon in urbanization, managing information—generating new things and patterns of Organization—is the actual driver behind urbanization.

Because people congregate in cities, the latter harness the densest and the most diverse information processing capacity. Not only does this relatively high information processing capacity ensure that they are able to maintain control over the channels through which goods and people flow on a daily basis, but their cultural (and, thus, information-processing) diversity also makes them into preferred loci of invention and innovation.

The super-linear scaling of innovation with city size enables cities to ensure the long-term maintenance of the information gradient that structures the whole system. It is due to a positive feedback loop between two of any city’s roles. On the one hand, most flows of goods and people go through towns and cities. That confronts them most intensely with information about what is happening elsewhere, and this – again – enhances their potential for invention and innovation. But the same connections enable them to export these innovations most effectively – exchanging some part of their information processing superiority for material wealth. Cities are demographic centers, administrative centers, foci of road systems, but above all they are the nodes in the system where the most information processing goes on. As such, they are the backbone of any large-scale social system. They operate in network-based “urban systems” which link all of them within a particular sphere of influence. Such systems have structural properties that derive from the relative position all the cities occupy on the information-processing gradient, and in the communications and exchange networks that link them to each other (White, Chapter 5 in this volume). Although the role of individual towns in such systems may change (relatively) rapidly (Guerin-Pace, 1993), the overall dynamic structures are rather stable over long periods of time. In the long run, the organisation of human habitat in networks of cities reduces the uncertainties of a closed environment by relying on more distant resources as well as by creating new ones. There is a shift from a human ecology toward another way of structuring the planet, entangling

territories (geographical structures based on continuity) in societal networks (that are based on connectivity).

The books of the Methodos series explore some of the new ways that emerge from the complex systems paradigm. Without entering into the debate about the possible emergence of a unified science of complex systems, we want here to develop a new theory of human social change, within a perspective that is informed by the recent developments of the complex systems paradigm. In this book, we will explain why we think that this paradigm can help us to identify the specificity of innovation and change in social systems.

Part I of this book: *From biology to society*, specifies how a new kind of organisation has emerged with the historical apparition of human societies. Although *Homo sapiens* is a biological species, whose individual elements do not in themselves differ from any other animal species in their biological organisation, and although social systems do share some properties with animal social organisations, two main radically new and distinctive features were created through the process that led to human social organisation. The first one is a self-monitored, directed (intentional) modality of social change. We shall demonstrate that this new kind of evolutionary driver is the result of the integration of new functionalities in social structures due to cultural processes. The second distinctive feature that is essential to our approach of social systems is that it is comprehensive: to shift from a static description of social structures to a dynamic one, we need to consider a variety of social interactions that are usually separated in disciplinary explanations of social systems. The modifications in social organisation that are directed at monitoring social changes, and that produce emergent patterns instantiated in organisations do affect a social system in every aspect and at all its levels of organisation. We describe how function, structure and process are affecting each other, and we build a dynamic, interactionist interpretation of the evolution of social systems.

In this attempt, it is important to determine which ingredients are necessary for developing a theory of human social innovation that is both general, and precise enough to be relevant. We believe that complexity theories are the necessary framework for developing a modern interpretation of change in complex systems. However, we question two principles that are part of the application of this theoretical approach to physical and biological systems. These are, firstly, the search for invariance and universality in processes. We demonstrate that human social change cannot be described in Darwinian terms, because something new has appeared, *i.e.* the fact that human societies are inherently responsible for their own innovation. This then leads us to question the applicability of the Darwinian approach of biological evolution to human social evolution, which we discuss in the first five chapters of this book.

Part II, *Innovation and urban systems*, and Part III, *Innovation and market systems*, develop examples of the application of these ideas and work out more precisely a number of aspects of this perspective. Urban systems are at the core of many important issues in contemporary societies. While cities concentrate a majority of the world population and human activities, the urban way of life may encounter limits that are increasingly perceptible in terms of the potential shortage

of environmental resources (mainly energy, soil and water), in terms of organising a livable social mix, and in terms of managing local systems that are threatened by the unpredictable effects of their increasing connectedness to a multiplicity of networks through globalisation.

Sustainable development, social cohesion and territorial cohesion have become challenging issues for the monitoring of urban systems, in modern information societies as well as in developing and poor countries. It is thus essential to develop a proper understanding of urban dynamics in its complex articulation of a typical hierarchical structure with monitored but highly decentralised innovation processes that periodically renew the functionalities of individual towns and cities in increasingly better connected and wider-spread systems of cities. This enlarged (comprehensive) view of urban spatio-temporal evolution and its connection to social innovation is developed in Chapter 6. In biology, scaling laws have been identified as useful methodological tools reflecting the effect of energetic and geometric constraints on the development and structural organisation of living systems. Two chapters are dedicated to the application of this methodology on urban systems, using different approaches to determine to what extent further urban growth may depend on similar effects or not.

Markets are perhaps less easy to coin as complex social systems whose evolution requires a specific social input in a general theory of complex systems. Indeed, they are very often analysed within a unique disciplinary framework, and sophisticated mathematical models derived from the principles of economic theory are sometimes thought of as very successful descriptions of even the most capricious fluctuations in stock exchange evolutions that would open the way to quasi physical theories of these very decentralised and complex systems. However, we claim that such interpretations can only operate in specific contexts, where the social ingredients of what constitute the markets are clearly defined. If we want to understand real market evolution, we have to develop a broader theory of markets as social organisation. A major point we insist on is to demonstrate how social interaction is linked to the invention of new artefacts and their functionalities in social systems. This is done in three chapters of section three, including a reflection on the implications of our perspective for social policies favouring innovation.

Our complexity perspective has several methodological implications that we develop in Part IV of this book, *Modeling innovation and social change*. Models are the indispensable steps that enable a fruitful dialog between theory and empirical studies. In dealing with social change, we need abstract narratives that can identify what is really changing and what remains invariant over time in the evolution of social systems. Within the continuous and more or less rapid flow of social innovation that manifests itself in the production of new artefacts, new institutions, technological inventions, scientific breakthroughs, but also new social practices, collective representations and beliefs as well as new modes of social interaction, what are the features that are significant for interpreting social change? What are the decisive configurations in structure, function and process that are necessary for driving that evolution? Among those, what are the possible levers for intelligent

monitoring, or partial control over the anticipated evolution? Models are useful tools for answering these questions.

The abstract description of social evolution can be translated into models. Mathematical and computational models help us to understand how social systems can share some features of their structure and evolution with other complex systems, as reflected for instance in structural power laws or scaling laws, whereas the parameters that are involved in these models take specific values which may imply a quite different qualitative evolutionary behaviour from natural or living systems. Because of that, we chose to develop not only analytical models of social change, but more flexible models that are no longer analytical but computational, including multi-agent models. These models allow the handling of both invariant features, including entities and rules whose properties represent stylised facts from observed empirical evolution, and creative aspects of social organisation when the nature of agents or artefacts is transformed through dynamic interactive processes. Although this can be partially modelled, the exogenous intervention of the modeller is still necessary to match such a creative evolution.

Another way in which computational models are helpful is that they enable the exploration of unexpected events and unforeseeable futures. Social changes do not operate in a predefined space where everything that might happen has been predefined. Mathematical models are not able to handle systems where something new can change the meaning of variables and parameters or create new categories. Whether the methods of artificial intelligence will be able to do so is still a matter of debate. We use AI here in modelling as a tool for exploring the limits between what can be endogenously produced from interactions within the model, and what has to be imported as exogenous knowledge by the model designer. Models are tools of random search in a space of not-yet-definite potential futures. Whenever possible, we have privileged data-driven simulation as a way of constructing models. This is the case for a variety of applications presented in section four (especially Chapters 13 and 17 on the dynamics of urban systems). The theory is injected in the model as abstract knowledge defining endogenous processes, and the results of simulation allow us to identify which specific processes have to be inserted from the outside to match observed or projected structures and functionalities.

As editors of this volume, we would not want to conclude this introduction without expressing our thanks to many people. First of all to the authors of the papers, who were our partners in the discussions and debates that led to these pages. Secondly to all those who helped us, at various times, to organize the many workshops and dealt with the inevitably complex administrative procedures, among whom we would like to mention particularly Irene Poli, Federica Rossi, Antoine Weexsteen, Martine Laborde, and the staff of the Santa Fe Institute. And thirdly to Callie Babbitt, who did all the final editing and typesetting of the manuscripts.

References

- Bettencourt, L.M.A., Lobo, J., Helbing, D., Kühnert, C., & West, G.B. (2006). Growth, innovation, scaling and the pace of life in cities. *Proceedings of the National Academy of Sciences USA*, 104(17) 7301–7306.
- Bettencourt, L.M.A., Lobo, J., & Strumsky, D. (2007). Invention in the city: Increasing returns to patenting as a scaling function of metropolitan size. *Research Policy* 36, 107–120.
- Girard, R. (1990). Innovation and repetition. *SubStance* 62/63, 7–20.
- Guérin-Pace, F. (1993). *Deux siècles de croissance urbaine*. Paris: Anthropos.
- Pumain, D. (ed.) (2006). *Hierarchy in natural and social sciences* (243 p). Dordrecht: Springer.
- Strumsky, D., Lobo, J., & Fleming, L., (2005). *Metropolitan patenting, inventor agglomeration, and social networks: A tale of two effects*. Santa Fe Institute Working Paper 05-02-004 (<http://www.santafe.edu/research/publications/wpabstract/200502004>).

Part I
From Biology to Society

Chapter 1

From Population to Organization Thinking

David Lane, Robert Maxfield, Dwight Read and Sander van der Leeuw

1.1 Introduction

Our species is still very young by biological time scales, and it is too early to know if we represent the cutting edge of a biological success story, like cockroaches or dinosaurs, or a brilliant but ultimately failed and short-lived experiment in niche construction and destruction. In the mere 200,000 or so years of *Homo sapiens*' story, and in particular in the approximately 50,000 years since we began to accrue the accoutrements of culture like language, art and multi-component artifacts, members of our species have populated a vast extent of the earth's surface and exploited for our own purposes an ever-increasing share of the planet's biologically utilizable solar energy. In the last few centuries, we have ravaged the stock of bioprocessed solar energy accumulated over millions of years, transformed minerals extracted from below the earth's surface into a huge variety of forms and new materials that satisfy what we regard as our needs, and increasingly concentrated the human population in urban spaces to which nearly all the raw materials necessary for human survival have to be imported from elsewhere. At the individual level, our first *Homo sapiens* ancestors managed to keep themselves going on the 100–300 watts their bodies were able to generate, assisted in their quest for survival by the handful of artifacts they knew how to make and use; in contrast, current residents of New York City mobilize on average about 10,000 watts to propel them through their daily rounds of activity, and the shops in their city offer them a choice of something like 10^{10} different kinds of artifacts to help them accomplish whatever it is they might feel inclined to do.¹

How have we managed to accomplish so much so fast? The main premise of this chapter is that we have done it through a new modality of innovation, through which human beings generate new *artifacts* that they embed in new *collective activities*,

D. Lane (✉)

Department of Social, Cognitive and Quantitative Sciences,
University of Modena and Reggio Emilia, Via Giglioli Valle, Reggio Emilia, Italy
e-mail: lane@unimo.it

¹ See Chapter 12. The number of artifacts refers to the number of different SKU's – that is, product labels used for distinct bar codes – on offer in New York, as estimated and reported in Beinhocker (2006). This is of course a rather crude measure of “kinds of artifacts”.

which are in turn supported by new *organizations* and sustained by new *values*. Over time, this new innovation modality gave rise to a positive feedback dynamic, which explains how we have generated so many transformations in our selves, our societies, our culture and our environment.

What is this new innovation modality? We begin by describing something it is *not*. Much of modern biology is based upon Darwin's theory of biological novelty, which analyzes the processes through which species come into being and are transformed, by means of mechanisms of heritable variation and selection. Given the tremendous scientific success of this theory, it is not surprising that many authors have sought to adapt it to other contexts. In particular, it is becoming increasingly fashionable to construct theories of innovation in human society and culture on a Darwinian foundation. We shall argue that this move is mistaken.

To be clear about the claim we are making, we need to define precisely what counts for us as a Darwinian foundation. An evolutionary theory seeks to understand a phenomenon by describing the processes that brought that phenomenon into being and that generate the transformations it successively undergoes.² A *Darwinian account* is a special kind evolutionary theory that, like the theory Darwin set out in the *Origin of Species*, rests on two foundational bases:

- it is characterized by what Ernst Mayr (1959, 1991) called *population thinking*; and
- it analyzes the transformation processes with which it is concerned in terms of two *fundamental* and *distinct* classes of mechanisms: variation and selection.

We discuss these two concepts, in the context in which Darwin introduced them, in the second section of the chapter, and we explain why they were such a great success in this context. We then describe some issues that must be resolved before they can be successfully applied in other contexts.

In the third section, we examine some features of human sociocultural innovation that distinguish it from biological evolution. We argue that these features are not consistent with the foundations underlying a Darwinian account. In particular, the distinction between variation and selection processes is difficult to maintain – and even when they can be distinguished, they frequently fail to be fundamental, since another kind of process, negotiation, underlies both of them. In addition, it is often impossible to identify a relevant “population”, which is the critical starting point for population thinking.

In the fourth section, we describe a shift in perspective, from population thinking to *organization thinking*. After analyzing some critical differences between

² In contrast, essentialist theories seek to explain phenomena by isolating their (unchangeable) essences and classifying all other aspects of the phenomena either as deducible consequences of these essences or inessential epiphenomena. For example, a noted neoclassical economist once told one of us, that “if it isn't about optimization and rationality, it isn't economics”. Essentialist theories tend to be associated with the “typological thinking” mentioned in the following section, just like Darwinian accounts are associated with “population thinking”.

biological organization and sociocultural organization, we propose an innovation modality alternative to the Darwinian account. We call this modality *organizational self-transformation*, and we argue that it bears a relation to sociocultural evolution similar to the Darwinian account's relation to biological evolution. We conclude by describing some features of organizational self-transformation, in particular the positive feedback dynamic to which it gives rise. This dynamic, which we call *exaptive bootstrapping*, has generated the proliferation of artifacts and organizations that construct our current sociocultural world – and us.

1.2 Darwin's Theory of Biological Evolution: What It Is, Why It Works

For over a hundred years before Darwin published the *Origin of Species*, students of natural history had collected a huge body of evidence illustrating what they regarded as a nearly perfect match between the morphological characteristics and behaviors of biological individuals and the environmental opportunities that these individuals exploited to earn their living. For example, the woodpecker's beak seemed *designed* to drill into bark to extract the insects the woodpeckers ate, while the hummingbird's beak and capacity to hover could hardly be improved upon as a way to sip nectar from flowers. Such matches between structure and functionality were generally interpreted as evidence for the existence of a benevolent Designer, who had constructed a world able to sustain in harmonious equilibrium all the products of His creation, including the many species of plants and animals that the natural historians were busy describing and classifying.

For these natural historians, a species was an immutable, ideal organization, exquisitely tuned to the exigencies of the environment in which members of the species sustained and reproduced themselves. This organization was more or less perfectly instantiated in all the individuals that belong to the species. Ernst Mayr (1959, 1991) labeled this way of construing biological organization *typological thinking*.

Darwin's first great accomplishment was to displace typological thinking by what Mayr called *population thinking*. Population thinking *un-reifies* the species. Instead of a timeless *kind*, the species becomes just an *aggregation of individuals*. The species has a beginning, in an act of speciation, and sooner or later it will have an end, when it becomes extinct. Its boundaries – who is in, and who is not – are determined by pedigree: once the species has come into being, no individual may enter unless its parent or parents are already members, and no member can defect, except through death. For Darwin, what counts for a species as *population* is not the commonalities that its members share, but the *variation* among conspecifics. That is, the species as population is characterized not by a shared organization, but by the *statistical distribution of those features that differ* among members of the population. Thus, not only is the *species* no longer identified with an *ideal* organization but the very *real* organization of an *individual* is disregarded in Darwin's story, and

the individual figures in it merely as a set of distinct features. From the point of view of population thinking, biological evolution is the story of the changes in the distribution in the population of these features over time, including the introduction of new features.³

The species as population is not composed of a fixed set of individuals: the members of the population change over time, through birth and death events. It is critically important to the success of Darwin's story that the features of new individuals are statistically related to the current distribution of features in the population. Fortunately, for features that are *heritable* this turns out to be the case. That is, absent changes due to the mechanisms of variation and selection described in the next paragraph, reproduction – whether it be sexual or asexual – does not change⁴ the distribution of differing heritable features, which, as we have seen, Darwin's theory takes as the variable of interest in the species transformation story.

Darwin analyzed the change in feature distribution over time by distinguishing between two fundamental classes of mechanisms through which change happens: mechanisms of *variation* and mechanisms of *selection*. The former introduce *new* variant features that individuals in the population may possess, while the latter determine which features will *increase in frequency* in the population over time. New heritable features persist in the population, unless they are eliminated by selection mechanisms.⁵ In *Origin*, he called attention to one particular selection mechanism, which he regarded as the primary determinant of the direction of change in feature distribution: *natural selection*. According to Darwin, individuals in the same species were always in competition among themselves to obtain the resources necessary to survive and reproduce. Thus, a heritable feature that helps individuals possessing it to get a reproductive edge over their conspecifics would tend to increase in frequency generation by generation. Such features came to be described as increasing the *fitness* of individuals who possessed them, and natural selection came to be conceived in terms of a stochastic process guided by a *fitness function*, whose value for a set of features represented the relative competitive reproductive advantage of an individual possessing this feature set.

According to the currently canonical version of Darwinian theory, the so-called neo-Darwinian synthesis that welded together ideas from Mendelian genetics to

³ Of course, this begs some key questions, which Darwin did not address: just what constitutes a new *feature*, and how can it be integrated functionally with previously existing ones; or which new features – or distributional changes in existing features – constitute the origin of a new species? Biologists since Darwin have addressed these questions, sometimes supplementing but never successfully challenging the basic suppositions underlying population thinking.

⁴ In expectation, with random mating in sexually reproducing species (that is, absent sexual selection, a selection mechanism). Of course, in sexually reproducing species, offspring will not have the same features as both their parents, and so the *realized* distribution after a new individual is born is not identical to what it was before the birth – it is conventional to refer to “random drift” as a selection mechanism when such random oscillations contribute to the directionality of change in the feature vector distribution.

⁵ Here, as in the previous footnote, we follow the convention that classifies random drift as a selection mechanism.

Darwin's evolutionary theory in the first three decades of the 20th century, the key variation mechanisms for heritable features in biological evolution derive from genetic operations during reproduction. Thus, a particular innovation is initiated as a new variant *genotype*, which, in interaction with pre-existing structural or behavioral features and particular environmental contexts, happens to result in a new structure or a new behavior at the *phenotypic* level. This new structure or behavior may get incorporated in some kind of process with already existing behaviors, structures and environmental features, and this new process may provide a survival or reproductive advantage to the individual who possesses it. If so, by natural selection, the frequency of individuals with the innovation will increase over time in the population, and the innovation will count as a success.⁶

From the moment of *Origin's* publication, Darwin's ideas were discussed everywhere, not just in the restricted circles of the natural historians.⁷ This initial high level of interest ensured that many other scientists would continue to probe these ideas – to support them, oppose them, extend them. It took several generations for the initial *succès de scandale* to develop into scientific orthodoxy, but by the 1930s, this too had been achieved. How can we account for the great success of Darwin's theory of biological evolution? We think four reasons, described below, were primarily responsible. The first two of these account for the unusual degree of interest that Darwin's ideas encountered in the first decades after *Origin* and kept them under the spotlight of intense scientific exploration, and the second two were responsible for the construction of the scientific consensus around the Darwinian account as a foundation for a theory of biological innovation.

1.2.1 R1: A Non-Teleological Explanation for Structure-Function Fit

Darwin's theory provided the first really plausible alternative to the hypothesis of intelligent design as an explanation for the extraordinary match in the biological world between structure and function. It differed in two fundamental ways from intelligent design. According to intelligent design, biological organization is *immutable* and *teleological*: each structural feature was designed *in order to* carry out a particular function, and once designed, it need not – and did not – change. In contrast, in Darwin's story, structural change was the *key* to the match between structure and function. New variants might appear in any possible direction in feature space, but

⁶ Note that *how* the genotype contributes to the construction of the phenotype is not part of Darwin's account. This important problem, under the rubric of "evo-devo", is currently a very hot topic in biological research. Evo-devo extends rather than revises the Darwinian account, and the separation between variation and selection mechanisms serves as a strong constraint on the kinds of structures that evolution can produce – for example, new variations, which occur at the genotypic level, cannot be directed toward the provision of specific functions at the phenotypic level.

⁷ See, for example, Desmond and Moore (1991), Chapter 33.

only those that provided an advantage with respect to the overarching functional imperatives – to survive and to reproduce – were retained through the generations.

The reaction to a non-teleological alternative to intelligent design was explosive. It empowered scientific radicals like Thomas Huxley and atheistic lay(wo)men like Harriet Martineau to take the initiative in dissolving the claims for divine causality that had subordinated natural history – and science more generally – to religion. Correspondingly, it deeply upset, even scandalized, proponents of traditional values, from scientists like Richard Owen and Louis Agassiz to clergymen like Samuel Wilberforce, who in an 1860 Oxford free-for-all discussion on Darwin's ideas asked Huxley from which side of his family was he descended from apes. Though debates of this sort generated substantially more heat than light, the idea of directional change without a directing intelligence became a central theme in the ongoing cultural battle between an emerging materialistic scientism and traditional deference to established authority and revealed religion. For scientists on both sides of this battle, finding evidence and arguments that bore on the plausibility of Darwin's ideas, whether to support or demolish theory, had a very high priority – and the guarantee of a large, very interested public should they accomplish a breakthrough in this quest.

1.2.2 R2: Demonstrating the Plausibility of Organizational Innovation from Gradual Feature Change: Two Analogies

However internally consistent Darwin's ideas were, they were certainly not directly demonstrable from empirical evidence. Indeed, to make his argument plausible, Darwin had to come to grips with a very difficult empirical problem: many natural historians were prepared to argue that the immutability of species was no hypothesis, but observable fact. From Aristotle to Linnaeus to Darwin, no observer of the natural world had seen a new species come into being. True, paleontologists had found what seemed to be fossilized remains of plants and animals that did not seem to belong to any known existing species, and enough progress had been made in the relative dating of geologic strata to convince most scientists that not all presumed life forms – living or fossil – had come into being in the same epoch. But the possibility of extinction and separate epochs of creation did not contradict the assumption that, once created, a species-as-ideal-organization didn't change. Where was the evidence for change?

Darwin answered this question with a spectacular rhetorical move, accomplished through the juxtaposition of two analogies. The first analogy was intended to show that biological form is indeed malleable.⁸ Darwin reminded his readers of the wonders achieved by plant and animal breeders – in particular, dog and pigeon

⁸ This analogy, as Darwin himself points out in the introduction to the sixth edition of *Origin*, had been introduced for this purpose by other authors before him. No other author, as far as we know, had coupled it as did Darwin with the analogy with Lyellian geology.

fanciers – in creating new breeds or varieties.⁹ Of course, the breeder plays an essential role here, since he selects individuals to reproduce on the basis of the features he wants to favor. In the wild, Darwin asserted, competition to survive and reproduce takes the place of the breeder. Obviously, this *natural* selection is much less intense than the breeders' *artificial* selection, and consequently will require much more time to accrue observable differences in bodily structure or behavior. Here, Darwin introduces his second analogy: we know, he claims, from the revolutionary work of Charles Lyell in geology what can happen when presently observable microprocesses act over huge spans of time. As he puts it in *Origin*, "a man should examine for himself the great piles of superimposed strata, and watch the rivulets bringing down mud, and the waves wearing away the sea-cliffs, in order to comprehend something about the duration of past time, the monuments of which we see all around us." Given enough time, small quantitative changes can produce large qualitative change. So, the same process of differential reproduction, favoring particular features, which we can observe when breeders produce new varieties, can operate over Lyellian "geological time" to produce new species by natural selection.

With this argument, Darwin changed the status of the immutability of species from an empirical fact to a mere hypothesis, and an increasingly implausible one at that – as the torturous evasions of Owen around this issue in the years after 1859, ridiculed in later editions of *Origin*, testify.¹⁰

1.2.3 R3: Carving Nature at the Joints: The Variation–Selection Dichotomy

When Darwin decomposed the evolutionary process into variation and selection mechanisms, he knew almost nothing about how the former actually worked. He was convinced, though, that they satisfied two properties, which were really all his theory at its most abstract level required of them: they generated variant features independently of their potential functionality, and these features were heritable. These properties sufficed to justify his analysis of the *directionality* of evolution – and hence, in his theory, the origin of species – solely in terms of selection mechanisms, in particular natural selection. The research initiated by the rediscovery of Mendel's work around the beginning of the 20th century resulted in the identification and detailed explication of the genetics underlying evolutionary variation mechanisms such as mutation and cross-over.

As a result, it became clear that Darwin's decomposition was not merely an analytic or conceptual move, but really carved nature at the joints. As Darwin had foreseen, they were functionally orthogonal: the directionality of evolution was supplied just by selection, not variation, mechanisms. But much more was true. The principal variation and selection mechanisms incorporated into the neo-Darwinian synthesis differed from each other with respect to their *ontological level* and their

⁹ Darwin admits that, like all analogies, this one is incomplete, since varieties aren't species: when interbred, their offspring are not sterile.

¹⁰ See Desmond and Moore, 1991, Chapters 33 and 34.

characteristic *spatiotemporal scales*. The variation mechanisms involved random changes in the genome, which took place essentially instantaneously, while selection operated in individuals interacting with their biological and physical environments and occurred on a time scale of many generations. Moreover, variants that produced viable individuals with new phenotypic features upon which selection could operate were exceedingly rare, so rare that the time between such events was sufficiently long that selection could in general process them one at a time, attaining equilibrium frequencies with respect to a new innovation without interference from another.

The significance of these level and scale differences between variation and selection mechanisms was two-fold. First, they made it possible to *parallelize* evolutionary research: laboratory scientists explored the genetic basis of variation mechanisms, while field researchers investigated the past history and present operation of directional change through selection. Secondly, the absence of strong interaction effects between the two classes of mechanisms very much simplified the work of theoreticians who sought to put together their effects to deepen and extend evolutionary theory – in particular, towards a quantitative theory.

1.2.4 R4: Mathematization

Since Galileo, the epistemological gold standard of science has been the construction of mathematical theories that provide succinct description and quantitative prediction for empirical phenomena. By the end of the 19th century, physics was enshrined as the king of sciences – in large part, because it had married the queen, mathematics. Other emerging sciences, like chemistry and economics, did their best to emulate the king in this respect. Biology seemed hopelessly behind. Despite the efforts of a few fringe players like D’Arcy Thompson, it still resembled its ancestor, natural history, much more than its successful rival science, physics. Though the theory presented in *Origin* was far more general and precise than any other yet introduced in biology, it was anything but mathematical. With the incorporation of Mendelian genetics into evolutionary theory, the situation changed. Ronald Fisher, Sewall Wright and JBS Haldane pioneered the quantitative theory of population genetics, which provided the beginning of what has become a flourishing mathematical foundation for Darwinian evolutionary biology, with a consequent upgrading of its scientific status. In this theory, natural selection is represented by a fitness function, with a natural interpretation in terms of expected offspring from individuals with alternative genomic configurations; the stochastic components of the models derive from genetic theories that can be calibrated with frequencies from genetic experiments.

1.2.5 Darwinian Accounts

As with other successful scientific theories, Darwin’s ideas have inspired scientists working on different problems than his. For example, biologists and cognitive

scientists have developed interesting and fruitful theories based on Darwinian reasoning to explain phenomena ranging from the construction of immunological and neural organizations during individual ontogeny to cognitive processes like perceptual categorization and even induction.¹¹ Such theories begin with the identification of the essential elements of what we will call a *Darwinian account*: a relevant *population*; and *variation* and *selection mechanisms* for features that vary among individuals in the population. If the members of the population change over time, there must be some mechanism, corresponding to reproduction in evolutionary theory, that guarantees the stability of the frequency distribution of features in the population over time, absent the operation of the identified variation and selection mechanisms.

One notable example of a Darwinian account is Edelman's theory of neuronal group selection, which describes the construction during an individual's ontogeny of a neural organization that can support perceptual categorization, the basis for many innovative context-specific behaviors that the individual may generate during its lifetime. In Edelman's theory, the population consists of repertoires of neuronal groups, "collections of hundreds to thousands of strongly interconnected neurons."¹² These groups arise according to mechanico-chemical processes of cell division, movement, death and differentiation, which guarantee that "no two individual animals are likely to have identical connectivity in corresponding brain regions." These processes thus constitute the theory's variation mechanisms. Selection operates on the neural groups, as synaptic strengths increase or decrease, within and between the groups, in response to patterns of activation correlated via re-entrant signaling with sensory and motor activity. In this theory, the repertoires of neural groups are constructed via the variation processes and remain stable (unless they disappear) during the subsequent operation of selection, so there is no need for any mechanism corresponding to reproduction. Edelman's theory shows that biological evolution not only follows the Darwinian account, but that it "engineers" systems which themselves operate coherently with that account.

In general, the relation between successful theories based on a Darwinian account and the phenomena they purport to explain shares the characteristics we claim account for success of Darwin's theory: macrolevel function-carrying structure emerges from non-teleological interaction microprocesses; proposed variation and selection mechanisms are ontologically distinct and causally independent; fundamental aspects of the phenomena of interest are expressible in tractable mathematical or computational models. It is of course *conceivable* that a Darwinian account for some class of innovation phenomena might succeed even if none of the conditions of success for Darwin's theory hold, but it would seem prudent to

¹¹ For example, clonal selection in immunology (Jerne, 1967), neural Darwinism (Edelman, 1987), classifier systems for induction (Holland, Holyoak, Nesbitt, & Thagard, 1989). In addition, there are some interesting and successful Darwinian accounts for particular social phenomena, for example Croft's (2001) theory of language change.

¹² Edelman (1987), p. 5.

assign a very low *a priori* probability to such an outcome – and to look elsewhere for the foundations of a theory of innovation in such a context.

Now, neither R1 nor R2 seem particularly relevant for human sociocultural innovation. In contrast to the biological world, the fit between structure and function in the social world is not so evident that it cries out for explanation. Rather, it is *function* itself that seems to need to be explained: what is the functionality associated with such social constructions as cathedrals, horror movies and dog shows? This is not a question that a Darwinian account is equipped to address.¹³ Moreover, as far as R2 is concerned, it is quite unnecessary to demonstrate the existence of large-scale organizational innovation in the sociocultural context: we all *know* that social systems, and the kinds of agents and artifacts that inhabit them change, sometimes drastically and suddenly. Nor will a Lyellian analogy work to relate observable microprocesses to large-scale sociocultural innovation, for two reasons. First, it is unlikely that all or even most large-scale sociocultural innovation proceeds by the gradual accumulation of changes induced by microprocesses. Second, it is even more doubtful whether these microprocesses are themselves sufficiently stationary over long time scales that they could generate large-scale changes, without undergoing such significant transformations that their “observability” becomes irrelevant to predicting long-term effects.

We also can – at least provisionally – claim that R4, mathematization, has yet to tell in favor of a Darwinian account for sociocultural phenomena. As far as we know, the attempts to provide such an account have yet to introduce any new mathematics, beyond variations on population genetics and evolutionary game theory; and the mathematics that has been applied has yet to produce the kind of verifiable predictions and unifying conceptions that marked the work of Fisher, Wright, Haldane and their successors.

Thus, the main issue for a Darwinian account of sociocultural innovation is that raised by R3: do the foundations of a Darwinian account carve nature (or in this case, society) at the joints? We take this issue up in the next section.

1.3 A Darwinian Account for Sociocultural Innovation?

To give meaning to the question that provides the title to this section, we need to say something about what we mean by “sociocultural innovation”. As it happens, the previous section refers to an interesting, if somewhat surprising, example of what we have in mind.¹⁴ The example comes from Darwin himself, through the analogy

¹³ Unless it imposes the Procrustean bed of subordinated functionality and defines the “master functionality” to which cathedral building, horror movies and dog shows are subordinated.

¹⁴ Of course, Darwin’s theory and the processes through which it became scientific orthodoxy offer another good example. Though at first sight it might seem very far removed from the dog fancy – after all, one seems to be just about concepts and the evidentiary standards of scientific research, while the other seems to be about the generation of new breeds of dog – in fact it shares all the key features that we identify with that example.

he introduces between natural and artificial selection, and in particular the “fancies” that were coming into prominence as a popular pastime and commercial enterprise in Darwin’s epoch – and which so intrigued Darwin for the sheer variety of animal forms the fanciers were able to generate.¹⁵ In our discussion of this example, we highlight four features that are signatures of sociocultural innovation processes. After we describe these features and explore some of their implications, we confront our findings with the foundational requirements for a Darwinian account.

1.3.1 Artificial Selection and the Dog Fancy

As we saw, Darwin discussed artificial selection in *Origin* to highlight the wide range of heritable features that nature provides as grist for selection’s mill. Depending on the selection criterion used, artificial selection can generate varieties from the same ancestral stock that are eventually as dissimilar as, say, tiny Maltese dogs and huge Saint Bernards. For the purposes of Darwin’s argument, what is being selected is just a free variable: it could be any feature, as long as it is observable, heritable and stable over time. But if we are interested in analyzing artificial selection as a process of *sociocultural* innovation, then of course we need to understand what kinds of features are employed as selection criteria, and how they become established.

According to Darwin, there is a single primary functionality that underlies all selection criteria in natural selection: reproductive potential – the capacity to produce the maximum possible number of surviving and reproducing offspring. All other evolutionary functionality is *subordinated* to this primary functionality. For example, individuals must survive to reproduce, hence survival takes on (subordinated) evolutionary functionality in Darwinian terms – so long as the features that ensure it do not compromise reproductive potential. Similarly, the woodpecker’s sharp strong beak makes food gathering more efficient, and so enhances the survival of a (proto-) woodpecker that has this feature; thus, a sharp strong beak has (subordinated) evolutionary functionality for the woodpecker.

For much of human history, such subordinated functionality probably provides an adequate first-order explanation for the criteria employed in artificial selection as well. People used the plants and animals they domesticated and bred as *tools* to help themselves (or the social organizations of which they were a part) to carry out functions related to reproduction or survival (of the relevant individual or social organization): ensuring an adequate food supply, for example, or defeating enemies in combat. Any features that rendered domesticated plants and animals particularly effective with respect to such functionality could become a selection criterion for artificial selection – that is, the basis for differential treatment by their human masters that enhanced the fecundity of individuals possessing these features, either directly

¹⁵ Though Darwin was particularly interested in the pigeon fancy, we will mainly concentrate on the dog fancy, which has continued to grow in social and economic importance to the present day, although somewhere in the process the label “fancy” has largely disappeared from common usage – even though “Dog Fancy” is the name of a popular magazine for fanciers.

(e.g. by determining which seeds were planted, which animals were allowed to mate, or which offspring were not intentionally eliminated) or indirectly (e.g. by food allocation practices).

At some point, though, in the history of human interaction with domesticated plants and animals, a different kind of selection criterion emerged. These criteria were no longer subordinated to reproduction or survival. For example, some plants were bred to enhance the beauty of their flowers. Similarly, some breeds of dog probably arose from selection criteria related to their capacities to provide pleasant companionship for their human masters: already in ancient Rome, tiny Maltese dogs were simply household pets, noted – and almost surely selected¹⁶ – for their affectionate temperaments and useless but luxurious silky coats.¹⁷

In England, in the mid-19th century, a new kind of non-subordinated selection criterion began to emerge, associated with a new kind of social activity, the dog fancy. In 1859, the same year that the first edition of *Origin* was published, the first official dog show was held in Newcastle, followed a few months later by another in Birmingham in which 80 dogs (and their human handlers) competed in 14 different classes. These competitions proved very popular, and they rapidly grew in number, as well as in the number of competitors and the size of the public who attended them. In the earliest competitions, the dogs were judged with respect to their competence in performing class-specific activities related to such subordinated functionalities as pointing, retrieving or herding. However, this quickly began to change, and in the final decades of the 19th century, the most prestigious competitions had a completely different kind of criterion: the winners were those animals that were judged to best exemplify the “standard” conformation (physical and temperamental) of their class! Thus, the selection criteria, both for the judges in conformation competitions and for breeders seeking to produce winning animals, was not only totally unrelated to Darwinian primary functionality, but depended on *attributions* about the *attributions of others*: what judges, and fanciers in general, believed were the ideal features of a particular class, and what determined how close they believed a given animal might be to this ideal.

Such selection criteria of course required considerable alignment among the attributions of the people who participated in the competitions – owners, breeders, judges and the public who paid to view the conformance competitions – about just what counted as a class and what constituted the ideal conformation for each class. And this raises an exquisitely social question: where did these attributions come from, and how did they come to be sufficiently aligned among fanciers to make conformation competitions possible, attractive and profitable?

¹⁶ Probably from ancestors selected for their subordinated functionality of eliminating vermin that could attack the master’s food supply.

¹⁷ A status they continued to enjoy in the Renaissance (and beyond). Carpaccio painted a beautiful diptych, in one frame of which a group of men hunt ducks in the Venice lagoon with the help of water spaniels – and in the other two women dreamily await their men’s return on the terrace of a Venetian palazzo, accompanied by a fluffy, cuddly Maltese.

It is important to understand that, with just a few exceptions, what we now know as dog “breeds”¹⁸ did not pre-exist the rise of the dog fancy and the conformance competition, even if dog *breeding* did, as we have seen. As we saw, the classes in competitions were initially defined by function, and sometimes also by size, not genealogy. For example, in the earliest competitions, spaniels were divided into “springer” and “field” classes, depending on whether the animals were trained to flush game or simply locate and retrieve it. Because of the requirements of their task, springers were generally larger and more agile than field spaniels, but it was perfectly possible for a dog to compete in the springer spaniel class in the same competition in which its sibling was entered in the field spaniel class. Moreover, because of the variety in form of dogs entered in the field spaniel class, some competitions began to distinguish between larger dogs, called “field spaniels” and smaller ones, called “cocker spaniels” (supposedly because they were used to hunt smaller prey, like woodcocks). As *conformance* competitions increased in popularity, the issue of judging criteria for classes like these became particularly rancorous, since winning these competitions did not depend on what the dogs did, but how they appeared.

The solution to these questions of rules and definitions lay in *organization* – and *negotiation* channeled by organizations. Fanciers, especially breeders, established societies that debated and established procedures for determining rules for entering, classifying and judging competitions and conventions for sponsoring or recognizing competitions based upon these rules. The first such societies were based upon interest in particular classes of dogs, but soon the desire for overall coordination led to the creation of a new national society, the Kennel Club, founded in 1873, which appropriated the responsibility to oversee all “official” dog fancy competitions. The following year, the Club published its first Stud Book, which included results of past competitions and rules and calendars for future ones. Moreover, the Kennel Club, working with the class-based societies, began to establish conformation *standards* for the classes it recognized. To handle entries in a standardized way, each individual dog was restricted to membership in one particular class. The basis of this enrolment quickly became genealogy. That is, the *classes* were transmuted into *breeds*. In 1880, the Kennel Club began to register dogs as “purebred” members of the newly standardized breeds.

All this, of course, goes exactly in the opposite direction to Darwin’s move to *unreify* species. In effect, the dog fancy societies, coordinated by the Kennel Club, reified the *breed* – and endowed each breed with an ideal organization, expressed in its published conformation standard. This reification of the breed was not a recognition of some existing underlying *natural* reality, but rather, through the activities of the clubs, the competitions, and in particular the Kennel Club’s breed registry, it *created* a new *social* reality.

¹⁸ Defined by a particular set of characteristics that “bred true” and backed up by certified breed genealogies.

This reality, although it is based upon typological thinking, is far from static. The process of determining what constitutes a breed and its associated standard was – and still is – ongoing, and it can be highly contested. The history of what is now known in the US as the *English cocker spaniel* breed (and elsewhere in the world, as simply the *cocker spaniel*) illustrates this point. Judges in early competitions in the field spaniel class favored larger dogs, and as a result most breeders selected for size. An article in a dog fancy magazine in the early 1880s described the result: small spaniels exhibited in the last few years were just “weeds and wastrels of larger breeds,”¹⁹ and this kind of dog was well on the way to extinction. Fortunately for the many 20th century admirers of cocker spaniels, one small spaniel, Obo, enjoyed considerable success in competitions in the mid-1880s, and his stud services became in high demand. A group of Obo-philic breeders broke away from the Spaniel Club to form their own Cocker Spaniel Club, which in 1902 drafted the first cocker standard and pushed it through the Kennel Club’s ratification procedure. With its reality acknowledged and its ideal organization described, the cocker spaniel was positioned to attain instantiation in competition-winning champions, a public attracted by these dogs’ appearance and “active, merry” disposition (as the standard proclaimed), an organized group of breeders willing to produce animals for this public to purchase, and – as a result – a future.

While standards may seem to their drafters to provide a clear vision of a breed ideal, they are in fact subject to interpretation. Breeders interpret standards through the dogs they bring to competitions. Judges award victory to the competitors who come closest to their own interpretations of what breed standards mean – which, in fact, can even change as particular individual dogs reveal previously unimagined potential to display such standard features as these, from the Kennel Club’s current cocker spaniel standard:

- “General Appearance: Merry, sturdy, sporting; well-balanced; compact . . .”
- “Temperament: Gentle and affectionate, yet full of life and exuberance”
- “Head and Skull: Square muzzle, with distinct stop set midway between tip of nose and occiput. Skull well developed, cleanly chiseled”
- “Eyes: Full, but not prominent . . . with expression of intelligence and gentleness but wide awake, bright and merry; rims tight.”

And the dog-fancying and -buying public tends to seek out puppies with champion pedigrees, or at least which (the sellers assure them) will resemble the images of champions they have seen and admired. So new instantiations of the standard can change attributions about what the standard means, which can change what kind of features instantiations tend to display, and so on and on. As a result of this recurring feedback between changes in attributions and changes in instantiations, current English cockers don’t look at all like Obo: for example, according to Caddy (1995), Obo was 10 inches tall and weighed 22 pounds, while the current Kennel

¹⁹ Quoted in Caddy, 1995.

Club cocker standards (drafted in 1992) decree a height between 15.5 and 16 inches and a weight between 28 and 32 pounds for males of the breed.²⁰

Sometimes, the gap between conflicting attributions about a breed can be so great that the only recourse is through a re-negotiation of the text of the standard itself. In the case of cocker spaniels, the Kennel Club's initial standards lasted for almost fifty years. However, during this period, some American breeders began to introduce significant changes in the configuration of their cockers. They were less interested in producing dogs well adapted for hunting than in satisfying the growing demand for household pets, and so they selected for features that appealed to the non-sporting dog-loving public: "cute", human-like facial features and a glamorous full coat. Other breeders, committed to their attributions of a cocker as a small but powerful sporting dog, tried to resist this trend, as usual organizing societies and functional competitions devoted to pushing their attribution of what an ideal cocker should be, but the judges in the big, prestigious conformance competitions awarded "Best of breed" and even "Best of show" prizes to the newcomers, and the "American" cocker quickly became one of the most popular dogs in the US. Moreover, the American Kennel Club cocker standard was revised to favor the new type of rounder-faced, smaller, full-coated dog. The dispute among cocker fanciers and their clubs in the US about what was the "real" cocker was finally resolved through the creation of a "new" breed with its own standard, the "English cocker spaniel" – which of course was essentially the same as the Kennel Club's cocker spaniel. The Kennel Club in the meantime revised its own cocker standard to eliminate the "undesirable" American innovations, but as the great popularity of the newer version in the US began to spill over across the ocean, it too finally admitted the "American cocker spaniel" as a new authorized breed, joining its "older" cocker spaniel cousin.

To non-fanciers, all these activities – forming societies, sponsoring competitions, drafting breed standards, maintaining registries – can seem like an eccentric and socially marginal exaggeration of man's long-time relationship with the dog. But this is illusory: the "pet" phenomenon – the reinvention of the domestic animal as a companion rather than a servant – and its emergence in the past century and a half as a mass movement in Europe and North America is a sociocultural fact with increasing political²¹ and economic consequences. Conformance competitions and breed standards lie at the heart of the modern dog industry. In the US, American Kennel Club registration can make the difference between a puppy selling for several thousand dollars – or being given away for free. The American Kennel Club itself is a large and powerful organization, which sits at the apex of a hierarchy of many societies and clubs dedicated to the purebred dog. In fact, since its inception in 1884, its members are not individual dog fanciers, but organizations. According to its 2006 Annual Report, the Club now has 594 member organizations, sanctions and regulates over 20,000 events annually (including over 1500 conformation competitions), sponsored by nearly 5,000 affiliated organizations, and it registers nearly

²⁰ The American Kennel Club standard height for male English cockers is 16–17 inches.

²¹ For example, the animal rights movement.

a million puppies each year as purebred members of the 157 breeds it currently recognizes.

1.3.2 Some Features of Sociocultural Innovation

In subsections 1.3.2.1– 4, we describe four characteristic features of sociocultural innovation that the dog fancy story illustrates. But first, we introduce several concepts that play a key role in formulating these features: *artifact*, *attribution*, and *agent*.

- By artifact, we mean something that human beings produce for the use of (generally other) human beings. We thus use the term in a very broad sense: a purebred puppy, for example, is for us an artifact. Artifacts may be physical, informational or performative. A breed standard, for example, is an informational artifact – as is Darwin’s theory of evolutionary biology or the text of *Origin*. To be useful to others, informational artifacts generally require some form of physical or performative instantiation: a printed copy of *Origin*, for example, is a physical instantiation of Darwin’s text. A conformance competition is an example of a performative artifact.
- We will use the term *attribution* to describe how people or social organizations represent to themselves the entities that inhabit their world. In particular, attributions specify the identity of social agents and artifacts: for agents, what they do and how they do it; for artifacts, how they are made and used, by whom. People, and social organizations, interact with things in the world on the basis of the attributions of identity they assign to them.
- By *agent*, we mean an organization of human beings and artifacts, in the name of which social action is initiated and executed. We defer our discussion of the concept of *organization* to the final section of this chapter.

1.3.2.1 The Emergence of Functionality: From Interaction to New Needs

According to a functionalist perspective, pre-existing *needs* lead social actors to participate in *interactions*, which *satisfy* the needs that induced them. In processes of sociocultural innovation, the causal arrow connecting social interaction and needs can point in the opposite direction as well. In this case, the operative causal chains have a few more essential links: from *interaction*, to new *attributions*, to new *values*, from which new *needs* emerge. Once the new needs arise, they can induce the formation of new patterns of social interaction, through which social actors seek to satisfy them. In this way, new *functionality*, which is not subordinated to any pre-existing functionality, can come into being in the sociocultural world.

In the example, the dynamic we just described was initiated by the new activity of canine competitions. Initially, as we saw, dogs competed on the basis of their functional competences, but soon new attributions came into being, through which

dogs were seen as embodying an ideal type, which came to be identified with its “breed.” These attributions were then incorporated in breed standards, which represented new values: the purebred dog was no longer merely a tool or even a valued companion for its master; rather, he was an instantiation of his breed, more valued and more valuable the better he conformed to the particular fancier’s interpretation of the dog’s breed standard. These values led to certain needs – viewing, owning, breeding particular types of purebred dogs – appropriate to a new social role, the dog fancier, needs that found expression in such activities as attending competitions, purchasing, training and showing a dog, operating a kennel, engaging with other fanciers in breed-based organizations in activities designed to enhance the breed’s status and value and in discussions about how best to interpret (and if necessary to revise) the breed standard.

The preceding paragraph thus links new activities with the emergence of new attributions, which in turn give rise to new values, and, for those individuals recruited into the new social roles opened up by activities around these values, new needs. Like most stories of sociocultural innovation, our dog fancy story begins *in medias res*, since it does not explain why canine competitions were initiated in the first place;²² and it concludes inconclusively, since it finishes with some new activities that we might expect (rightly) will lead in their turn to new attributions, new values, new needs – and new activities. In sociocultural innovation, it is often the case that *one new thing leads to another*. We will discuss the modality through which this positive feedback dynamic plays out in section 1.4.3. Suffice it here to say that the “reverse functionalism” we describe here plays a critical role.

1.3.2.2 Agents and Artifacts: The Reciprocity²³ Principle

Our dog fancy story is about artifact innovation: the emergence of conformance competitions, breed standards, and purebred dogs. As we saw, these artifact innovations were accomplished by new kinds of agents – in particular, local and regional breed societies and national kennel clubs. These agents structured the negotiations and generated the rules through which participants to the fancy, human and canine, were recruited, their respective roles defined, and their activities coordinated. On the other hand, these agents’ activities were also made possible, indeed constructed around, an array of artifacts: from breed registries and registered dogs to

²² Though we could go farther back, and talk about an increasingly economically and politically challenged aristocracy developing cult activities that reinforce its social status – in particular, fox hunting, the new attributions about dogs and their breeding that masters of the hunt and their acolytes developed, the adoption of these cult activities and attributions by socially ambitious members of the rising bourgeoisie and their extension beyond the circles of those who actually participated in the hunt itself; and other streams of change in attributions, values and needs.

²³ We use the neologism reciprocity to indicate the relationship of reciprocal causation between transformations in agent space and transformations in artifact space. The correct nominative form “reciprocity” has already been appropriated to mean something quite different in the social sciences.

the blue ribbons and trophies that provided incentives and acknowledged success for show competitors and to the built sites, furnishings and rule books that structured agent interactions, from committee meetings to association conventions to dog shows.

We claim that this intertwining of innovations in artifacts and in the organization of the agents that make and use them is very general. We can express this claim in the form of the following **reciprocity principle**: *the generation of new artifact types is mediated by the transformation of relationships among agents; and new artifact types mediate the transformation of relationships among agents*. In particular, the reciprocity principle implies that any causally convincing narrative about artifact innovation will constantly jump back and forth between transformations in the space of agents and transformations in the space of artifacts. The proper domain for such a narrative, and for a theory of artifact innovation, is thus neither of these: rather, it is agent-artifact space (Lane, Malerba, Maxfield, & Orsenigo, 1996).²⁴

1.3.2.3 Tangled Hierarchies

Around half a century ago, Herbert Simon developed a theory of organization for complex systems, based on the idea of nearly decomposable hierarchy. According to Simon, complex systems are composed of entities and processes arranged in a sequence of nested hierarchical levels. Entities are recursively structured, in that level n entities are composed of components, which were entities of lower levels. Inclusion is strict: each level $n - 1$ entity can be a component of only one level n entity. Processes involve series of interactions among entities. Near decomposability implies that processes too can be localized hierarchically: they consist mainly of interactions among entities *at the same hierarchical level*. Moreover, each level is characterized by a particular spatial and temporal scale for its entities and processes respectively. In particular, this permits scientists to study (Simonian) complex systems level by level: to follow processes at level n , properties and configurations of entities at level $n + 1$ can be regarded as *constants*, since the processes through which they change are too slow to matter at the characteristic time scale for level n processes;²⁵ and level $n - 1$ processes are so rapid relative to this time scale that the

²⁴ It is perhaps worth reminding the reader that many current theories of artifact innovation do not respect this seemingly obvious consequence of the reciprocity principle. For example, the vast literature on S-shaped adoption curves from “innovation diffusion theory”, as well as most neoclassical economic research on “technological innovation”, regard only processes on agent space; while work in evolutionary economics around the idea of “technological trajectories”, as well as many attempts to explain such supposed empirical regularities as Moore’s and Wright’s Law, have to do just with transformations of artifact space.

²⁵ This of course does not preclude downward causation: different values of these “constants” can lead to very different outcomes for the processes under consideration, and the “constants” differ over spatial and temporal changes relative to *their* level’s characteristic spatiotemporal scales.

scientist can assume they are at their equilibrium values as far as level *n* processes in which he is interested are concerned.²⁶

We have already seen how the two fundamental classes of processes in the Darwinian account of biological evolution occupy two different levels in a Simonian hierarchy: variation processes happen at the genotypic (molecular) level, while selection processes take place at the phenotypic (individual) level. As we saw, the difference between the entities and time scales for these two classes of processes is critical for the Darwinian account's success in cutting nature at the joints.

The social world is also characterized by different ontological levels. In the dog fancy story, breeders belong to breed societies and other organizations that sponsor competitions, while these societies and organizations may be members of higher-level agents like the American Kennel Club. But social organization is not in general Simonian, for three principal reasons. First, strict inclusion of entities rarely holds. For example, the same breeder may be a member of many different dog fancy organizations. Second, social processes are often not localized in single ontological hierarchical levels,²⁷ and as a result near decomposability fails. For example, to follow the story of the emergence of the cocker spaniel, we have to move *back and forth* between events that involve Obo and his descendants, individual breeders and judges, competitions, breed clubs, and the Kennel Club. Third, even when processes can be assigned mainly to a single hierarchical level, the correlation between hierarchical level and intrinsic temporal scale (i.e. the "larger" the entities, the slower the processes) does not necessarily hold in the social world. That is, processes involving higher-level entities need not be slow relative to processes restricted to lower-level entities.²⁸

These non-Simonian properties of social organization have important implications for social science: they undermine the strategy of studying social systems "level by level".²⁹ In particular, they make it very unlikely that anything like the argument for the fundamental distinction between variation and selection processes that worked in the case of biological evolution – namely, that they involved

²⁶ Which again does not imply the absence of "upward causation": indeed, the relevant "equilibrium" might include the very *emergence* of the level-*n* entities under study from interactions among lower level entities!

²⁷ Despite the best efforts of social scientists to define them so: for example, the attempt to build a (sub)disciplinary divide between microeconomics and macroeconomics.

²⁸ Indeed, with the technological advances in communication and information processing over the last several centuries, and the capacity of higher-level organizations to exploit these technologies, social processes are now taking place at ever larger spatial scales, involving interactions among new higher-level entities, with increasingly *rapid* time scales.

²⁹ In fact, it is interesting to observe that the social sciences are organized very differently from the physical and biological sciences: the latter tend to divide their material by hierarchical level (elementary particle physics, atomic physics, condensed matter physics, chemistry; molecular biological, cellular biology, physiology, various whole organism specialties, ecology), while the first-order divisions in social sciences are functional rather than hierarchic-structural: anthropology, political science, sociology, economics.

interactions amongst entities at a different hierarchical level and hence with different characteristic time scales – could apply to sociocultural innovation.

1.3.2.4 Negotiation Structured by Rules Structured by Negotiation

A conformance competition is a set of social interactions, involving breeders, presenters, judges, dogs and the public. The interactions that constitute the competition follow a set of rules, which determine what, when and how each of the participants in the competition is allowed or required to act. Some of the rules are explicit (like the movements and poses through which the presenter exhibits her dog to the judges), some not (for example, when the members of the public should applaud and when they must be silent). Without rules to channel the interactions among the participants, a conformance competition – indeed any structured social interaction event – cannot happen.

Where do the rules come from? Many of them, and all the explicit ones, are determined through processes of negotiation, which take place within the organizations that sponsor and sanction the competition. In these negotiations, members of the organization, drawing on their own prior experience from their participation in the dog fancy as well as other domains of their lives, and directed by their attributions about what a conformance competition ought to be, propose alternative interpretations and argue about their relative merits, until the relevant group reaches closure. How these negotiations proceed – who is allowed to say what, to whom, for how long, in which illocutionary mode, and how closure can be attained – of course depend on some other set of rules, which are determined by the same organizations in which they happen, based perhaps on rules followed by other organizations and encoded in such manuals as Robert's Rules of Order. And the process through which these rules are determined proceeds through negotiation, structured by rules. . .

We can even think of the conformance competition itself as a kind of negotiation. Each presenter is offering her interpretation of what a breed standard means – the interpretation in this case is in the form of a suitably groomed dog going through its prescribed paces. The judges carry out their evaluations, in the light of their own interpretations of the standard – which sometimes might change, usually slightly, as the competitors enact the interpretations they embody. The result of this mute negotiation is a judgment, which will have its effect on the choices the presenters and others make for future competitions, about which dogs to show and how to show them – that is, on their interpretation of what the breed standard means.

From this point of view, our dog fancy story – and sociocultural change in general – is nothing but a story of negotiations structured by rules structured by negotiations, if we are willing to consider as a negotiation process any structured confrontation among social agents over alternative attributions about the structure of the agent-artifact space they jointly inhabit. As the lady said to William James about the succession of mutually supporting turtles that support the turtle that holds up the world on its back, it's negotiations all the way down.

In the next two subsections, we argue that a Darwinian account is foundationally inappropriate for phenomena of sociocultural innovation with features like those we have just described.

1.3.3 Are Variation and Selection Separable and Fundamental?

Think of the establishment of the cocker spaniel breed as a sociocultural innovation. What processes can we identify that produced this innovation and accounted for the multiplication of tokens of the ‘cocker spaniel’ type in the UK around the beginning of the 20th century? At least these: the particular genetic combination that produced an exemplary small spaniel, Obo, who could be perceived by some judges and breeders as representing the ideal properties of the spaniel breed as currently conceived, and who begat progeny sufficiently similar to him in appearance; the formation of an organization of breeders dedicated to breeding and showing small spaniels; the drafting of a breed standard restricted to small spaniels; the approval of this draft standard by the Kennel Club; the admission of registered dogs conforming to this standard in prestigious conformance competitions; the initiation and diffusion of the attribution that a small hunting dog makes a good household pet.

In this list, there are some elements that we could identify as variation processes: the genetic processes that give rise to an Obo, the introduction of a new standard that breeders and competition organizers can adopt as a basis for action; a new attribution of hunting dog as household pet. And some that seem like selection: judging in dog shows, choosing which kind of breed (if any) to raise, show or buy as a pet. But in both cases, the situation is very different from the case of biological evolution. The variation processes listed above happen at very different levels of organization and involve completely different mechanisms: where is the analytic bite in labeling them with the same term? Moreover, while a lot of selection is going on in this story, the criteria through which the selecting happens are changing on nearly the same time scale as the selections themselves, at least in the first and crucial stages of the establishment of the breed. The change in criteria, which depend upon new attributions of identity and functionality, themselves involve variation (the generation of new attributions) and selection (or more precisely aligning attributions among heterogeneous agents), so variation and selection are inextricably intermingled within what we have initially labeled as selection processes. The more intermingled they are, the less analytic value there is in considering them as distinct processes, since they cannot be analyzed separately – as was possible in the case of biological evolution, due to the ontological and temporal separation of the relevant entities and processes.

Worse, if we return to the list of processes contributing to the establishment of the cocker spaniel breed, several of the most important of them do not seem to be decomposable into variation and selection components, intermingled or not. For example, the formation of a new organization dedicated to promoting the (proto)breed was certainly a critical step towards the establishment of the breed. But it represents a variation only in the trivial sense that anything new is a variation, and it leads to no subsequent selection events, except again in the trivial sense that it itself didn’t

disappear. Rather, it is a construction, achieved through negotiations among a group of heterogeneous agents with aligned directedness (Lane & Maxfield, 1997), who were able to project the effect that such an organization might have in inducing changes in the attributions of other relevant agents about small spaniels and in carrying out activities consistent with their aligned directedness, such as drafting a breed standard and lobbying for its adoption by the Kennel Club. Moreover, the approval of this standard by the Kennel Club and the consequent establishment of a certified breed registry are organizational transformations, achieved through negotiations structured by the Club's rules – and these transformations, essential as they are to the innovation process under consideration, cannot be classified as either variation or selection processes, again except in the most trivial and analytically useless sense of these terms.

To summarize, variation and selection processes do not seem to carve this sociocultural phenomenon at its joints. While both kinds of processes do occur, they are not distinct with respect to ontological level and time scale, and, in particular because of the endogenous generation of new attributions of functionality and hence criteria for selection, they intermingle in a way that make them analytically inextricable. Moreover, other kinds of processes, in particular organizational transformation achieved through structured negotiations, seem even more fundamental in achieving the kind of sociocultural innovation in which we are interested. Indeed, if we look carefully at our list of possible candidates for variation and selection processes in the establishment of the cocker spaniel breed, we see that almost all of them are actually brought about through underlying processes of organization transformation and structured negotiation. These conclusions seriously undermine the possibility of a Darwinian account for this kind of phenomenon.

1.3.4 Where Is the Population?

Darwin succeeded because he un-reified the species. Un-reifying higher-level entities is undoubtedly a good thing: only real historical entities should play causal roles in accounts of historical processes. But not all higher-level entities are as causally inefficacious as species end up after their un-reification by population thinking. In particular, the dog fancy is supported by a multilevel set of organizations – from breeders, to breed societies, to national Kennel Clubs – which together form a system that is itself a kind of organization, as we explain Section 1.4. As we have seen, processes like establishing new breeds are enacted at the *system level* and rely on negotiations and other forms of structured interactions within and among the component organizations that comprise the system. These processes may transform the structure and functionality of the organizations of which the system is composed. Thus, to regard these higher-level systems as mere aggregations that passively monitor changes in frequency distributions of their components' properties is to ignore the most salient features of the dynamics of multilevel organizational change. Indeed, in our dog fancy example – and in the examples of innovation in urban and market systems discussed later in this book – it is rarely the case that the processes

we wish to understand can be localized to a single level of organization, never mind to a population of entities all inhabiting the same hierarchical level, as population thinking requires. As we argue in the next section, in these cases what we wish to study are examples of *organizational self-transformation*. When an organization transforms itself, where is the population?

1.4 Organization Thinking

We agree with Herbert Simon that complexity science stands in need of a theory of organization. As we have already observed, Simon's proposal for such a theory, based on his idea of nearly decomposable hierarchies, turns out to be unsatisfactory to account for important aspects of human social organization. Several strands of recent complexity research offer great promise in developing a deeper theory of organization: complex networks, modularity, degeneracy, scaling laws, and renewed approaches to hierarchy. While we cannot yet offer such a theory, it is our hope that some of the ideas presented in this book might contribute to its construction. In the meantime, in this final section of the chapter, we present some concepts and proposals that are intended to illustrate how organization thinking can provide a foundation for a theory of sociocultural innovation.

1.4.1 *What is Organization Thinking? Some Concepts*

The worlds that scientists study consist of a flux of energy, matter and information. The flux is generated by transformations, through which the patterns constructed from energy, matter and information change. These transformations result from interactions among these patterns. We call these patterns *organizations*, and the relations among energy, matter and information that organizations construct through their interactions we call *organization*. Organization thinking seeks to understand how these patterns form and transform through interaction.

We can describe organization in terms of the relationship among three fundamental aspects: *structure*, *process* and *function*. The *structure* of an organization describes its parts (energetic, material and informational), the interaction modalities among its parts, and the modalities through which the organization interacts with other organizations. The *processes* associated with an organization describe the transformations (in organization) in which the organization may participate. The *function* of an organization provides *directedness* to its actions, through its role in determining *which* processes the organization enacts, when it is in a context in which it is possible to enact more than one process.³⁰ None of these elements are

³⁰ Obviously, the concept of function is irrelevant for organizations that are never in such contexts – or if, when they are, chance rather than the organization determines which process is enacted. This is the case for physical systems. Function begins with biology.

necessarily static; indeed, organizations may have processes through which they themselves transform some or all of them.

Processes are supported by structure. To participate in a process, parts of the organization must engage in a sequence of interaction events, each of which requires some particular interaction modality. Instantiating the structural support for a given process may require the activation of *management processes*, of which there are three principal types: *recruitment*, which induces (perhaps even forms) the parts that will participate in the process; *differentiation* or *specialization*, which provides these parts with the requisite properties and interaction modalities; and *coordination*, which controls that requisite interactions happen in the right spatiotemporal order to achieve the appropriate transformation.

It is helpful to describe structure in terms of three subcategories. We will use the terms *representations*, *rules* and *relationships* to describe these subcategories, even though in some contexts some common meanings of these terms may carry inappropriate connotations. *Representations* comprise what we may call the organization's cognitive or classification system: they provide an organization with its view of the world it inhabits. *Rules* determine the organization's behavioral – that is, interactional – possibilities. *Relationships* arise from the history of the organization's interactions with other organizations; they describe how an organization is linked to these other organizations.

1.4.2 What's Special About Human Sociocultural Organization?

To begin to answer this question, we use the concepts described in the previous paragraph to describe three idealized types of organizations and the worlds they inhabit. These types are meant to represent in a very simplified, even caricatured way, physical, biological and human sociocultural organizations respectively, so we call them P-, B- and S-organizations.

In a world composed of P-organizations, the rules for each organization have the form of associated fields of forces, which taken together determine, perhaps with some randomness, the kinds and outcomes of the interactions in which the organizations engage with each other. The organization of such a world emerges from the interactions so determined. P-organization processes depend on structure (positions and forces) and chance. There is no need to introduce representations, function or management processes in a description of such a world.

B-organizations actively monitor the contexts in which they find themselves. They do this through representations in the forms of *categories*, whose levels or values the organization can register. These categories are then employed in condition-action rules of the form “If(cat) then(act),” where “cat” describes a context in terms of categories observable by the organization, and “act” describes a particular interaction modality that the organization can enact.

B-organization processes consist of chains of interactions arising from rules of this form; such chains can form if each act reliably generates the cat-condition for the next rule in the chain. When various parts of the organization are responsible for

different actions in such a chain, we can describe say that these parts communicate via “signaling” – since the transformations in context that result from one part’s act “signal” the condition that recruits the next part to make its contribution to the process enactment. There is no need for semantic interpretation here: the signal is not intended by the sender to refer to “something else”, which the receiver must infer correctly if it is to respond “correctly.” Indeed, the signaler need not have any representation that would indicate to it the existence of the receiver or the nature of the response to the signal, and vice versa.

Since certain contexts may trigger more than one rule, some of which might indicate mutually impossible interactions, B-organizations must have a way of deciding which triggered rule they will implement. This is the role of function: we can imagine function as some component part of the organization that assigns values to rules, which guide the management process that selects which eligible rule to enact. In this way, function provides directedness to the organization’s interactions.

S-organizations represent their contexts by modeling them: that is, they “populate their worlds” with entities, to which they assign attributions of identity (what kind of organization they are, what they do); moreover, they can operate in the “putative mode”, in which they use their representations (attributions of identity, plus a stock of narrative forms that express what happens when different kinds of entities interact through particular kinds of modalities) to “simulate” real interactions, projecting changes in context that result when particular sequences of interactions among entities take place. As modelers, S-organizations employ organization thinking: their attributions of identity to S-organizations include attributions of functionality. That is, in their models, S-organizations *do* things because they *want* something. In particular, they make attributions about *their own* functionality, and they monitor whether what they do tends to produce what they want. When it doesn’t, they seek to generate new interactions that, according to what they experience in putative mode, may do better.

So far, what we have described could be read as a simplified description of what human beings do – indeed, if we except the attributions of functionality to others, what some other animals do as well. But our description is not meant to be restricted to the individual level, and our claim is that human sociocultural organizations are the only *supraindividual* organizations that are capable of sustaining the kind of representation we have described. We will justify this claim soon, but first we want to emphasize how important it is. Almost all of what we human beings have achieved, in terms of the incredible expansion of artifact space and all its attendant phenomena as we sketched in the opening paragraphs of this chapter, have been accomplished by *organizations*, not by individual human beings – whatever we might conceivably mean by an individual human being operating independently of the sociocultural organizations of which he is a part. We are claiming that the capability of human sociocultural organizations to innovate depends on the representations, rules, relationships, management processes and function associated with these organizations, which are different from, and have vastly more transformative and generative capability than, those at the individual level. Just as a single neuron may contribute to

the expression and interpretation of a concept, but is not equipped to itself express or interpret one, so individual human beings contribute to the formation and enactment of organizational representations and processes, but cannot form or enact them themselves.

If representations of S-organizations do not reside in individual brains, where are they – and how are they related to interaction rules? Clearly, they are distributed among many brains – and also memory artifacts, including books, manuals and memos, and more recently computerized databases. The processes through which they are generated, modified on the basis of experience and exercised in the putative mode are many, but the most important and oldest of these is negotiation. Indeed, negotiation is the process that underlies many management processes in S-organizations. We discussed negotiation – and its recursive relationship with rules – in the previous section. Here, it will suffice to contrast it with signaling, which plays a homologous role in B-organizations. In negotiation, as opposed to signaling, semantics counts. The message that passes between sender and receiver may be completely novel in its form, and yet the sender anticipates that the receiver will be able to interpret it – indeed, that it will have the same meaning for the receiver as it had for the sender. On what is this expectation based, and how can it be right? The answer to this question is of course very complex; we describe a possible beginning point in the next chapter, based upon cognitive capabilities at the individual level and coordination requirements at the social level. Here, we note only that there is a bootstrapping involved, based upon the generative structures of language (and not only: certain kinds of joint action and communication by other symbolic representation systems based for example on pictures or number have similar negotiation efficacy) together with past experience in which negotiation led to mutually satisfying joint action. Most important, though, is the fact that many individuals who are part of the same organization sufficiently share the attributions and narrative forms on which their messages are based that sender and receiver can be sufficiently aligned around the meaning of the messages they exchange to generate agreement on appropriate joint action.

So negotiation can generate novel possibilities for joint action, at least when representations of the negotiating parties are sufficiently aligned to develop mutually comprehensible projects – and their directedness sufficiently aligned to make the projects attractive to all of them. But negotiations can also lead to the formation of new representations, when the negotiators' existing representations are not so closely aligned. This is due to the fact that for any human being or S-organization in general, it is very difficult to conceive of the possibility that its representation of the world is not the world! Representations change when this fact (since it must always be a fact, the world being what it is) becomes evident. Of course, occasionally the world has a way of making it evident, by producing big (and in general unpleasant) surprises. Another such occasion is when meaning breaks down during negotiations – that is, when one party to a negotiation uses an attribution or frames a narrative in a way that makes no sense in (or if it makes sense, is contradictory to) the framework of the representations of others. This sort of semantic uncertainty can lead the negotiators to “open up” attributional

space that previously was closed, and explore different possibilities that transform existing representations – not necessarily, of course, leading to agreement among them.³¹

In a supra-individual S-organization, the putative mode operates in general through negotiations (not only between people, but between people and mathematical and computational models, and among people, models, data analytic algorithms and data). Thus, as organizations explore their action possibilities in the putative mode, the generative capabilities of negotiation may result in the discovery of new interaction modalities and new entities with which to interact. As a result, rules for S-organizations necessarily have a more fluid, open-ended character than the condition-action rules that we posited for B-organizations. Indeed, for S-organizations, rules are better described in terms of *permissions*: who may (or may not: unusually, here we use the word permission in a negative as well as positive sense) interact with whom, about what, in which interaction modality. Many permissions are expressed explicitly in S-organizations, through the hierarchical command structure that is part of the structure supporting various coordination processes in many of these organizations. Many more are shared attributions within the S-organization, not explicitly stated but accepted nonetheless. Particular S-organizations (including of course individuals) may of course arrogate to themselves the permission to engage in an interaction; when they do, this permission may be contested by other S-organizations and then negotiated, with the negotiations structured by the rules of the organization responsible for coordinating action among the relevant disputants.

There is of course a great deal more to be said in answer to the question raised by the title of this subsection. Some of these issues will be addressed in subsequent chapters of the book; a general answer awaits future research. Even in the stylized form of this comparison between three caricatured organizational types, it should be evident that S-organizations are constantly transforming themselves and their relationship with other S-organizations, through the negotiations in which they engage and the permissions they arrogate to themselves or grant to others. We have already seen examples in Section 1.3.1 of the fact that S-organizations also generate new attributions of functionality, for themselves and for other S-organizations. They may then attempt to realize this new functionality by means of the development of new artifacts. We now return to this theme, concluding the chapter by describing a positive feedback dynamic for innovation in agent-artifact space.

1.4.3 Positive Feedbacks in Agent-Artifact Space: Exaptive Bootstrapping

We already noted in our dog fancy story that one new thing leads to another: innovations occur in cascades, and involve transformations not only in artifact types, but

³¹ See Agar (1994) and Lane and Maxfield (2005) for extended discussions of this idea.

in organizational forms and attributions as well. In this section, we sketch the theory of exaptive bootstrapping, which explains how such cascades happen. The theory, based on organization thinking, provides a qualitative description of a positive feedback dynamic in agent-artifact space, which we claim accounts for the explosive growth in that space that characterizes human sociocultural change, particularly over the past several centuries. The recognition of the importance of the positive feedback dynamic for artifact innovation and its implications for organizational innovation, including the growth of cities, underlies almost all the research discussed in this book and represents its principal organizing theme.

We begin by distinguishing between two different kinds of invention activities: those that are intended to deliver an existing functionality “better-faster-cheaper” than the artifacts that currently do so, and those that are designed to deliver *new* kinds of functionality. An innovation cascade can be initiated by either type of invention, and in any cascade, both types are present.

In our dog fancy story, the cascade was initiated by the first dog show in 1859, which was intended to deliver a functionality not previously provided by other performative artifacts. For an example of a cascade that began with a better-faster-cheaper invention, we recall one of the most important innovation cascades in human history, which began with the invention of printing by movable type. This was a “better-faster-cheaper” innovation: Gutenberg’s workshop figured out how to produce multiple copies of a manuscript more quickly and cheaply than was possible with the previous method (hand-copying). But almost immediately, the first printing enterprise, headed by Gutenberg’s ex-partner Fust and ex-assistant Schoeffer had to solve a series of organizational and business problems that required new attributions of functionality: for agents, who had to pay up front for the paper for over a hundred copies (soon hundreds to a thousand or more) of a text, before selling any of them, and needed to work out new techniques for financing, selecting, marketing and selling their products; and for artifacts – what kinds of texts to print, and how to present them, in order to attract new customers who could not afford hand-copied manuscripts, but could pay enough for the right kind of printed book. And the solutions that the early book producers developed to these problems established new kinds of texts (and hence “reading functionalities”) that in turn induced the development of better-faster-cheaper improvements and novelties, in both the physical and informational forms of books.

Though typically innovation cascades contain both types of innovation, we claim that the positive feedback dynamic depends on the existence of the second kind – in particular, on the role of new attributions of functionality in bringing these about. The theory of exaptive bootstrapping posits the following stages for the positive feedback dynamic:

1. New artifact types are designed to achieve some particular attribution of functionality.
2. Organizational transformations are constructed to proliferate the use of tokens of the new type.
3. Novel patterns of human interaction emerge around these artifacts in use.

4. New attributions of functionality are generated – by participants or observers – to describe what the participants in these interactions are obtaining or might obtain from them.
5. New artifacts are conceived and designed to instantiate the new attributed functionality.

Since the fifth stage concludes where the first begins, we have a *bootstrapping* dynamic that can produce cascades of changes in agent-artifact space. These cascades inextricably link innovations in artifacts, in organizational structure, and in attributions about artifact and organizational functionality.

Exaptation happens between the third and the fourth stage in this process, whereby new attributions of functionality arise from observing patterns of interaction among agents and already existing artifacts. The idea here is that artifacts gain their meaning through use, and not all the possible meanings that can arise when agents begin to incorporate new artifacts in patterns of use could have been anticipated by the designers and producers of those artifacts: the combinatorial possibilities are simply too vast when a variety of different agents intent on carrying out a variety of different tasks have available a variety of different artifacts to use together with the new ones – not to mention that the designers and producers do not share the experiential base and the attribution space of all the agents that will use the artifact they produce, in ways that depend on their experience and attributions, not those of the artifact’s designers and producers! Meaning in use is one thing – the *recognition* that that meaning might represent a functional novelty is another. For this to happen, some participants in (or observers of) these patterns of interaction must come to understand that something more is being delivered – or could be delivered, with suitable modifications – to some class of agents (perhaps, but not necessarily, including themselves) other than what the participants were thinking to obtain through the interactions in which they were engaging – and which these agents might come to value. Thus, the generation of new attributions of functionality is grounded in an *exaptation*: from the interactions between existing structures (agents and artifacts), new functionality emerges. It may then become recognized by appropriately situated and motivated agents, and (re)cognized as a new attribution of artifact functionality.

To illustrate the stages described, consider the following example from the early days of printing. In this example, stage 1 corresponds to the printed book, and stage 5 to the printed advertisement. The linking stages can be summarized as follows. Before printing, almost all manuscripts were produced in response to orders from a commissioning agent. Not surprisingly, this was initially the case also for the first printing firm, established in Mainz using the printing technology developed by Gutenberg and his co-workers, which was headed by the financier Johann Fust and the printer Peter Schoeffer (Gutenberg himself was an early example of an inventor who failed to make the transition to innovating entrepreneur). Fust and Schoeffer had one important client, the archdiocese of Mainz, which commissioned many works from them, including religious books, references in canon law, and texts for the new humanistic school curriculum in which their clerical workers

were trained. Fust and Schoeffer realized early on that they could probably find purchasers for additional copies of these books. They faced the problem of how to reach these potential purchasers and convince them to buy the printed books. One organizational solution to this problem that the firm explored was to hire traveling representatives, which constituted stage 2 of the exaptive bootstrapping cycle. These representatives of course visited fairs and festivals, but they also stopped at towns along their route. When they did so, they would have to make known to potential purchasers their whereabouts and their wares – cycle stage 3. One approach that the firm took to this problem was exapted from their primary ongoing activity, in cycle stage 4: they conceived the idea of using printing, the same technology they employed to produce their wares, to enhance distribution. The new artifact type they developed (stage 5) was the printed advertisement. Their earliest surviving printed advertisement dates from 1469. It is a one page broadside, which begins as follows: “Those who wish to purchase for themselves the books listed hereafter, which have been edited with the greatest care and which are set in the same Mainz printing type as this announcement . . . are invited to come to the dwelling place written in below” (Lehmann-Haupt, 1950). Thus, the advertisement attests not only to the nature of the wares (the list of books that it provided), but also to their quality (the “same Mainz printing type as this announcement”). Note that the name of the inn where the representative could be found had to be hand-written, as it changed with time and town. The printed advertisement instantiates the new attribution of functionality: the possibility of mass-circulating information about a product to recruit potential purchasers. Other instantiations of this attribution, for other classes of products, followed, and the circulation of printed catalogues soon became an important means of disseminating product information and organizing exchange activities.

Innovation cascades involve many cycles of the exaptive bootstrapping process. In addition, these cascades also include processes that are purely adaptive: given an attribution of functionality and an artifact that realizes it, apply a known technology to improve the artifact or its method of production to render it better (according to the values associated with the given attribution of functionality), faster or cheaper. Such processes do not require the generation of new attributions of functionality. Note, though, that better-faster-cheaper invention is not necessarily purely adaptive. Many require new attributions of functionality as well: for example, Gutenberg had to exapt a variety of techniques he had learned as a jeweler in quite different contexts, even with different materials, for the new functionality of type-casting. In such cases, not only the exaptation of new attributions of functionality, but also organizational transformations like those in stage 2 are required, for example in assembling a team of agents that collectively embodies the different competences necessary to achieve a complex better-faster-cheaper invention – and in developing the procedures whereby this team can sufficiently align their directedness and then attributions about each other and the artifacts with which and towards which they work to accomplish what they have come to intend to do together.

References

- Agar, M. (1994). *Language shock: Understanding the culture of conversation*. New York: Harper Collins.
- Beinhocker, E. (2006). *The origin of wealth: Evolution, complexity, and the radical remaking of economics*. Boston, MA: Harvard Business School Press.
- Caddy, J. (1995). *Cocker spaniels today: Book of the breed*. Lydney, UK: Ringpress Books.
- Croft, W. (2001). *Explaining language change: An evolutionary approach*. Harlow, UK: Longman Linguistics Library.
- Desmond, A., & Moore, J. (1991). *Darwin: The life of a tormented evolutionist*. New York: Warner Books.
- Edelman, G. (1987). *Neural Darwinism: The theory of neuronal group selection*. New York: Basic Books.
- Holland, J., Holyoak, K., Nisbet, R., & Thagard, P. (1989). Induction: Processes of learning, inference and discovery. Cambridge, MA: MIT Press.
- Jerne, N. (1967). Antibodies and learning: Selection versus instruction. In: G. Quarton, T. Melnechuk, & F. Schmidt (Eds.), *The neurosciences: A study program* (pp. 200–205). New York: Rockefeller University Press.
- Lane, D., Malerba, F., Maxfield, R., & Orsenigo, L. (1996). Choice and action. *Journal of Evolutionary Economics*, 6, 43–76.
- Lane, D., & Maxfield, R. (1997). Foresight, complexity and strategy. In: W. B. Arthur, S. Durlauf, & D. Lane (Eds.), *Economy as a complex, evolving system II* (pp. 169–198). Redwood City, CA: Addison-Wesley.
- Lane, D., & Maxfield, R. (2005). Ontological uncertainty and innovation. *Journal of Evolutionary Economics*, 15, 3–50.
- Lehmann-Haupt, H. (1950). *Peter Schoeffer of Gernsheim and Mainz*. Rochester, NY: Leo Hart.
- Mayr, E. (1959). Darwin and the evolutionary theory in biology. In: B. Meggers, *Evolution and anthropology: A centennial appraisal* (pp. 1–10). Washington, DC: Anthropological Society of Washington.
- Mayr, E. (1991). *One long argument: Charles Darwin and the genesis of modern evolutionary thought*. Cambridge, MA: Harvard University Press.

Chapter 2

The Innovation Innovation

Dwight Read, David Lane and Sander van der Leeuw

2.1 Introduction

As humans, we are the only species that reflects consciously on our existence and how we came to be. Such musings have led us to formulate many different scenarios that see us as coming into existence through a creative act by forces outside of ordinary experience. However, within the domain of scientific reasoning, any appeal to such extraordinary forces is excluded. We therefore seek a natural account of how a species as complex as ours, capable of formulating and realizing the widely diverse forms of social systems that we know, could have arisen. Such an account must be embedded in the Darwinian paradigm for evolution, which has been fundamental to our understanding of the way in which biological reproduction can drive change from simpler to more complex biological forms.

The Darwinian evolutionary argument owes its success to Darwin's realization that the engine of reproduction – necessary for the continuation of life forms – is both the location of innovation in the traits that make up an organism, and the driver for change in the distribution of traits in a population. By coupling innovation (in the form of mutation in the genetic material transmitted) with differential reproductive success, Darwinian evolution connects patterning expressed at the level of the individual (novel traits) with patterning expressed at the aggregate level of a population (frequency of traits). Both are components of a single system in which change is driven by the environmental and social conditions responsible for differential rates of reproductive success.

That Darwinian evolution can account for changes in the frequency distribution of a mutation-induced trait in a population is not in question here. Less clear, though, is whether macro-level patterning in the organization of traits within an organism is a reflection of the micro-level of trait occurrence and trait frequency distribution. Or to put it even more broadly, whether collective functionalities arising from systematic

D. Read (✉)

Department of Anthropology, University of California at Los Angeles, 341 Haines Hall,
Los Angeles, CA 90095-1553, USA
e-mail: dread@anthro.ucla.edu

organization of the behavior of individual organisms emerge solely through a Darwinian process that changes the frequency distribution of traits in a population. If, instead, there should be innovation that allows for organizational change through endogenous processes acting on an assessment of current organizational functionalities, then a fundamentally non-Darwinian form of evolutionary change will have come into play. We argue that such an “innovation innovation” did take place during hominin evolution and that it is the basis for the forms of social organization we find in human societies today.

Conceptually, we will develop the argument in two parts. First we discuss the organizational implications of a process – enculturation – that is critical for the transmission of the cultural framework at the core of human social organization. We will show that even though change in social organization begins with properties arising from Darwinian evolution, once enculturation became the means of transmission of cultural resources, our species acquired the ability to construct and transmit forms of social organization in which individual functionality derives from organizational functionality. Such organizational functionality is subject to endogenous change by the individuals involved, and can therefore introduce new functionalities to cope with changing conditions in a manner independent from Darwinian evolution at the level of individuals.

Then, we will discuss a possible evolutionary pathway that may have led to this fundamental change in the basis for societal organization. That pathway leads from social learning, through face-to-face interaction, to the ability to anticipate patterns of behavior from a system of conceptualized relations among group members. This pathway, we argue, arises from a shift in cognitive abilities that enables (1) categorization on the basis of conceptual relations between individuals, and (2) construction of new relations through recursive reasoning (such as “mother of a mother”). Social organization based on such a system of conceptual relations decouples societal organization from biological kinship. Behaviors can be associated with relations, and thereby become part of the interaction of individuals in a network that is itself constructed through the composition of relations. Thus, distribution of social behaviors among a society’s members is no longer dependent upon Darwinian processes enacted at the individual level.

2.2 Organization of Behavior and Collective Functionality

Although the change in cognitive capacities that enabled this fundamental change in humans’ social organization arose through Darwinian evolutionary processes, the changes themselves imply a fundamental shift in the basis for social structure, from the phenomenological domain of traits to the ideational domain of concepts and relations among concepts. Rather than arising from a genetic substrate, these relations are structured and organized through systems of rules that are part of the informational structure we refer to as culture. They provide the framework within which human behavior takes place and frame the interpretations made by societal members of the behaviors of others.

Any social system must combine (1) a means through which the organization of behavior gives rise to collective functionality and (2) a means to perpetuate the social system long enough for functionality benefits to accrue to group members. Behavioral organization is expressed through patterns of interaction among group members. This interaction combines behaviors that are differentially expressed among group members with collective coordination to ensure group functionality, which in turn adds to the functionality of the individuals concerned. *Individual functionality* refers to the consequences for an individual of the range of actions in which he/she is involved. Such consequences can be material when actions are directed towards the phenomenological environment and behavioral when they directed towards the actions of others. *Group functionality* refers to the actions and consequences, including those at the ideational level, that accrue to individuals through membership in a group organized as a social system.

To illustrate this, consider a female/male dyad that forms in a sexually reproducing species – even if temporarily – for the purpose of sexual intercourse. The dyad forms a social group since the action of one member of the dyad affects – and responds to – the behavior of the other. Moreover, the behaviors engaged in by both members are coordinated during their interaction. Group functionality therefore refers to functionality not available to an individual outside of the social group, such as sexual reproduction, and the result of that functionality, namely the production of offspring, adds to individual functionality by increasing individual reproductive fitness.

Whatever may be the coordination of behaviors expressed through, and the functionalities derived from, a social system, both ultimately arise from individual properties that are consequences of Darwinian evolution. But as we are here concerned with proximate rather than ultimate explanations, we need to focus on the means of transmission of the basis for such behaviors from one individual to another. In fact, we can distinguish three modalities for social organization, according to the mode of transmission of behaviors: genetic transmission, individual learning, and enculturation. And as one moves from genetic transmission to enculturation, there is a gradient from individually enacted to socially constituted behaviors that defines the conditions under which effective behavior transmission can occur.

2.2.1 Social Organization and Mode of Behavior Transmission

We next consider some of the implications of each mode of transmission for the organization of behavior and the functionality that arises from it.

2.2.1.1 Social Organization Based on Genetically Transmitted Behaviors

Consider first genetically transmitted behaviors (see Fig. 2.1), i.e. behaviors whose expression arises primarily, if not entirely, from specification at the genotypic level. For the individual, such behavior is enabled through genetic transmission and does not require any social unit other than a copulating couple. If at all, the

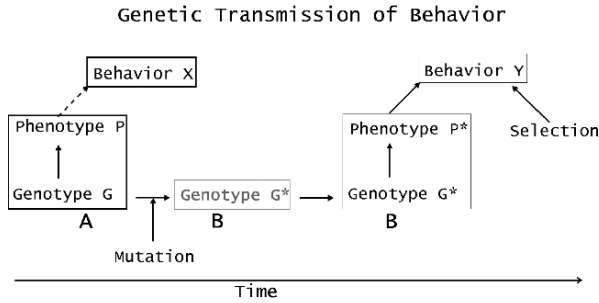


Fig. 2.1 A mutation in the transmission of genetic information from A to B leads to a changed genotype that gives rise to a changed phenotype with a different behavior Y. Selection acts on individual B via the fitness consequences of behavior Y. For diploid, animal organisms a social context other than a reproducing dyad is not necessary for transmission to take place

social dimension arising from organized interaction of individuals comes into play only after genetic transmission has taken place. Coordination of behaviors, and thus the predictability of the behavior of one individual with respect to that of others, derives from differentiated distribution of genetic material over individuals, and is therefore dependent upon the system that structures that distribution.

Social insects are prototypic examples of social systems based on genetically transmitted behavior. The organization of a colony of social insects is derived from a reproductive system that has been co-opted by a single female, the queen, so that her fitness is determined by the functionality of the colony as a whole. The individual functionality of all other females is shifted away from reproductive behavior to behavior that serves the functionality of the colony as a social system. Less extreme than the social insects are mammalian and primate social systems where social organization may also be framed around genetically transmitted behaviors, but without the extreme co-option of reproduction that occurs in the social insects. Old world monkey social units, for example, are often based on a dominance hierarchy for females constructed around genetic mother – daughter linkages that give the social unit cohesion and stability. Such linkages are emphasized through a residence pattern in which female offspring remain in their natal group and sexually mature male offspring migrate to other groups (Pusey & Packer, 1987). Such a social system is based on genetic transmission of behavior, and thus stays within the framework of Darwinian evolution – expanded to include fitness based on interaction with biological kin and sexual selection – since coordination of behavior is embedded in the genetic system (Mitani, Watts, Pepper, & Merriwether, 2002). Perpetuation of such a social system primarily depends on maintenance of a mating system and a pattern for the distribution of adult individuals across residence groups.

2.2.1.2 Social Organization Based on Learned Behaviors

When the linkage between the genotype structure and its phenotypic expression is more relaxed, behaviors are increasingly expressed through individual properties at

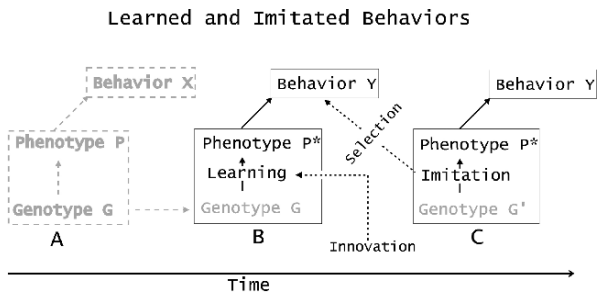


Fig. 2.2 Learning is a source for innovation in behaviors. In the diagram, B is the genetic offspring of A. Behavior X otherwise associated with genotype G and phenotype P has changed to behavior Y in individual B through learning by B. A learned behavior can spread through a population through imitation. Individual C, with possibly a different genotype G', takes on phenotype P* and behavior Y through imitation of B. The imitation process has characteristics due to prior Darwinian selection that frame the conditions under which imitation will occur. Selection arises through evaluation of B as a possible target for imitation by C according to whether the conditions under which imitation will occur are satisfied. Evaluation by C may be based on characteristics of B (“imitate successful individuals”) or it may be based on the consequences that arise from doing behavior Y (“imitate behaviors that lead to a positive reward”) independent of the other characteristics of B. Imitation selection is decoupled from fitness selection

the phenotypic level, and individual learning plays an increasingly important role in the formation and organization of social groups (see Fig. 2.2). Coordination of behavior may now derive from individual learning as well as genetic transmission. Primates are a prototypic example of such social organization, worked out through face-to-face interaction among group members. Among the Old World monkeys, extensive interaction between an infant female and her female biological kin play an important role in determining her position in the female dominance hierarchy as she matures into adulthood. Such social organization requires that offspring be engaged in interactions in which individual learning takes place, so that a new group member becomes incorporated into the group behavior patterns upon which social cohesion is based, and from which collective functionality arises. Consequently, the continued existence of the social group depends upon social interaction, even though the phenotype is being developed through individual learning. Individual learning also leads to behaviors that become part of an individual’s phenotype separate from behaviors acquired through social interaction.

Individual Learning and Darwinian Evolution

When behavior is derived from individual learning, innovation and change can arise outside of genetic mutation, through novel learned behaviors (see individual B, Fig. 2.2). Though innovation by learning plays an analogous role to that of mutation in genetically transmitted behaviors, fitness selection acting on mutations does not have a simple counterpart with individually learned behaviors. While change at the aggregate level of a population can occur as a consequence of individually

learned behavior, thereby affecting the individual learning of other group members, this is not Darwinian selection in the strict sense. However, when there has been selection for changes at the phenotypic/cognitive level that enable one individual to imitate the behavior of another (in the sense of functionally repeating the behavior in a manner consistent with obtaining the outcomes associated with it, rather than merely mimicking it (Shettleworth, 1998; Tomasello, 2000)), this would constitute a strict analogy to fitness selection. But such selection directly derives from how, and under what conditions, imitation takes place (Boyd and Richerson 1985) and is only indirectly related to measures of fitness such as reproductive success (see individual C, Fig. 2.2). Hence, though the means by which innovation and selection take place *are not* identical to mutation and fitness-based selection, if the structural property of innovation leads to patterning at the level of the individual, and operates independently of selection that leads to change at the aggregate level, it *is* functionally the same. Therefore we will include innovation in individual learning and selection through imitation under the umbrella of Darwinian evolution.¹

Individuation Versus Social Cohesion

Regardless of the means by which the organizational structure and its transmission are achieved, all social groups face two problems: (1) how to accommodate (or reject) novel behaviors introduced through mutation or innovation, and (2) how to coordinate group behavior so as to reduce individual conflict within a group. Variation and novelty are problematic for social organization because they introduce behavior that is unpredictable by other group members, and may thus interrupt the coordination from which collective functionality arises.² Social systems, though, may have to accommodate novel behaviors due to biological selection for more complex neurological systems that can process external information and generate novel behaviors (involving learning from past experience and from interaction with other individuals, including imitation) (see Fig. 2.2, individuals B and C). One means to accommodate such novel behaviors is through the cognitive ability to predict, with sufficient accuracy, the behavior of other group members so that an individual may modify its behavior in anticipation of the behavior of others.

With the advent of the primates and especially the evolution of the pongids, the ability to make predictions about behaviors of other group members, including third party group members, has become a regular part of the cognitive repertoire (Tomasello, 1998). Hence, more complex forms of social organization, with

¹ In contrast to social organization based on genetic transmission of behaviors, where group membership is framed around common genetic ancestry, the boundary of a social group organized through interaction of group members is more complex and may bring into play conflicting factors regarding group boundaries. These may be accommodated by the conditions under which individuals can transfer from one group to another.

² At one extreme, accommodation of individual differences and individual learning has been resolved in the negative by the social insects through reducing individual variation by the queen controlling reproduction and the absence of individual learning.

a collective functionality that ensures a fitness payoff for individual group members, have become possible. This interdependence between cognitive capacity and complexity of social organization has been discussed, following seminal papers by Chance and Mead (1953), Kummer (1967) and Humphrey (1976), by a number of researchers under the rubric “Machiavellian Intelligence” (see papers in Whiten and Byrne (1988), Byrne and Whiten (1997b)). Social complexity has been seen as a driving force for increased cognitive capacities among the pongids and hominins (and possibly other social mammals) (Dunbar 1995). Increased cognitive capacity for varied behaviors and the capacity to modify one’s own behavior in expectation of the likely behavior of others has also been posited as an impetus for increased individuation of behavior. But individuation poses problems for social cohesion because it increases the complexity of the social field in which individual group members interact (Read 2004) and augments the potential for conflicts between individuals that disrupt the social units. In the absence of sufficient mechanisms for controlling conflict (Flack, Girvan, de Waal, & Krakauer, 2006) or resolving it (see de Waal, 2000), such conflict can lead to smaller, less diverse and less integrated social units (Read, 2005).

By individuation of behavior we mean expansion of the total behavioral repertoire of group members to the point where the behavior of an individual targeted for interaction cannot be induced accurately from experience with the behavior of other group members. Assume we have a group G of n individuals and let $B_i = \{b_{i1}, b_{i2}, \dots, b_{im_i}\}$ be the repertoire of behaviors that can be engaged in by individual g_i in G , which may have an impact on the functionality of behaviors of other group members. Lack of individuation will correspond to low diversity across the sets B_i and extensive individuation will correspond to high diversity across the sets B_i .

We can then define as a simple society one in which the sets B_i have low diversity, different individuals exhibit essentially the same range of behaviors, and knowledge about one individual’s behavior can successfully be applied to predict the behavior of other individuals (Read, 2004). The group size of simple societies (such as a school of fish or a herd of ungulates) will tend to be scale-free due to low diversity of behaviors across individuals. In such societies, an individual’s social field will be determined by the variety of behaviors within a single, summary behavior set B , rather than by the number of individuals. As a first approximation, if we assume that individuals can cope with all the behaviors in the behavior set B , the complexity of the social field will tend to be independent of the size of the group.

In these terms, a complex society can be defined as one with high diversity of behaviors so that the experience one individual has with another individual may only of limited use in predicting the behavior of yet other individuals. Under these circumstances, the degree of social coherence will be related to the number of individuals for whom behavior is predictable, which will in turn be related to the total number of individuals; hence, social coherence in complex societies based on individual learning will not be scale-free and, all other things being equal, social

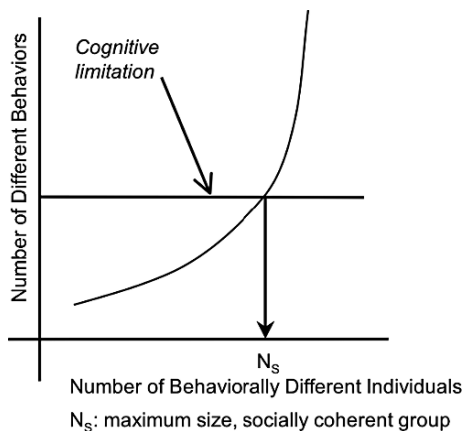


Fig. 2.3 Relationship between degree of individuation (number of behaviorally different individuals) and social complexity (number of possible, different individual, dyadic, triadic, etc.) behaviors. Limit to cognitive capacity for dealing predictably with different behaviors places an upper bound on the maximum size of a socially coherent group

coherence will decrease with group size.³ In a complex society, therefore, the complexity of the social field will scale with the number of individuals in the group *plus* the number of dyads (since a dyad can form a temporary alliance vis-à-vis a third individual), the number of triads and so on. As noted by Byrne and Whiten (1997a: 11): “a monkey, taking the probable actions of a third party into account, is facing a more challenging world than an animal that only interacts dyadically” Even if we only take dyads into consideration, the complexity of the social field in a complex society will scale with n^2 , n being the size of the group.

Individuation and Cognitive Limitations

Ability to cope with diversity in behavior depends in part on individual cognitive capacities, and this is reflected in the non-human primates by a positive correlation between innovative behavior and executive (neocortex) brain size (Reader & Laland, 2002). Every species has an upper bound to its cognitive capacities, and (as is shown in Fig. 2.3) that in turn bounds the size of a complex society dependent on coordination of behaviors for its coherence (see also Dunbar (2003)). Hence, in the absence of any new mechanism that enables accommodation of the complexity inherent in increased individuation, the latter will lead to a decrease in the mean group size of coherent social units. The effectiveness of face-to-face interaction for accommodating individuated behaviors diminishes rapidly with increased individuation. This has occurred with the non-human primates who depend extensively on

³ The definition of simple and complex societies introduced here is consistent with the concept of complex systems developed under the rubric Complex Systems Science (see Bourguine & Johnson, 2006).

Table 2.1 Meat Sharing (1 eland), *!Kung san* Hunter-gatherer group, Kalahari Desert, Botswana

Genetic Kin ¹	Number of sharing instances	Non-genetic Kin	Number of sharing instances	Number of different residence groups
Biological Parent	2	By marriage	20	6
Biological Sibling	3	Uncertain	6	
Biological Cousin	9	Other	20	
Other	2			

¹ Biological relation inferred from kin term usage
 Data from Marshall (1976), pp. 300–302

Table 2.2 Food Sharing (Bananas), *Pan troglodytes* (Gombe Stream Reserve)

Group size: <i>n</i> = 37	Biological mother/offspring Connection	No biological connection	Total
Number of dyads	33	625	658
Number of sharing instances	360 (mo → o) 31 (o → mo)	47 (m → f) 17 (other dyad)	457
Rate of sharing (based on numbers in bold)	11.8 instances/dyad	0.1 instances/dyad	Ratio of rates: 100:1

Data from McGrew (1992), pp. 107–108 and Fig. 5.10

face-to-face interaction for social integration. Data on the pongids (orangutans, gorillas and chimpanzees), *Ceboids* (New World monkeys) and *Cercopithecoids* (Old World monkeys) suggest that the pongids show increased individuation as well as a reduction in the size of their social units (see Tables 2.1 and 2.2 in Read, 2005). One pongid, *Pongo pongo* (Orangutan) has reverted to solitary foraging while another, *Pan troglodytes* (Chimpanzee), has developed various kinds of unstable, generally small male groups within a larger, open community of conspecifics (Mitani et al., 2002 and references therein), while females need not be part of any social group. Other pongids have worked out still different “solutions” to the social coherency problem arising from increased individuation of behaviors.

The diversity of solutions to increased individuation suggests that the latter has brought the pongids, our closest non-human primate relatives, up against an evolutionary “barrier” caused by increased individuation. The two primary mechanisms for social integration – familiarity of individuals with one another through face-to-face interaction and biological kin selection for social behaviors between genetically related individuals – apparently cannot cope with the degree of increased individuation and the consequent range of possible patterns for social interaction among wild chimpanzees. Though a chimpanzee community is made up of males with greater biological kin affinity within communities than between them (Morin et al., 1994; Vigilant, Hofreiter, Siedel, & Boesch, 2001), social behaviors are not determined through kin selection (Goldberg & Wrangham, 1997; Mitani, Merriwether, & Zhang, 2000). Instead, they are responsive to age and rank of chimpanzees (Mitani

et al., 2002), that is to small sub-groups in which social learning through face-to-face interaction can occur despite increased individuation of behavior.

We must conclude that face-to-face learning can be overwhelmed by the amount of interaction (and possibly by cognitive overload) needed to maintain social cohesion in the presence of highly individuated behavior, and that biological kin selection may run into limitations due to difficulties in identifying more distant kin, or because biological kin may not be available or suitable (Mitani et al., 2002).

2.2.1.3 Social Organization Based on Enculturation

We next consider the much more complex case of social organization based on a process by which the ideational aspect of the phenotype of an individual (which we refer to as cultural knowledge) develops through what cultural anthropologists have called *enculturation*. As noted by the anthropologist Conrad Kottak (2004: 209):

Enculturation is the process where the culture that is currently established teaches an individual the accepted norms and values of the culture or society in which the individual lives. The individual can become an accepted member and fulfill the needed functions and roles of the group. Most importantly the individual knows and establishes a context of boundaries and accepted behavior that dictates what is acceptable and not acceptable within the framework of that society. It teaches the individual their role within society as well as what is accepted behavior within that society and lifestyle.

Enculturation is the cultural analog to the transfer of genetic information from parent to offspring. Though we may focus on transfer of a single genetic trait for analytical purposes, humans are endowed with a genome made up of 23 pairs of chromosomes, and the ensemble of chromosomes that constitutes our genome is transferred via DNA duplication. Our genome contains both individual genes and organizational information governing the development of an organism through epistatic and other effects among genes. In a similar manner, cultural knowledge transferred from cultural parents and other culture bearers to offspring through enculturation is not just transfer of specific cultural information such as a particular norm or value, but transfer of the complete conceptual framework through which behavior is produced and interpreted by individuals. Just as genetic behaviors are not transmitted directly through sexual reproduction but indirectly via the genetic basis for phenotypic traits, cultural behaviors are not transmitted directly through enculturation. What is transferred is the ideational basis through which culturally based behaviors are constructed. And just as novel behaviors may arise out of individual experience interacting with the cognitive capacities that unfold during phenotypic development, novel behaviors may also arise out of an individual's evaluation of one's cultural development that unfolds during enculturation.

Enculturation begins at birth through interaction between the already enculturated mother and the newborn child and between the newborn child and other enculturated people in its environment (see Individual D in Fig. 2.4).⁴ Such interaction is

⁴ Enculturation is not mechanical information transfer, but a complex process of interaction between less encultured and more encultured individuals (Vinden, 2004), and is analogous to the

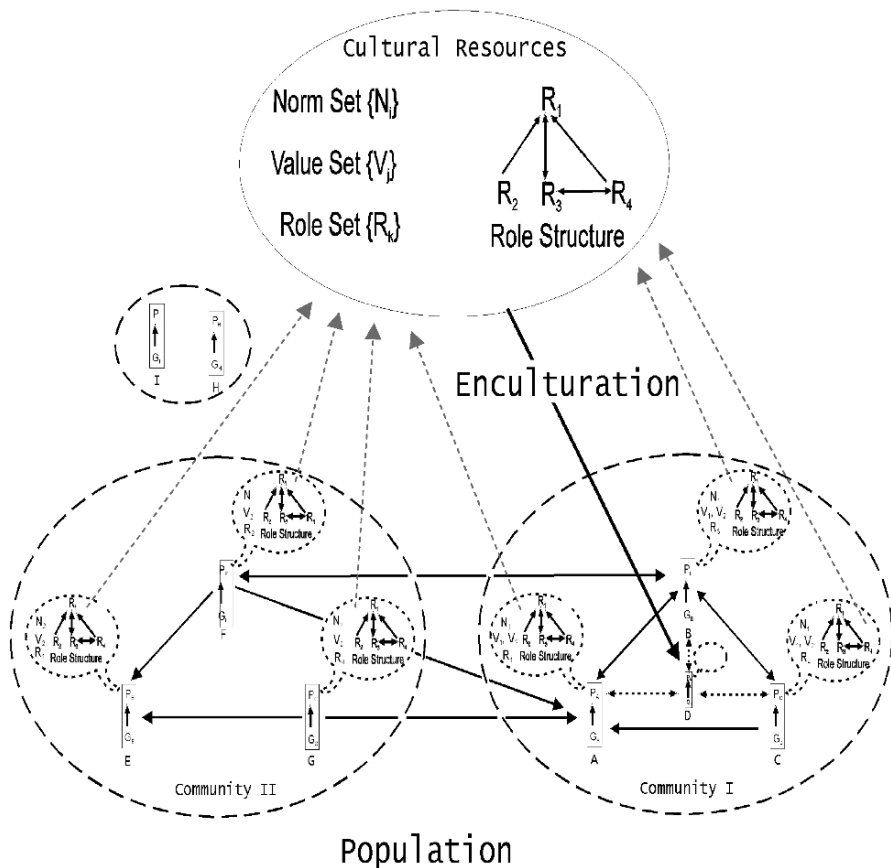


Fig. 2.4 Population of nine individuals (A–I), subdivided into two communities, each determined by shared norms and values but differing between the two communities. Oval at top of figure includes cultural components summed over all individuals (*gray dashed arrows*). Two individuals, H and I, do not share cultural components with the members of the two communities. Individuals A–C and E–G interact according to the one’s role/identity and in accordance with the role structure shared by these individuals (independent of community). Individual D is being enculturated through interaction with A, B and C (dashed arrows) who are already enculturated into the same cultural framework and are members of the same community. The actual pattern of interaction among the individuals need not have a network structure identical to the conceptualized role structure

process of language acquisition. Just as language acquisition continues throughout one’s life, enculturation continues throughout one’s life. Just as languages can be modified (though in structurally constrained ways) during the process of language acquisition, cultural information systems can be modified in structurally constrained ways during the process of enculturation. Just as language acquisition involves the learning of complex semantic and syntactic systems with multiple levels, enculturation involves the learning of complex cultural knowledge systems with multiple levels. Just as language acquisition is error prone and involves means for error correction, so does enculturation.

crucial as it provides the child with “all of the cultural resources that inform both its cognitive processes and the events to which they are applied” ((Schwartz, 1981: 14), shown at the top of Fig. 2.4) and does so by involving the developing child in

...information-rich, culturally structured events. . . Those events are structured by other enculturated persons and by the child as participant. The child learns by acting, acting upon, but also by being acted upon and by acting in pre-structured and other-structured scenes and events . . . The child is immersed in richly structured events upon the natural structure of which, a ‘second nature,’ cultural form, is superimposed (1981: 14–15).

Through enculturation, a child learns the conceptual basis underlying behavioral interaction, allowing it to understand the potential implications of one person’s behavior for another person. This is achieved by transmitting, at one level, the behavior in isolation and, at another level, the meaning given to that behavior as part of a culturally constituted conceptual system. Just as language can convey meanings expressed through linguistically constructed speech acts, behavior can convey meaning expressed through culturally constructed patterns of behavior. When we interact through a role we take on, we are not merely engaging in a behavior, but we are communicating to an audience (with whom we are interacting) the identity and meaning of the role we have taken on. Thereby we are also communicating information about the kinds of behavior that we are likely to engage in, and the likely kinds of behavior we expect in return. For this communication to be meaningful to both sender and receiver, all must participate in the same cultural framework (see Fig. 2.4, Community I and Community II).

Categorization of Behavior and Conceptual Systems

We can illustrate the implications of enculturation for social organization with an example based on categorization of individually formulated notions of “friend-like” and “enemy-like” behavior. In one form or another, categorization is wide-spread because it is a basic means to differentiate behavior according to the kind of entity with which an organism needs to interact, and thus increases the average utility of the interactions of an organism by enabling it to direct appropriate behavior towards the entities it categorizes. Hence, we may assume that some form of categorization is in place even in the absence of cultural categorization resources.

How would this work in practice? Consider three interacting individuals, A, B and C, who are not currently drawing upon cultural resources, and two kinds of behavior: friend-like behavior (including cooperation) and enemy-like behavior (including non-cooperation). At the level of dyads, the individuals may work out interaction patterns based on prior experience: A, through prior interaction with B that had negative consequences, may be induced to exhibit “enemy-like” behavior towards B in the future. Similarly B, through prior friendly interaction with C, may be led to exhibit “friend-like” behavior towards C. The two dyads together represent a simple example of a partial social structure worked out through learned behavior. It is partial because A and C have not yet learned how to interact with each other.

Cultural Computational System and Behavior

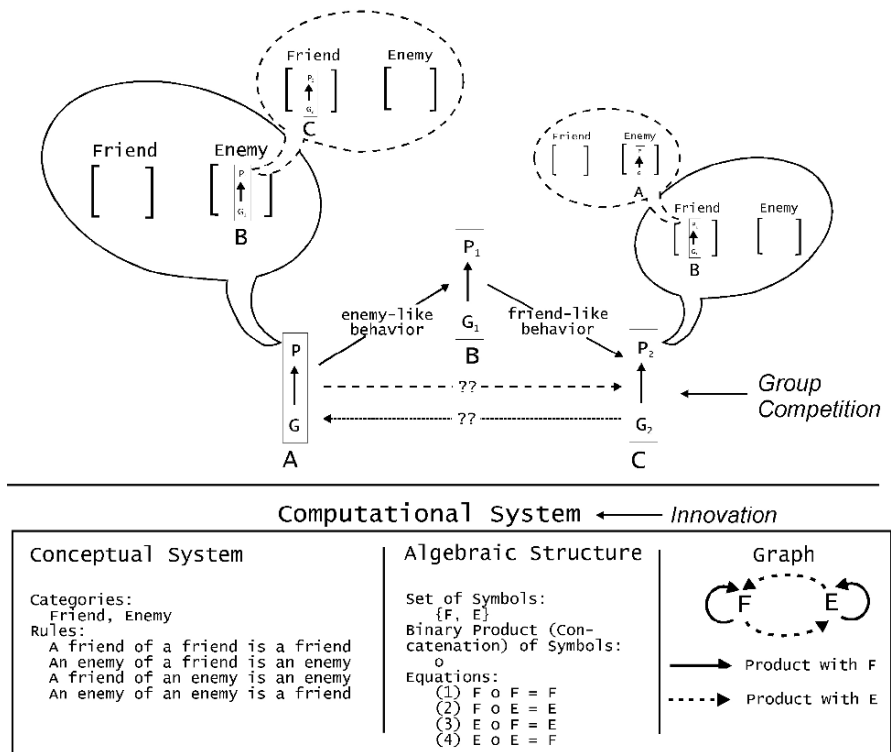


Fig. 2.5 Cultural computation system constructed through cultural rules. *Top of Figure:* First order categorizations (solid “thought clouds”) by A and C of individual B based on their respective interactions with B. Second order categorizations (dashed “thought clouds”) by A and C based on observations of the other individual’s interactions with B. *Bottom of Figure:* Cultural rules linking categories are in the left box. Rules permit computation of culturally proper behavior based on categorizations of individuals. Computations by A and by C will lead to consistent expected and actual behavior when A and C use the same conceptual system through enculturation into the same community. Middle box shows how the conceptual system can be modeled as an algebraic structure. Right box graphically represents the algebraic structure. Innovation can occur at the level of structure and/or in behavioral instantiation of the conceptual system. Competition arises at the group level through group benefit arising through functionality emerging from group organization based on the conceptual system

Now assume that individuals not only learn behaviors, but categorize other individuals on the basis of behavior patterns. Assume A distinguishes two categories, “friend” and “enemy”, and has categorized B as an “enemy” since A interacts “enemy-like” with B (see the “thought cloud” outlined with a solid line pointing to A in Fig. 2.5). If C then categorizes B as a “friend” on the basis of B’s behavior towards C, we again have two dyads constituting a partial social structure, but this time based on categorization. For the two dyads to become a full triad, A and C need

to be able to make reasonably accurate predictions about each other's likely behavior in response to their own actions. Learning to do so through interaction is, of course, not problematic so long as the repertoire of behaviors and number of individuals in the group is not too large, as discussed above.⁵ But when that repertoire has become too large for successful exploration via dyadic interactions, it is more likely that the wrong behavior is adopted. How does A, then, decide on a way to act towards C that is consistent with the way C will likely act towards A but *without* prior interaction?

One solution involves *constructing a computation system* based on the two categories "friend" and "enemy" under the assumption that individuals have both a concept of "self" and a "theory of mind." The former implies that A is "consciously aware" of her/his own categories. The latter means that A believes that other individuals have the same categories ("friend" and "enemy") as she/he does (a second order belief). This is indicated in Fig. 2.5 by the "thought cloud" with a dashed border pointing to individual B, who is categorized by individual A as an "enemy". Further, A, through observing the behavior of B towards C (or by other means), believes that B has categorized C as a "friend". We assume that the same occurs with C (as indicated by the "thought cloud" with a dashed border pointing to individual B as categorized by C as a "friend").

Up to this point we have simply extended the repertoire of beliefs that individuals have about others either on the basis of experience with the other individual (first order categorization) or on the basis of projecting onto the categorized individual *the result of* first order categorization (second order categorization). Extending the repertoire of beliefs does not, in and of itself, lead necessarily to any specific behaviors. That A believes B categorizes C in B's "friend" category does not indicate what behavior A should exhibit towards C. Indeed, the meaning of a category such as "friend", and the behavior to be derived from it, depend upon the definition of the category, i.e. upon the "meaning" of the behaviors involved to the person doing the categorization. In addition, categorization of this kind does not require any coordination between individuals and one person's criteria need not match the criteria of another, even if both conceptualize a category labeled "friend." Finally, there is nothing emanating from the process of categorization that necessarily entails two categories, "friend" and "enemy". One category could be dropped or never defined without affecting the other category.

Cultural Rule System and an Algebraic Model

Now let us expand the two isolated categories, "friend" and "enemy", to a system of categories by using four rules that conceptually link them and thereby form a structure. The four rules are shown in the lower left part of Fig. 2.5. They were chosen because in some societies these four rules are part of the conceptual system linking categories such as "friend" and "enemy", and determine the semantic meaning of the categories when they are used to guide behavior. These four rules

⁵ For completeness, we also assume that the process of categorization and predictions of behavior are dynamic in that they are subject to updating through future interaction experience, but this part of the argument will not be explored here.

determine an algebraic structure – namely a set of symbols, a binary product defined over the set of symbols, and a set of equations that indicate when a symbol product may be simplified to a shorter symbol product, or even to a single symbol. The algebra is formed from the correspondences “friend” \leftrightarrow F, “enemy” \leftrightarrow E, “of a” \leftrightarrow o (the binary operation o) and Rule $i \leftrightarrow$ Equation i ($1 \leq i \leq 4$), as indicated in the bottom middle of Fig. 2.5; the algebraic structure has the symbol set {E, F}. The binary operation may be defined as the concatenation operation for this symbol set, and the four equations indicate when a pair of concatenated symbols may be reduced to a single symbol. The algebraic structure may be graphed by using an arrow to indicate the symbol (F or E) that is produced when a product is made of a symbol with either E or F (lower right side of Fig. 2.5).

If the system of four rules is part of the resources with which individual A has been enculturated, it allows A to make a somewhat more complex computation. Notably, if A categorizes B as “friend”, and believes that B categorizes C as “enemy”, then A calculates C as “friend of enemy,” which reduces to “enemy” via the third rule. A should then categorize C as “enemy”, and exhibit enemy-like behavior towards C. However, thus far the conceptual system only generates a behavior that A should exhibit towards C if A is consistent with A’s conceptual system, but A still has no way of knowing whether or not C is likely to exhibit enemy-like or friend-like behavior towards A. But if both A and C have been enculturated with the “friend/enemy” conceptual system, C will make the computation that A is an enemy of B and B is a friend of C, hence C should direct enemy-like behavior towards A. The result is that A will exhibit enemy-like behavior towards C, and C will (independently) exhibit enemy-like behavior towards A. Thus if A and C are both enculturated with the same conceptual system, A’s expectation about C’s behavior will be accurate and vice-versa.

Observe that the computation system is “useful” for any individual only if all other individuals in the group share the same computation system. Among those individuals sharing the same computation system, and to the extent that behavior is made in accordance with the computations, a consistent and predictable pattern of behaviors will emerge.

The conceptual system will also have implications for the social organization of a group. If we interpret “friend-like” as cooperative behavior, the group will partition into two subgroups where all individuals in a subgroup cooperate with all other individuals in that subgroup and if we interpret “enemy-like” as non-cooperative behavior, individuals in one subgroup will be non-cooperative with individuals in the other subgroup. Under conditions such as those in this example, therefore, a possible evolutionary outcome is that, even though all individuals share the same conceptual system, a single group fissions into two non-cooperative groups as the conceptual system becomes part of the enculturation of individuals.

The Community Boundary Problem and a Computational Basis for its Resolution

How do community members identify other community members except through prior interaction? If a group consists of individuals who share the same conceptual system and individuals who do not (individuals H and I in Fig. 2.4), then, without

prior interaction, a member of a community would not know how to identify the other members (or the non-members). How was this so-called “group boundary problem” resolved in hominin evolution?

Apparently, a special, universal computation system evolved that provides the basis for computing community boundaries of similarly enculturated individuals. That system is built out of the semantic terms we use to define, and to refer to, kin. The kin in question need not be biological kin, because the domain of kin in human societies is culturally, and not biologically, constructed (Read, 1984, 2001, 2005).⁶

Cultural kin are determined through a computation system in the form of a kinship terminology based on genealogical instantiation of kin terms. The system makes it possible to compute from the perspective of one individual, A, whether another individual, B, is among A’s cultural kin, and reciprocally that A is among B’s cultural kin, when both individuals share the same kinship terminological system. But if A and B have enculturated the same kinship terminology, then one may assume that A and B share all other cultural resources that are transmitted through enculturation. Hence an effective behavior strategy becomes: “First determine if an individual is within your kinship domain. If so, assume that person shares your cultural resources and act accordingly. Do not interact with persons who are not within your kinship domain.”

This is precisely the strategy for behavior that occurs, for example, in hunter-gatherer (and other kin-based) societies such as the !Kung san who live in the Kalahari Desert in Botswana. For the !Kung san the word for stranger (*dole*) is also the word used for something that is harmful or dangerous and someone who is not a kin is a stranger, thus bounding social interaction to one’s kin. One’s kin are determined through their kinship system that enables individuals to compute whether they are kin.

The computation system is similar to that of the “friend-enemy” example in that it is also made up of a set of symbols (the kin terms) and a binary product defined over the kin terms. The products for pairs of kin terms can be elicited from users of the terminology and are based on their kin term usage. For example if, in the American kinship terminology, I refer to someone as “uncle” and that person refers to someone as “son”, then I (properly) refer to the last person as “cousin” and so we have as a product for the pair of kin terms “son” and “uncle”: “son of uncle is cousin”.⁷

⁶ This does not mean that cultural kin and biological kin do not overlap as the conceptual basis for cultural kinship ultimately derives from biological reproduction, but cultural kin are constructed through an abstract computation system that removes any causal linkage between biological reproduction and the construction of cultural kin (Read, 2001).

⁷ Product definitions are specific to a terminology as different societies may have non-comparable terminologies; that is, a term in one terminology may not be equivalent to any term in another terminology. Not all kin term products yield a kin term; e.g. for the American kinship terminology Father of Father-in-law is not defined as there is no kin term for the person one’s father-in-law refers to as father.

We have now identified the two basic elements through which human societies are able to circumvent what appear to be severe limitations on group size, while maintaining social cohesion in the face of increased individuation when social organization is based on a combination of genetically based and individually learned behaviors. These elements are (1) a mode of transmission at the ideational level, namely enculturation, that is comparable to the process of sexual transmission at the bio-chemical level of a genetic system, and (2) a means other than prior interaction to identify individuals with whom behavior based on an enculturated ideational system will be appropriate, namely computation of individuals who are mutually kin to each other.⁸ Since the computation system is transmitted through enculturation, reciprocity in kin identification is equivalent to identification of individuals enculturated with the same complex of cultural resources. In addition, reproduction among cultural kin involves cultural restrictions on the conditions under which procreation may occur. Kin computation has consequences similar to restricting sexual reproduction to conspecifics. Both provide a means to sufficiently bound variation in the pool of transmissible traits so that functionality is not lost in the transmission.

Enculturation and Selection: Two Modalities for Selection

Selection provides the balance between what can be transmitted and the implications the latter may have for functionality. In reproductive transmission, a random mutation – which may be deleterious to the functioning of an organism – can be transmitted and selection is then virtually an automatic by-product of the consequences that mutation has on the reproductive success of the receiving organism. In transmission by imitation, such transmission is subject to evaluation of novel learned behaviors, thereby introducing consequences for the phenotype that would not arise under reproductive transmission. Nonetheless this process still is a form of selection activated by individual functionality.

Once enculturation transmission is in place, functionality shifts from individual functionality to group functionality because what is being transmitted addresses group and not individual properties. In addition since the cultural resources are themselves constructions of a past social group, and were carried forward through enculturation, they can be modified by the current social group on the basis of its understanding of its current versus its desired functionality for the individuals involved. Hence, they are subject to change and modification in a manner that is poorly described by a combination of population-based selection and reproductive and imitation-based transmission.

Selection may not be the proper term for this kind of modification given its technical use in Darwinian evolution for processes that lead from individual to aggregate level change in patterning (such as reproductive success). Selection at the cultural resource level has to do with change in a specific cultural resource (such as a kinship

⁸ Appropriate behavior and actual behavior need not coincide for a variety of reasons and discordant behavior can be a signal carrying information about the trustworthiness, or reliability, of one's kin (Biersack, 1982).

terminology, a role structure, a norm or a value), and with its implications for the functionality of a social system of individuals enculturated under that changed set of cultural resources. Yet, this is selection in the non-technical sense and so the term selection will be used here for the evaluations that lead to change in cultural resources. Context should make it clear which meaning of selection is being used.

Selection in the non-Darwinian sense can occur in two ways. Firstly, as *selection that leads to a conceptual system* – which may range from a system enabling computation of categories of individuals (e.g., a system with an algebraic structure), to a system where the roles and the role structure are learned rather than computed, and where individuals may be marked according to the roles they take on (e.g., by special clothing or “scripted” behavior). Such selection will favor cultural resources that can be transmitted in a manner that allows errors to be corrected (e.g., in the form of algebraic and other formal structures), since maintenance of functionality depends on faithful transmission. Functionality obtained from a cultural resource *does not accrue to the group without the enculturation of group members* with that cultural resource. One such functionality that is especially relevant to early hominin evolution is the capacity to facilitate social cohesion in the face of a level of behavioral heterogeneity that would limit socially coherent groups to a size/density suboptimal for the exploitation of their resource base. Another functionality is the use of hunting techniques only available to larger, socially coherent groups. In either case, *the functionality is group functionality, and not individual functionality*, and individuals benefit from the functionality by being a member of the social group.

Secondly, when a shared conceptual system is in place, such as when one group is in competition with another group for access to resources, *group competition at the level of the social group is itself a form of selection*. Such *group competition* depends on the functionality of the social group, and thus on the effectiveness of the shared conceptual system. Hence a group and its organizational structure are in competition with another group and its organizational structure.

Enculturation and Group Competition

If several similar societies (e.g., hunter-gatherer societies) are in competitive equilibrium, under what conditions will a change in organizational structure in one society lead to a new configuration in another? The Lotka-Volterra model of group competition implies (see Read (1987) for details) that two groups in competition will have an equilibrium attractor state when the feedback from the growth of group A has a greater effect on reducing future growth of group A than on growth in group B, and vice versa (see Fig. 2.6, intersecting lines). For example, two hunter-gatherer groups can be in equilibrium even when regions from which they obtain resources partially overlap.

Unlike reproduction transmission, where an increase in individual fitness of any trait in a fitness equilibrium leads to replacement of another traits by the fitter one, group competition involves just a shift in the position of the equilibrium point between the two groups when organizational change in one group makes it slightly more competitive (e.g. by resulting in a slightly higher carrying capacity).

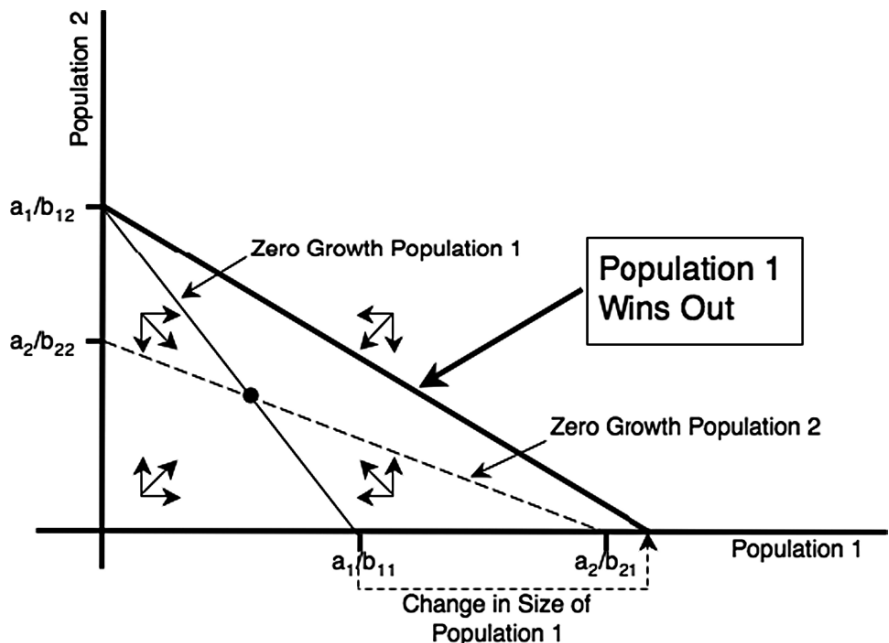


Fig. 2.6 Phase space graph for two populations in competition modeled via the dynamic model $dP_1/dt = a_1P_1(1 - b_{11}P_1 - b_{12}P_2)$ and $dP_2/dt = a_2P_2(1 - b_{22}P_2 - b_{21}P_1)$. Dashed line and light solid line show when P_2 or P_1 , respectively, have zero growth. Intersection of these two lines determines a stable equilibrium between the two populations (solid dot). Only when, say, Population 1 makes a qualitative change in, say, its carrying capacity does the configuration now shift to one in which Population 1 (heavy solid line) wins out in competition with Population 2

Replacement will only occur if one group is substantially more competitive than the other, but not otherwise (see Fig. 2.6, heavy solid line). In effect, a threshold value has to be crossed that is comparable in magnitude to the value for the measure in question. This is precisely the pattern we see in organizational structures for human societies. Hunter-gatherer societies generally consist of up to about 500 people. When replaced by, say, a society with a tribal form of organization, the latter will have a population size substantially larger – often an order of magnitude larger (e.g., 5,000 or more people). Hence, group competition acting on organizational change will give rise to qualitative shifts in organizational structure. This implies that we should see a “step sequence” rather than smooth transitions for organizational structures.

Figure 2.7 summarizes the argument we have developed in this section. The solid vertical line represents the constraint that separates non-human primate social organization from the social organization that developed in the hominins by the time *Homo sapiens* appeared. On the pongid (left) side of the diagram, the combination of increased individuation and social organization based on face-to-face interaction leads to a decrease in the size of social units and the size of social units may have

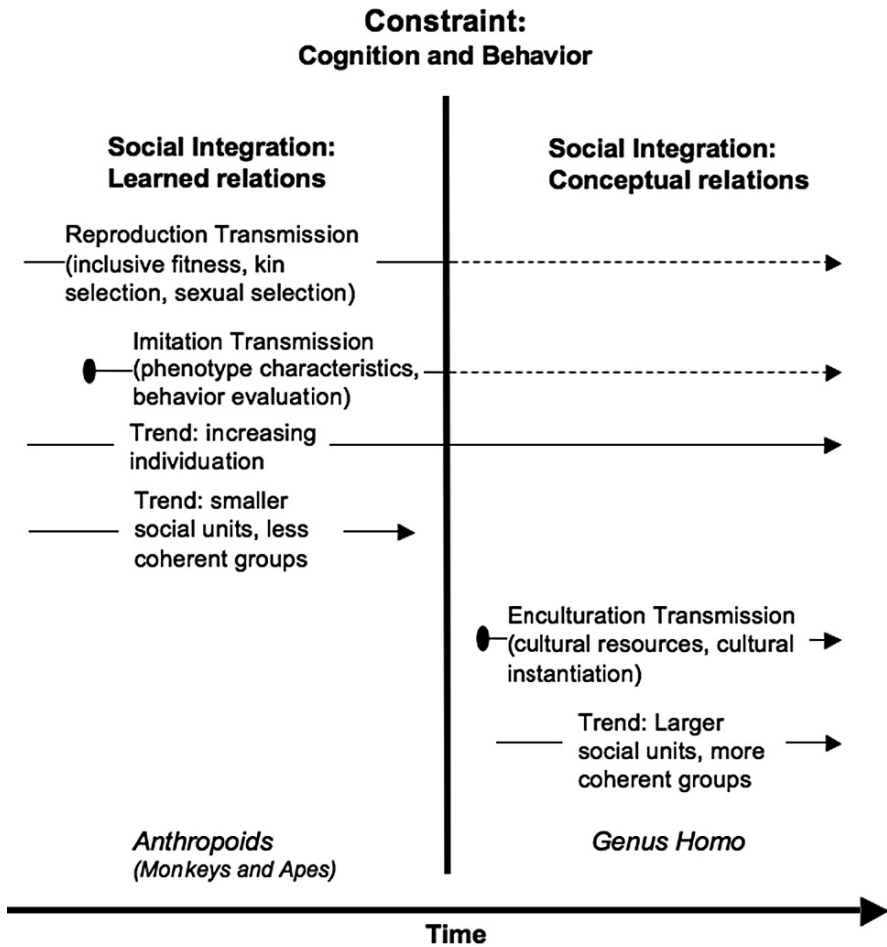


Fig. 2.7 *Vertical line*: Evolutionary barrier for social integration based on individual (reproductive and imitation) transmission due to increased individuation and cognitive limitations. *Solid ovals*: New functionality introduced through Darwinian evolution of cognitive capacities

been sub-optimal for efficient exploitation of resources in the (spatially and temporally) heterogeneous east African open woodland/savannah environment in which they developed. On the hominin (right) side of the line, we posit the advent of a computation system of conceptually formulated, dyadic roles – the basis for cultural kinship systems – that link one individual with another through reciprocal behavior without requiring lengthy prior face-to-face interaction between these individuals. Hence, in so far as expected patterns of behavior are associated with these relations, social integration and group coherence no longer depend on extensive interaction (Rodseth, Wrangham, Harrigan, & Smuts, 1991).

The key conceptual abilities that were needed to make this shift (to be discussed in the next section) do not appear to be present in the non-human primates. Hence, the transition depended upon the introduction of new conceptual abilities. It also shifted the basis for social integration from aggregate level change in a population of interacting individuals, to change in an organizational structure(s) in which individuals are embedded. This shift away from individual fitness-based selection had its precursor in direct transmission of behavioral phenotypes through imitative behavior and individual learning (as indicated in Fig. 2.7). But neither imitation transmission nor reproductive transmission, nor a combination of these, suffices to account for the forms of social organization (and the cultural systems) that eventually arose with the hominins and now characterizes our species.

2.3 Decoupling Social Systems from a Genetic Basis: A Pathway from Darwinian to Non-Darwinian Evolution

We assume that the trend towards increasing individuation of behavior, and the negative impact it had on the coherence of social structure, set the conditions for selection in favor of a new mode of individual interaction. But we are not arguing that increased social complexity, which appears to be part of the evolutionary pathway leading to the pongids and the beginning of hominin evolution about 8 million years go, was due *only* to the trend towards increasing individuation. There were other changes, such as a shift to a frugivorous diet (reducing the range of vegetal resources and the spatial and temporal predictability of resource distribution), which may have led to behavioral changes that added to social complexity. Hunting for meat, a socially complex activity, also becomes part of the behavioral repertoire of *Pan troglodytes*. We are merely using the trend towards individuation to highlight the likelihood of initial Darwinian selection for the cognitive ability to engage in some form of social interaction not dependent on extensive, prior interaction among individuals.

A “bottleneck” limited the domain of possible solutions, so that a solution to the problem of integrating individuation with social cohesion did not arise even with the cognitive capacity of the chimpanzees, despite the 8 million years that elapsed between our common primate ancestors and the development of modern chimpanzees. It is only with hominin evolution that these inherent limitations on the size of a coherent social group were circumvented, by means other than elaboration of face-to-face interaction and imitation of behaviors.

For social interaction to be systematic, it must be reciprocal, ongoing and not just episodic. As a consequence, biological means for facilitating social interaction depend either on some form of biological kin selection or on a way to identify individuals predisposed to engage in reciprocal and positive social interaction. Behaviors repeatedly directed towards either non-biological kin or non-reciprocating individuals favor selection for non-social interaction (e.g. “cheating”) on the part of the recipient, since the latter benefits without engaging in positive, reciprocal behavior.

A limit to the size of primate social groups integrated through biologically based (rather than learned) social interaction arises from the relatively few means available to primates for identifying biological kin. Such means may arise indirectly from interaction with biological kin, primarily between female genetic parents and offspring, or between biological siblings raised together. But, even if more distant biological kin can be “identified” through patterns of behavior (as in the case of sexually maturing females remaining in their natal troop), the effectiveness of such biological kin selection in biasing behavior in favor of reciprocal social interaction decreases exponentially with genetic distance. Identification of non-biologically related individuals predisposed to engage in reciprocal social interaction is even more problematic.⁹

While we cannot yet identify the precise conditions laying the cognitive foundations for a new mode of social interaction that eventually became decoupled in its implementation from a biological substrate, the existence of such a new mode is evident from the fact that our species, *Homo sapiens*, has found the means to accommodate both increasingly individuated behavior and larger social units in a coherent and effective manner that addresses the collective interests of the group. The magnitude of this shift in the basis for social integration is evident in three major differences between our non-human primate relatives and ourselves.

First, even in the smallest and simplest of modern human societies, namely hunter-gatherer societies, the number of individuals integrated together is between one and two orders of magnitude larger than the size of social units found among the pongids. For example, the !Kung san hunter-gatherers mentioned above live in groups of about 30 individuals (Lee, 1990) integrated together as a single society of about 500 individuals, whereas pongids such as the chimpanzees have unstable social units with around 6–20 males (Nishida & Hiraiwa-Hasegawa, 1987). The individuals in a single hunter-gatherer society are divided into residence groups that may be spatially isolated yet allow for frequent, non-disruptive movement of individuals from one group to another during the life cycle of individuals. In contrast, non-human primates are typically organized into small social units between which there are often antagonistic relations (due to territoriality) and for which change in social unit residence of adults (especially males) is usually highly disruptive and infrequent.

Secondly, hunter-gatherer societies have a pattern of food sharing that need not be structured around biological kin relations and may involve different residence groups (see Table 2.1), whereas food sharing among non-human primates is not common and when it does occur, almost exclusively concerns vegetable food exchange between females and their offspring or occasionally from a male to a female (see Table 2.2, entries in bold), apparently as a way to gain access to a female for reproduction (McGrew, 1992).¹⁰

⁹ The difficulty of circumventing this limitation can be seen in the fact that the non-human primates have not yet found a way to resolve this biological limitation even after more than 8 million years of Darwinian evolution.

¹⁰ Patterns of food sharing among chimpanzees are more complex when meat sharing is taken into consideration and extensive variation occurs among hunter-gatherer groups with regard to when

Third, in contrast to a bounded primate troop based on face-to-face interaction, the social boundary of a hunter-gatherer society is defined by the set of individuals who can mutually determine that they are cultural kin to one another, i.e. by identifying shared kin relations through a kinship terminology. Kinship terminologies are cultural constructs (Parkin, 1997; Read, 2001); hence the people we identify as our relatives are culturally specified and can include non-genetic kin, who only bear an indirect relationship to biological kin relations. Kinship terminologies differ from one society to another in a manner analogous to differences between languages, making conceptual distinctions in one terminology that need not be matched in another kinship terminology. Hence, the way in which one society is socially structured via cultural kin relations need not have its counterpart in another society.¹¹

2.3.1 Four Cognitive Capacities

The evolutionary pathway of our hominin ancestors necessarily starts with cognitive changes introduced or elaborated through Darwinian selection. The decoupling arises because subsequent consequences of those cognitive changes made it possible to construct social relations between individuals independent of biological kin relationships. The cognitive changes are four-fold: (1) a concept of ‘self’, (2) a “theory of mind”, (3) categorization based on the concept of a relation between individuals and (4) recursive composition of relational categories. We first briefly indicate what is meant of each of these cognitive properties and the extent to which their precursors can be found among the non-human primates. Then we consider in more detail how (3) and (4) made it possible to construct conceptual relations between individuals that can be organized into a computation system that serves among other things as a means for identifying a group of individuals that form a community through enculturation.

(1) The “concept of self” implies cognitive awareness of one’s own existence, or identity, in contrast to the existence of others. It entails that seeing an image of oneself is cognized as a representation of oneself, and not of a conspecific. To test whether non-human primates have a concept of self, researchers have placed a mark on a target individual and then registered whether the individual responds to the mark upon seeing her/his image in a mirror. If it does so by attempting to touch its location, it is assumed that the individual is linking properties seen in the image with those of his/her own body. Based on this criterion, some of the pongids, such as chimpanzees, have a concept of self, but their evolutionary precursors, the Old World monkeys, do not. As experimental evidence for a concept of self is substantial for the chimpanzees, we assume that a concept of self was already present in a primate ancestor common to chimpanzees and hominins.

food sharing occurs and the relationship between giver and receiver. These two examples should just be viewed as illustrating the qualitative difference in food sharing among primates versus human groups.

¹¹ Not all kinship terminologies are unique. Unrelated societies can have identical kinship terminologies.

(2) Having a notion of a “theory of mind” means not only that one has awareness of one’s own properties, whether they be physical, behavioral, or cognitive, but that one is able to conceptualize that other conspecifics may also have the same properties or mental representations. In particular, when an individual is aware, for example, of its own actions in response to external stimuli, then upon seeing another individual act in a similar manner under the same circumstances, the first individual can conceptualize that the other is doing so for a similar reason. Experimental work with human infants has established that a theory of mind is in place in humans by around 3–5 years of age (Hughes, 2004). But whether any of the non-human primates have a theory of mind is less clear. There is no evidence for a theory of mind among any of the Old World monkeys. Some have argued that experiments with chimpanzees show behavior patterns consistent with a theory of mind (Povinelli, Nelson, & Boysen, 1990; Premack & Woodruff, 1978; Woodruff & Premack, 1979), though others (e.g., Savage-Rumbaugh, Rumbaugh, & Boysen, 1978; Tomasello, 1998) have challenged that interpretation. Given this uncertainty, we assume that even if the common ancestor for the chimpanzees and ourselves did not have a theory of mind, the cognitive ability to do so most likely arose early during hominin evolution, so we will assume that a theory of mind is already in place among our hominin ancestors.

(3) The next cognitive capacity of concern is *categorization of relations between individuals*. Category formation with respect to *properties* of phenomena external to an individual is virtually ubiquitous among organisms, though the means by which it occurs may vary from cognitive to chemical. At some point, though, this capacity is extended to categorization based on *relations between pairs of individuals*. Recent work in primatology has begun to document the importance of such categorization as a causative factor in social behavior: “The individual model has therefore been replaced by a relational model. . .” (de Waal, 2000: 588). Correspondingly, categorization by non-human primates may possibly incorporate categories based on relations and not just properties of objects or individuals. However, the cognitive capacity to conceptualize a category of relation such as “mother” or “daughter” will arise only if the capacity is biologically grounded and a positive fitness benefit accrues from behaviors associated with the categorization. For this to be the case, the behaviors need to be directed towards one’s biological kin. Hence the category of relation that is being conceptualized also needs to be biologically accurate.

The extent to which categorization based on relations occurs among the non-human primates is unknown except for one experiment with long-tailed Macaques showing that they are capable of categorization based on biological mother/offspring and sibling relations (Dasser, 1988a, 1998b).¹² We recognize that the small number

¹² For the mother/offspring experiment the target macaque was primed with pictures of biological mother/daughter pairs. When presented with a choice between pictures of a novel biological mother/child pair and a biological mother/non-child from the same troop, the target macaque consistently (12 out of 12 trials) selected the biological mother/child pair. Similarly, when presented with a novel female and prompted to select between that female’s offspring or a non-offspring the primed macaque almost always selected the former (18 out of 20 trials). The experimenter writes:

of trials, the absence of other experiments that replicate these results, and the fact that only two out of many trained individuals responded in this manner (Cheney & Seyfarth, 1999) imply the evidence of the ability of macaques to form categories on the basis of the biological mother/offspring relation (rather than, for example, on the basis of individual characteristics of behavior between a female and her offspring), is suggestive only and not definitive. For our purposes, therefore, we will use the experimental results to suggest *only* that the ability to categorize on the basis of a relation, even if on a sporadic basis, may already be present in the early ancestry of the hominins.

The trend towards individuation may play a role in a shift towards categorization based on relations, as increased individuation decreases the likelihood that different individuals will have similar behaviors, hence making it more difficult to categorize on the basis of behavioral features. When different female/offspring pairs engage in the same behavior, for example, categorization could be based on the behavior and not the relation. With increased individuation and greater novelty in behavior on the part of a female/offspring dyad, categorization, if it occurs at all, should increasingly be based on the relation, as the relation may be the only constant factor across the instances of biological mother/offspring pairs. Hence, we posit that the trend towards increased individuation also increased the likelihood that categorization would take place on the basis of the relation between pairs of individuals. That this kind of categorization did arise eventually among our hominin ancestors is not in question, as it is the basis for the cultural kinship systems that arose sometime in our hominin ancestry; the only uncertainty is when and under what conditions the capacity arose.

(4) The last cognitive capacity, namely *recursive reasoning*, does not occur among the non-human primates (Hauser, Chomsky, & Fitch, 2002). Even the simpler cognitive task of learning a phase-structure grammar is beyond the capacity of non-human primates such as the tamarins (Fitch & Hauser, 2004). Experimental work on chimpanzees (Spinozzi, Natale, Langer, & Brakke, 1999) demonstrates clearly that regardless of their ability to work with, and attribute meaning to, symbolic representations, and regardless of the extent to which chimpanzees may have the cognitive capacities that are the precursors for language ability, they lack the cognitive capacity to reason in a fully recursive manner. The inability of chimpanzees to reason recursively may relate to the fact that their short term working memory capacity is too small for recursive reasoning (Read, 2006).

This recursive capacity relates to the ability to form a composition of relations and thereby to generate a new relation through recursion. To illustrate, let the two-place predicate $M(-, -)$ represent the biological mother/biological daughter relation defined over a set of biological individuals, S , so that, for all x, y in S , $M(x, y)$ is true when, and only when, y is the biological mother of x . We can define a new relation, MM , by the two-place predicate $MM(-, -)$, where $MM(x, z)$ is true if, and only if there is a y in S for which $M(x, y)$ and $M(y, z)$ are both true. The relation MM may

“Mother-offspring pairs were differentiated from any other pair . . . cues other than the relation between individuals do not plausibly account for the result” (Dasser, 1988b: 91).

be constructed recursively from the relation M as follows. Since there is a single y for which $M(x,y)$ is true, let $y = M(x, _)$, hence we can think of y as the outcome of applying the predicate $M(x, _)$ to the set S . Similarly, we can let $z = MM(x, _)$ when $MM(x,z)$ is true since there is a single z for which $MM(x, z)$ is true. Then $z = MM(x, _) = M(y, _) = M(M(x, _), _)$, hence MM can be reconstructed recursively by applying the M relation to the outcome of the M relation. (This form of recursion is the basis for genealogical tracing via the mother/offspring relation that we will discuss in the next section).

Even if non-human primates apparently are not capable of this type of conceptualization, the mental capacity to conceptualize recursively did arise at some point in hominin evolution. Since it is a biologically grounded capacity, there must have been a fitness benefit for it to arise in the first place. Undoubtedly the development of that capacity does not refer to any single event but rather to a series of cognitive changes. Though languages of modern *Homo sapiens* make extensive use of recursion, the capacity for recursive thinking may also have arisen as part of non-linguistic activities such as tool making (Read & van der Leeuw, 2008), and then exapted as part of the development of linguistic capability. Precursors for recursive thinking can be seen in hominin stone tool making (but not pongid stone tool use such as nut-cracking), where an action such as the technique of flake removal from a stone object is done repeatedly. Each subsequent flake removal repeats the same action on the object produced from the previous step.¹³ Regardless of the specific trajectory that led to the cognitive ability to engage in recursive thinking, such a pathway was followed during hominin evolution. We are here dealing with what may have occurred during hominin evolution after recursive thinking is already in place.

2.3.2 Theory of Mind and Recursion

We begin the pathway from biologically framed to non-biologically framed evolution by assuming that we have a set of individuals, S , with the four cognitive properties we have just discussed. For the sake of illustration, we will focus on a single relation, namely the M (“mother”) relation, but the argument applies to any relation that has become the basis for a categorization of dyads among the individuals in S . For the reasons discussed, we assume that the M relation is initially based on

¹³ Some of these changes may have been triggered by selection for the cognitive ability to make more effective and efficient stone tools. Tools made by our hominin ancestors have varied in conceptual complexity from one-dimensional (e.g., an edge, which typifies stone tools dating about 2 million years ago) to three-dimensional (e.g., blade making in the Upper Paleolithic) conceptual control over the process of stone working (Pigeot, 1991; Read & van der Leeuw, 2008; van der Leeuw, 2000). For example, the making of tools from Upper Paleolithic blade cores, Middle Paleolithic disc cores, and Middle and Lower Paleolithic Levallois cores and bifaces involves a recursive technology, but the earlier one-dimensional and two-dimensional technologies of the Olduvian (beginning c. 1.7 mya) were iterative but not recursive. This development from iterative to recursive technologies provides evidence for a time range within which the cognitive elaboration of recursive reasoning was introduced through selection for an increase in hominin working memory.

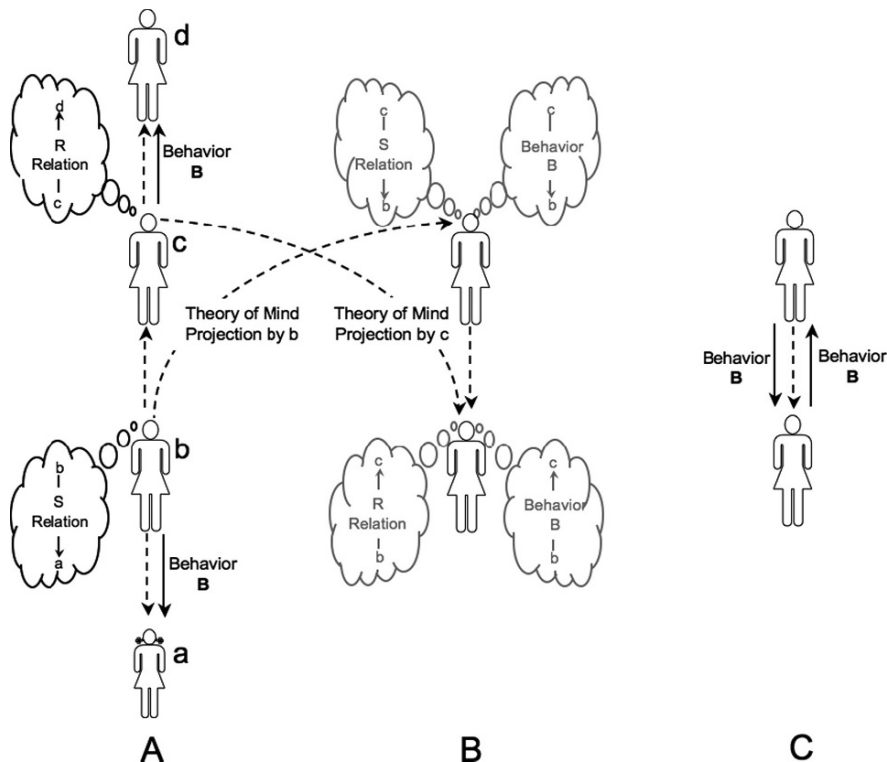


Fig. 2.8 (A) Individual *a*, biological daughter of *b*, conceptualizes a mother relation and (B) projects, via the Theory of Mind, the same relation concept to her biological mother *b*. (C) By composition of relations, individual *a* constructs a relation linking her to individual *c*, the female *a* believes to be the target of the mother relation she has attributed to *b*

categorization of actual biological mother/offspring relations. In Fig. 2.8(A), female *a* is the biological daughter of female *b*. We assume that as the daughter perceives that she and her biological mother are a dyad, she conceptualizes an instantiation of the *M* relation between herself and her mother. This is indicated by the “thought cloud” in Fig. 2.8. By virtue of the Theory of Mind, she believes that her mother also instantiates the same *M* relation between herself, *b*, and a female *c* perceived by *a* to be the biological mother of *b*. Thus the (*b*, *c*) dyad is believed by *a* to be an instantiation of the *M* relation perceived by her mother, *b* (see Fig. 2.8B). From the perspective of *a*, this instantiation is a belief since *a projects onto her mother her own belief* that her mother also perceives an *M* relation. The thought cloud in Fig. 2.8B for female *b* is dashed and in gray to indicate that this is the relation that *a* believes (correctly or not) is held by her mother.

By recursion, *a* can now construct the *MM* relation that instantiates her relationship with female *c* (see Fig. 2.8C). It differs in a crucial way from the *M* relation: it is constructed from the *M* relation through recursive reasoning and not

from categorization of actual “biological grandmother/biological granddaughter” dyads. Instead, categorization is a *consequence* of the new relation built on recursive reasoning, and categorization thus encompasses all those instances where, by virtue of the Theory of Mind, *a* projects onto another individual the relation *MM*. In other words, following the Theory of Mind, if individual *a* has constructed an *MM* relation instantiated by the (*a*, *c*) dyad, then individual *a* can project that *MM* relation onto other individuals and perceive other dyads as instantiations of the *MM* relation. Thus, the *MM* relation is not based on the biological relation of individual *c* to individual *a*, but on what individual *a* *believes* to be the case about the relation of *c* to *b*. That belief may be erroneous, but that does not affect the construction of the *MM* relation since there is no external reality against which the construction can be falsified. *Recursion of relations leads to decoupling of constructed relations from the biological basis for conceptualizing the relations involved in forming the constructed relations.*

2.3.2.1 Reciprocal Relations

From the perspective of the mother (*b*), *a* will be in a biological daughter relation *D* with respect to *b*. If *b* perceives both an *M* relation with *c* and a *D* relation with *a* (see Fig. 2.9A), and projects the *D* relation onto *c* (see Fig. 2.9B), then *b* will simultaneously be an instantiation of the projected *D* relation. Hence, *b* will

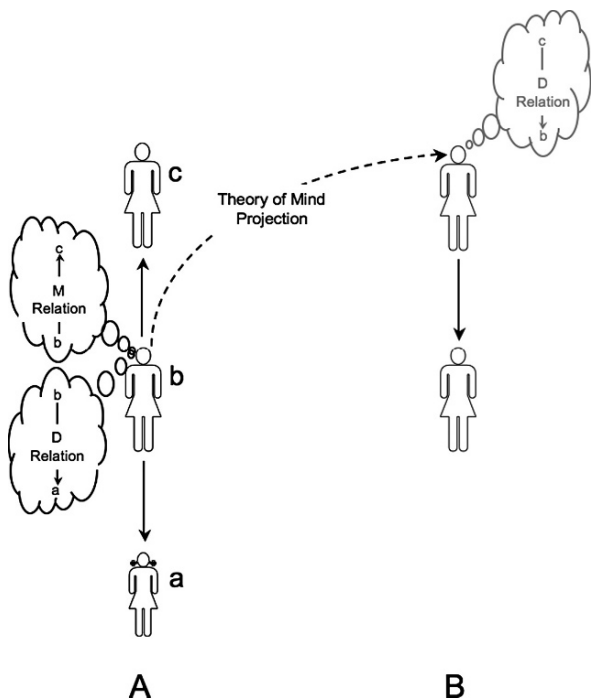


Fig. 2.9 (A) Individual *b* conceptualizes an *M* (mother) relation to *c* and a *D* (daughter) to *a*. (B) Individual *b* attributes the *D* relation to *c*, hence *b* believes that *c* has *b* as a target for the *D* relation, a precursor for a reciprocal social relationship from *b*’s perspective

perceive not only that b has an M relation to c , but also that c perceives a D relation from c to b . Consequently b will believe that b and c are conceptually linked to each other. Hence the precursor for a reciprocal social relation from b 's perspective, namely that b not only perceives a relation with c , but *also* believes that c perceives a reciprocal relation with b , is in place.

The same kind of pattern may arise with any relation R that b has with c and the corresponding reciprocal relation S that b may believe to have with a . The projection of the relation S onto c will have b as an instantiation of that relation from b 's perspective, and so b will perceive that b has a relation R with c and will believe that c perceives a relation S between c and b , regardless of what is the actual biological relation of c to b .¹⁴

2.3.3 *Functionality of the Projected Relation*

Just as with kin relations, where the evolutionary importance lies not in the biological kin relationship *per se*, but in a biologically based behavior associated with the biological relation, the importance of perceiving a relation R lies not in the relation *per se*, but in behaviors and/or motivation for behaviors that can be associated with the relation, and that lead to social interaction. Following Talcott Parsons (1964: 5), we distinguish here between interaction and *social* interaction:

It is a fundamental property of action thus defined that it does not consist only of ad hoc 'responses' to particular situational 'stimuli' but that the actor develops a *system* of 'expectations' relative to the various objects of the situation. These *may* be structured only relative to his own need-dispositions and the probabilities of gratification or deprivation contingent on the various alternatives of action that he may undertake. **But in the case of interactions with social objects a further dimension is added. Part of ego's expectation, in many cases the most crucial part, consists in the probable reaction of alter to ego's possible action, a reaction which comes to be anticipated in advance and thus to affect ego's own choices** (Italics in the original, bold added).

Parenting is an example of interaction that need not involve *social* interaction since a may engage in parenting behavior towards b by virtue of being a parent and without any necessary expectation of reciprocal behavior. Or, from the perspective of the child b , b has expectations about how a will act (feed, comfort, etc.) regardless of how b might or might not perceive his or her interaction with a . In general, behavior such as altruism introduced through selection based on biological kinship *is not social if there is no anticipation that the behavior will be reciprocated* in some manner. But behavior based on cultural kinship *is social* as the conceptual system that

¹⁴ The relation R could be "is a friend" or "is an ally" given the evidence that non-human primates modify behavior on the basis of the relationship between the two individuals in a dyad (reviewed in Silk, 2003). However, a constructed relation such as "is a friend of a friend" would presumably imply an expectation of friendly behavior on the part of the individual so identified, and the validity of the constructed relation for behavior may thus be subject to empirical verification – in contrast to a constructed relation such as "mother of a mother."

structures cultural kinship (the kinship terminology) is based on reciprocal relations and expected reciprocal behavior. If a recognizes b as a cultural kin then the kinship relation entails that a also recognizes that a is a cultural kin from b 's perspective by virtue of the reciprocal property of the kinship terminology. Therefore a can expect reciprocal behavior on the part of b .

To activate a social relation between a and b depends, then, on some understanding by a and b that they are at least conceptually linked to one another.¹⁵ Otherwise, there is no reason to expect reciprocal behavior. For kin-based societies (e.g. hunter-gatherer societies), there cannot be a social relation between individuals a and b unless they have first established that they are (cultural) kin – which means that they are both part of each other's conceptual domain, and have acknowledged it. An extreme example of this is the fact that, in the past, among the Waorani in South America, if a person b , came to a 's village and b did not have a kin relation with a , then the matter of determining whether a social relation is possible was resolved by a killing b (Davis & Yost, 2001).

Theory of Mind projections may trigger any kind of behavior that one individual might engage in vis-à-vis another if it is connected with a relation R between the two individuals. If b has a relation R with c , and a reciprocal relation S with a , and the biological relations among a , b and c are indeterminate, then b may exhibit behavior B towards any individual that is a target of the relation S conceptualized by b (see Fig. 2.10A). If b also projects the relation S and the associated behavior (B) onto c (see Fig. 2.10B), b believes that c will reciprocate with behavior B since she/he is a target of the relation S that b believes to be held by c . Then, b may engage in the same behavior towards c in the belief that c will reciprocate with that behavior (see Fig. 2.10C). We now have a basis for interaction to become social interaction: *one individual acts towards another individual under the belief that the other individual will act in a reciprocal manner. Further, and critically, this basis for social interaction is decoupled from any requirement of biological relations among the individuals in question.*

2.3.4 Mutual Recognition as a Basis for Social Interaction

While the projection of a behavior linked to a relation may lead to the belief that such behavior will be reciprocated, such reciprocal behavior need not actually occur unless the other individual has constructed a complementary belief system and behaves accordingly. Cheating, used here in the sense that the behavior is not initiated despite having the complementary belief system, is always possible and if b acts

¹⁵ A conceptual linkage between individuals is not a necessary pre-requisite for social interaction. Social interaction can arise from genetically based behavior if the selection for the behavior by a is correlated with possible responses by b . In probabilistic terms, non-social interaction would be a behavior B where the unconditional probability of a doing B in the presence of individual b is the same as the conditional probability of a doing B knowing that b does some behavior B' .

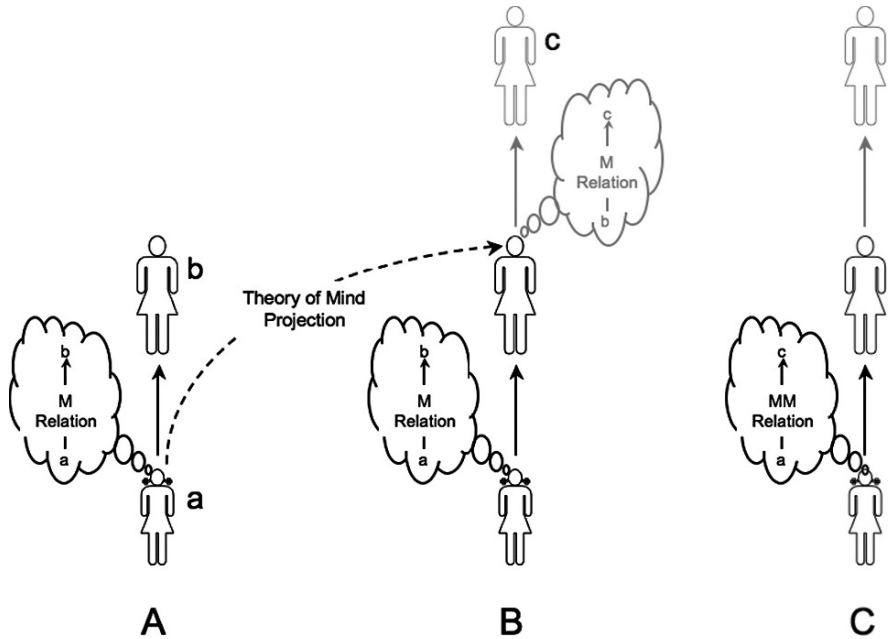


Fig. 2.10 (A) Individual *b* conceptualizes a relation *R* with *c* and a reciprocal relation *S* with *a*. In addition, *b* directs behavior *B* towards individual *a* when *a* is the target of the *S* relation conceptualized by *b*. (B) Individual *b* projects relation concept *S* to individual *c* and *b* is the target of the relation *S* believed by *b* to be a relation concept held by *c*. (C) Individual *b* directs behavior *B* towards *c* due to *b*'s belief that *b* is a target of the *S* relation held by *c*. That is, *b* believes *c* will direct behavior *B* towards *c* since *b* directs behavior *B* towards *a* due to *b*'s relation *S* with *a*, hence *b* expects *c* to direct behavior *B* towards *b*

towards *c* simply under the *belief* that *c* will reciprocate, then *b* has actually initiated conditions that favor cheating by *c*.

Actual, as opposed to potential, social interaction depends upon engaging in reciprocal behaviors. If both parties believe that the other will reciprocate, then the foundation for continued social interaction will have been laid. For this to occur, it suffices that *c* associates behavior *B* with the relation *R* and *b* associates the same behavior with the reciprocal relation *S*, as shown in Fig. 2.11.¹⁶ Under these conditions, both *b* and *c* will independently construct the belief the other will reciprocate with the behavior *B*. When each individual, based on his or her own beliefs, then engages in behavior *B* towards the other, the beliefs are reinforced by the actual behavior of the other individual.

This argument has two important implications. First, the pool of individuals who have the reciprocal belief system illustrated in Fig. 2.11 will increase if the behavior *B* is associated with both the relation *R* and its reciprocal relation *S*. When the

¹⁶ For the sake of clarity, the reciprocal relations have not been drawn for each of the individuals *b* and *c*.

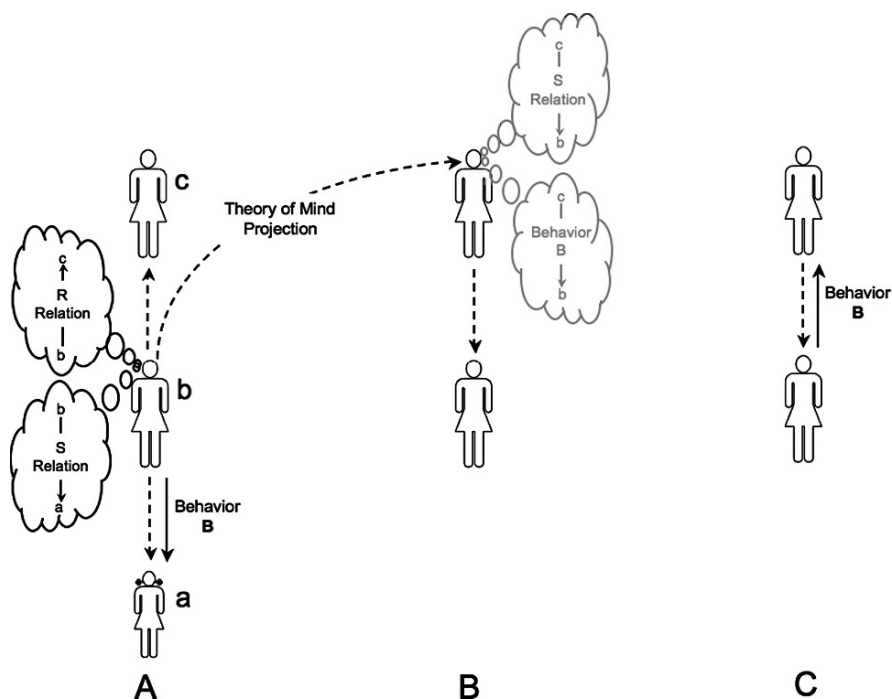


Fig. 2.11 (A) Individuals *b* and *c* each share the same conceptual pair of reciprocal relations (only one relation from each pair shown for clarity) and each associates] behavior *B* with a relation, with individual *c* directing behavior *B* to individual *d* and individual *b* directing behavior *B* to individual *a*. (B) Each of *b* and *c* projects their conceptual relations onto the other individual. (C) Each of *b* and *c* directs behavior *B* towards the other individual on the basis of one's belief that the other individual will reciprocate with *B* or *B*-like behavior. The beliefs of both *b* and *c* are reinforced by the behavior of the other individual

behavior is associated with both these relations for *b* in Fig. 2.11, then when *b* interacts with *c* it will be irrelevant whether *c* only associates the behavior with *R*, or only with *S*, or both. In any of these situations the conditions for arriving at Fig. 2.11C are satisfied. Translated into cultural kin relations, this means that a behavior will be associated with the kin relation *R* and the reciprocal kin relation *S* if whenever ego engages in a behavior *B* with respect to someone in the kin relation *R* to ego, then ego is equally willing to engage in that behavior with respect to someone in the reciprocal kin relation *S* to ego. For example, if ego as an adult engages in cooperative behavior towards his/her children, then ego should equally engage in cooperative behavior towards her/his parents since the parent-kin relation is the reciprocal for the child-kin relation.¹⁷

¹⁷ This differs from what would be predicted under biological kin selection with inclusive fitness for adults. Even under inclusive fitness there is little or no direct fitness benefit to be gained by directing even cooperative behavior towards parents instead of towards one's offspring.

Second, to realize the functional benefit of reciprocal behaviors, individuals must recognize the kind of relation with which a behavior is associated in comparable ways. Agreement between actor and recipient with respect to enactment of a behavior requires that both the actor and the recipient associate it with the set of reciprocal relations between them. The likelihood that that occurs depends on the degree of coordination among group members with regard to the relations they recognize and the behaviors associated with them. This, in turn, requires a role system for the patterns of behavior engaged in by individuals.

The coordination problem is thus solved through enculturation involving kinship terminologies by virtue of the fact that the system of kinship terms (1) is a computation system through which kin relations may be calculated in a simple manner, (2) is a generative computation system, and (3) implies that reciprocity for all kin relations follows from reciprocity of the generating kin relations.

By a computation system is meant that two individuals a and b can determine the kin relation they have to each other when there is a third individual, c , and a and b both know their relation to c via the kin terms they each use to refer to c :

... [Maori kin] terms permit *comparative strangers* to fix kinship rapidly without the necessity of elaborate genealogical reckoning — reckoning that typically would be impossible. With mutual relationship terms all that is required is the discovery of one common relative. Thus, if A is related to B as child to mother, *veitanani*, whereas C is related to B as *veitacini* (sibling of the same sex), then it follows that A is related to C as child to mother, although they never before met or knew it. *Kin terms are predictable. If two people are each related to a third, then they are related to each other* (Sahlins, 1962: 155, emphasis added).

The computation system is generative in the algebraic sense that there is a subset of the set of kin terms from which all other kin terms can be generated through use of (a) the binary product for kin terms and (b) a set of structural equations based on just the generating kin terms. Empirically, the kin terms serving as generators refer to family relations. For some terminologies, the terms that express the mother and father relations are the generators; for others, the terms expressing the brother and the sister relations are also generating terms. The latter terminologies have very different structural properties than those for which the terms expressing the brother and sister relations are not generators (Bennardo & Read, 2005; Read, 2001).

Reciprocity of kin terms follows from reciprocity of the generating kin term relations. It means simply that when, say, the mother relation is conceptualized, then the relation from her to her offspring is also conceptualized, or when the sister relation is conceptualized, then the reciprocal sibling relation is also conceptualized. Hence just knowing another person is one's cultural kin is a sufficient basis for social interaction as defined by Talcott Parsons and the set of persons who can mutually recognize each other as cultural kin can form a social system that does not depend on prior face-to-face interaction to adopt appropriate behavior towards one's cultural kin that will likely be reciprocated.

The fact that individuals can compute on the fly whether they have a cultural kin relation implies that the size of the group of socially interacting individuals is limited by the connectedness of mating/marriage networks. Pragmatically, the limit appears to be about 500 individuals, the modal size of hunter-gatherer groups independent

of ecological, climatic, and other environmental conditions and independent of the relative abundance or scarcity of resources. The modal 500 individuals typically do not form a single residence group but are subdivided into smaller residence groups in a manner consistent with resource distribution and methods of resource procurement, yet they maintain coherence as a system of socially interacting individuals through the kinship system expressed through the kinship terminology. Transmitting the cultural kinship computation system as well as associated appropriate behaviors depends on developing children interacting with enculturated adults engaging in those behaviors, who perceive that the children's well-being depends upon their being enculturated with the knowledge that these are behaviors to be directed towards one's kin.

The remaining piece of the evolutionary pathway to be worked out concerns the reasons why a computation system for cultural kin relations should have a generative structure. In theory, composition of relations could simply be carried out for two or three products and no more, or some potential products might simply not be recognized. But in reality, a kinship computation system is logically consistent and can be modeled isomorphically as an algebra (see, for example, Bennardo & Read, 2005).

2.3.5 Cross-Generational Decay of the Functionality of a Set of Relations

We now assume that we have a set S of conceptual relations and an associated set of reciprocal conceptual relations, R , along with a set B of behaviors that are part of the repertoire of a cohort of socially interacting individuals. We assume (1) that some functionality, f , is obtained from the behaviors as a consequence of the shared set of relations, S . We further assume (2) that the total functionality, f , that can be obtained from the behaviors associated with S varies directly with the number of relations in S , which determines the size of the network of interconnected individuals, and (3) that the larger the network the greater the total functionality, f , that can be obtained from interaction among the individuals. But the functionality is subject to an upper bound due to resource constraints (given the mode of resource procurement), and organizational stress (given the mode of societal organization¹⁸).

However, not all individuals connected via relations in the set S will also be individuals for whom the conditions of Fig. 2.10 apply, so we need to distinguish between the potential functionality that is based on the size of S and the actual functionality based upon the actual pattern of interaction among the connected individuals. In mathematical terms, the number of relations in the set S is the cardinality of the set S . Let us refer to f , the functionality associated with S , as the *potential*

¹⁸ The constraints identify conditions under which evolutionary change in either the mode of resource procurement or the organizational structure can occur via cultural group selection as discussed below.

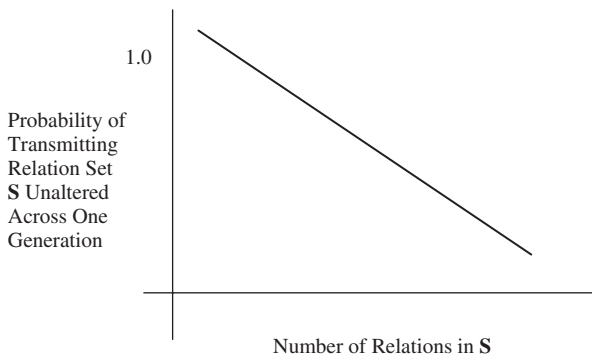
functionality of the set S . It will vary with the cardinality of S . The *actual functionality*, f^* , of S will depend upon the number of co-ordinated individuals in the cohort C . Individuals a and b are coordinated (see Fig. 2.11) if whenever individual a perceives a relation R between a and b with associated behavior B , then individual b perceives the reciprocal relation S between b and a with associated behavior B' . Coordination is thus a consequence of enculturation.

Cohort C will change membership through time due to demographic effects, and to maintain the same functionality, f^* , the new members of the cohort must learn the relations in S and the categories of persons associated with these relations. However, the probability of the set S being transmitted without loss or alteration of content from one generation to the next decreases with the cardinality of the set S (as is shown schematically in Fig. 2.12). Consequently, the actual functionality f^* associated with the set S will decay through time and the greater the cardinality of the set S , the more rapid the decay.

2.3.5.1 Syntactically Organized Versus Syntactically Unorganized Sets of Relations

The decay due to transmission arises from the number of distinct and independent relations that need to be learned.¹⁹ For some relation sets, each relation must be transmitted, for example, when the probability that relation S is in the set S cannot be determined from a subset of S that does not contain S . In such cases, the number of relations that need to be learned is the same as the cardinality of the set S . For other relation sets it is possible to predict the occurrence of a relation S from a subset of S because the set S is *syntactically organized*, i.e. there is a set of rules and initial conditions that suffice to generate the full set of relations from a set of generating relations that is a proper subset of S . When this is the case, then just the rules and

Fig. 2.12 Schematic graph showing a declining probability of faithfully transmitting a set of relations from one generation to the next



¹⁹ It is important from our perspective that there is a qualitative difference between non-human primates and humans, in that humans easily infer simple finite state and phrase-structure grammars, but primates such as the tamarins can only infer patterning in the form of a finite state grammar and not a phrase-structure grammar (Fitch & Hauser, 2004).

the generating relations need to be transmitted. For example, if S consists of the relations $\{M, MM, MMM, MMMM\}$, then one only needs to transmit the relation M , the rule $x \rightarrow xM$, (where x is either M or the output of the rule), and the constraint that the rule can only be applied to strings of length < 4 . If the rules are simple enough to be transmitted faithfully, the relation set will be transmitted faithfully and there will be no loss in the functionality through time (see Fig. 2.13). The decrease in the probability that the relation set S is transmitted faithfully when there is an increase in the cardinality of S and when S is not syntactically organized, drives the relation set S to an algebra-like structure. Due to loss of functionality through transmission error, we can therefore expect change in the content of S to continue until S has a syntactically organized configuration (an algebra-like structure²⁰).

Relations (by virtue of the fact that they may be composed to form new relations) allow a simple syntactic organization consisting of an indefinitely large set of relations S that is generated from a limited set of rules. The generating relations, M , F , S and D (i.e., Mother, Father, Son and Daughter) and the recursive generative rule: “If x and y are either generating relation or the outcome of this rule, then xy is also a relation” generate an infinite set of relations that includes all of (and more than) the relations involved in genealogical tracing from one individual to another.²¹

A rule-based system of genealogical relations not only reduces the likelihood of transmission error, but also provides a way to correct transmission errors due to (a) incorrect calculation of the rule or (b) transmission errors due to direct transmission of relations (rather than through inference from the set of rules and generating relations). Transmission errors can be corrected through majority agreement on the set S as long as the likelihood that a majority of individuals all simultaneously learn the

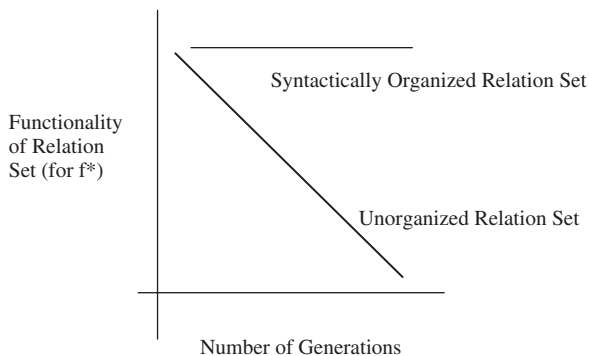


Fig. 2.13 Schematic graphs comparing fitness change for syntactically organized versus unorganized sets of relations

²⁰ The likelihood that there will be faithful transmission can undergo something like a phase transition when changes in the content of S introduce patterning in S so that it now has syntactic organization, thus enabling a switch in transmittal of S from learning each relation in S to learning a set of rules from which S can be reproduced.

²¹ Genealogical tracing is usually limited to either tracing up only, tracing down only, or tracing up and then tracing down. More complicated tracing patterns such as tracing up, then tracing down and then tracing up are excluded from the way genealogical relations are calculated.

same error is low. We can distinguish between systematic and non-systematic errors. Systematic errors are cases in which it is erroneously assumed that a sequence of relations has a “natural” continuation (for example assuming that the genealogical tracing sequence “up, up, down, down” might be followed by two more ups, which runs afoul of the fact that “up, down, up” sequences are usually not allowed). Non-systematic errors occur for other reasons.

The probability that all novices independently make the same non-systematic errors when learning the set S is obviously very low, and individual non-systematic errors can be eliminated through consensus agreement on the content of S . Systematic errors, however, increase the likelihood that a majority of individuals all simultaneously believe that an erroneous relation occurs in the set S , so that consensus would change S to a modified form, S^* . Because a systematic error led to S^* , the set S^* may also be transmitted faithfully. Hence, a form of evolution has taken place, but one that is neither a Darwinian fitness-based evolution, nor a dual inheritance form of evolutionary change. Rather, the evolution from S to S^* has been driven by a change in the set of relations that maintains (or possibly even increases) the likelihood that the set S^* is transmitted faithfully.²²

2.3.5.2 Organizational Implications

Faithful transmission of the conceptual system that underpins its organizational structure is, of course, necessary for organizational continuity in a population. But the fact that the functionality of the system depends upon the social interactions enabled by the relation set complicates our understanding of change during the transmission process. In a Darwinian evolutionary framework, an individual’s conceptual system and cognitive repertoire would have to be included in the person’s phenotype, and we would consider change in the frequency of people with that relation set as due to the manner in which it is transmitted. That approach has two serious flaws.

²² Evolutionary change can also take place when a difference in the decay rates of sets of relations will lead to a difference in the rate of loss of functionality obtained from a relation set. Suppose we have two populations with relation sets S and T in competition (in the Lotka-Volterra sense) over resources. Assume (1) that the (potential and actual) functionality is initially the same for both sets of relations, (2) that the respective carrying capacities of these populations relate directly to the functionalities obtained from their relation sets, and (3) that the competition parameters are such that the two populations are initially in equilibrium. Then, if the decay rate of relation set S is slow compared to the decay rate of set T , the population with set T will find the functionality obtained from its relation set, and thus its carrying capacity, decreases more rapidly than is the case for the other population. If the reductions in carrying capacity are proportional to the rates of decay of the relation sets, then the population with the slower rate of decay of its relation set will win out in competition with the other population. But because the decay rate for a set of relations S is, at least in first approximation, determined by the cardinality of S , and not by the functionality f^* of S , evolutionary changes in functionality caused by differences in decay rates will not be driven by Darwinian fitness since the differences in decay rates are not themselves due to differences in Darwinian fitness but to the cardinality and syntactical organization of a set of relations.

The first, already noted, is the assumption that we are dealing with individual traits, whereas because the functionality arises from a relation set, we must consider a dyad to be the minimum unit. The second flaw is that the functionality accruing to an individual arises from a network of socially interacting individuals, and is not an individual property. It arises from the consequences for an individual of the organizational structure encoded in the conceptual system.

As discussed in Section 2.2, transmission of the organizational structure for a population via enculturation is not simply the sum of trait transmission among individuals. Once a conceptual basis for an organizational system is established, it is maintained by the fact that individuals are born into an ongoing system of enculturated individuals that have, from birth, operated in accordance with the organizational system. Hence, the new-born become enculturated with that system. In large-scale social systems, enculturation need not be uniform across all individuals (see Fig. 2.4); instead enculturation may be regional in scope, but the difference is one of degree and not of kind. It is difficult to see how individuals, except under extraordinary circumstances, could not, to one degree or another, be enculturated in the conceptual system underpinning the society in which they are born and raised.

Organizational systems, as a whole, become larger or smaller as a consequence of demographic changes, including recruitment. The magnitude and form of demographic changes relate to the functionality of the organizational system, hence the changes are due to group properties rather than individual properties and their associated relative fitness values. Variation in individual fitness values, while affecting the frequency distribution of individual traits linked to them, does not determine the demographic trajectory of a population and its organizational system.

2.4 Conclusions

The cognitive evolution we have outlined begins, as it must, as Darwinian evolution driven by individual fitness and extended through biological kin selection and inclusive fitness to behaviors of interacting biologically related individuals. We have argued that the introduction of relational concepts based on categorization of dyadic interactions into the conceptual repertoire of individuals, in conjunction with the cognitive capacity to form recursively defined compositions of relations, had the effect of decoupling the emerging system of conceptually formed categories of relations from its foundation in behavior among biological kin. But the importance of this shift lies not just in what became conceptually possible once the cognitive capacity for recursive reasoning was in place, but also, if not primarily, in the functionality that was thereby obtained. One aspect of that functionality is the manner in which behaviors could now be associated with relational concepts, hence allowing for the extension of kinds of behavior (e.g., co-operation, altruism, reciprocity) otherwise restricted to interaction between biological kin to distantly related (or even non-biologically related) individuals on the basis of the conceptual relation that links

them. To achieve this functionality, though, the conceptual system must be shared by the individuals involved.

What our hominin ancestors worked out, then, was not simply an elaboration on the form of social organization that characterized our common ancestor with *Pan*, our closest non-human primate relative, but social organization based on an entirely different modality. The shift, we have argued, is from social organization based primarily on reproductive and/or individual learning/imitation transmission to social organization based on enculturation transmission of a conceptual system in which relations among societal members are worked out (even though it may draw upon properties transmitted through reproductive transmission or individual learning/imitation transmission). Thereby, the properties of the social system became decoupled from its genetic foundations. The transition assumes Darwinian evolution led to the cognitive capacity for constructing a system of conceptually based social relations. Once in place, the subsequent shift in the basis for social organization gave rise to further evolutionary change driven by the need for social organization to be understood in essentially the same manner by each of the society's members. Change in the basis for social organization supplemented evolution driven at the individual level by an evolutionary process that is not driven by properties of individuals as individuals. This change causes functionality to arise from the form of social organization and the behaviors of individuals *as part of their social system*. Whereas behaviors are seen as the driving force for social organization based upon individual learning/imitation, the social organization and the conceptual system(s) upon which behaviors are based become a driving force for an increasingly complex behavioral repertoire whose transmission among community members was enabled by enculturation. Enculturation transmission can incorporate the increasingly elaborated conceptual dimensions discussed by van der Leeuw et al. in Chapter 3.12 that are the basis for our interaction with the material world through our production and use of artifacts. It can also encompass innovation in information processing and in forms of communication central to evolutionary changes in the scope and complexity of human social organization that are detailed in that chapter.

Reproduction transmission concerns the *implications of individual change for summary population properties*, where the population is determined by the mating patterns leading to reproduction. In contrast, enculturation transmission concerns *the implications social groups have for individuals* by providing the context for interaction between enculturated adults and newborn offspring that leads the child to internalize an ensemble of cultural resources from among those with which s(he) is necessarily involved from birth. Enculturation enables interaction in accordance with a framework for behavior that is functional because other people share the same cultural resources. Comprehensive participation in such cultural resources is virtually a prerequisite for their functionality, and enculturation is the transmission process that enables this.

Enculturation with these cultural resources enables the formation of coherent groups of socially interacting individuals without first requiring experiential learning through prior face-to-face interaction, making social interaction among "strangers" feasible. Identification of a cohort of individuals with shared cultural

resources further depends on a cultural resource we know as a kinship terminology. It enables an individual to determine if s(he) has a cultural kin relation to another individual and vice-versa, so that both must share a wider ensemble of cultural resources individually obtained through enculturation.

Consequently, it is not surprising that the boundaries of social interaction begin with the boundaries of the group of individuals who can mutually recognize one another as cultural kin.²³ The social group thereby was freed from the constraints imposed by the conditions that individual learning through face-to-face interaction impose on social interaction (or the constraints of highly structured mating systems for social interaction to be based on genetically transmitted behaviors). In this manner, a social group can take on functionality far exceeding the forms of social organization available to the non-human primates. New functionality could now be introduced through change in the organizational basis for societies, as expressed through change in cultural resources and tested through cultural group competition (Read, 1987).

Acknowledgments We would like to thank Henry Wright for comments on an earlier draft of this manuscript.

References

- Bennardo, G., & Read, D. W. (2005). The Tongan kinship terminology: Insights from an algebraic analysis. *Mathematical Anthropology and Culture Theory*, 2.
- Biersack, A. (1982). The logic of misplaced concreteness: Paiela body counting and the nature of the primitive mind. *American Anthropologist*, 84, 811–829.
- Bourgine, P., & Johnson, J. (Eds.). (2006). *Living roadmap for complex systems science*. Vol. Version 1.22.
- Boyd, R., & Richerson, P. (1985). *Culture and the evolutionary process*. Chicago: University of Chicago Press.
- Byrne, R. W., & Whiten, A. (1997a). Machiavellian intelligence. In R. W. Byrne & A. Whiten (Eds.), *Machiavellian intelligence II: Extensions and evaluations* (pp. 1–15). Cambridge: Cambridge University Press.
- Byrne, R. W., & Whiten, A. (Eds.). (1997b). *Machiavellian intelligence II: Extensions and evaluations*. Cambridge: Cambridge University Press.
- Chance, M. R. A., & Mead, A. P. (1953). Social behaviour and primate evolution. *Symposia of the Society for Experimental Biology*, VII, 395–439.

²³ Other cultural constructs that involve enculturation can give rise to similar functionality. Lane and Maxfield (2005) observe that “narrative structures are cultural facts of narrative communities . . . [and] narrative . . . rules of syntax . . . are abducted from the many stories instantiating them that circulate in the narrative community, to which members of the community begin listening as infants and continue listening, and then telling, throughout their lives” (p. 13). They consider “narrative logic [to be] a local . . . solution to the problems posed by ontological uncertainty” (p. 15) that arises, in their case study of Echelon (a Silicon-valley start-up company), when innovation is poised to replace current technology. Narration, they argue, gives the narrator “a sense of direction . . . and an understanding of his own character and that of the other actors with whom he is enmeshed” (p. 15); that is, it helps to reestablish group coherency despite the uncertainties introduced as a side effect of innovation.

- Cheney, D. L., & Seyfarth, R. M. (1999). Recognition of other individuals' social relationships by female baboons. *Animal Behaviour*, *58*, 67–75.
- Dasser, V. (1988a). A social concept in Java monkeys. *Animal Behaviour*, *36*, 225–230.
- Dasser, V. (1988b). Mapping social concepts in monkeys. In R. W. Byrne & A. Whiten (Eds.), *Machiavellian intelligence: Social expertise and the evolution of intellect in monkeys, apes, and humans* (pp. 85–93). New York: Oxford University Press.
- Davis, E. W., & Yost, J. A. (2001). The creation of social hierarchy. In R. B. Morrison & C. R. Wilson (Eds.), *Ethnographic essays in cultural anthropology* (pp. 82–120). Belmont: Thomson Publishers.
- de Waal, F. B. M. (2000). Primates – A natural heritage of conflict resolution. *Science*, *289*, 586–590.
- Dunbar, R. (1995). Neocortex size and group size in primates: A test of the hypothesis. *Journal of Human Evolution*, *28*, 287–296.
- Dunbar, R. (2003). The social brain: Mind, language, and society in evolutionary perspective. *Annual Review of Anthropology*, *32*, 163–181.
- Fitch, W. T., & Hauser, M. D. (2004). Computational constraints on syntactic processing in a nonhuman primate. *Science*, *303*, 377–380.
- Flack, J. C., Girvan, M., de Waal, F. B. M., & Krakauer, D. C. (2006). Policing stabilizes construction of social niches in primates. *Nature*, *439*, 426–429.
- Goldberg, T., & Wrangham, R. (1997). Genetic correlates of social behaviour in wild chimpanzees: Evidence from mitochondrial DNA. *Animal Behaviour*, *54*, 559–570.
- Hauser, M. D., Chomsky, N., & Fitch, W. T. (2002). The faculty of language: What is it, who has it and how did it evolve? *Science*, *298*, 1569–1579.
- Hughes, C. (2004). What are the links between theory of mind and social relations? Review, reflections and new directions for studies of typical and atypical development. *Social Development*, *13*, 590–619.
- Humphrey, N. K. (1976). The social function of intellect. In P. P. G. Bateson & R. A. Hinde (Eds.), *Growing points in ethology*. Cambridge: Cambridge University Press.
- Kottak, C. (2004). *Window on humanity: A concise introduction to anthropology*. New York: McGraw-Hill.
- Kummer, H. (1967). Tripartite relations in Hamadryas baboons. In S. A. Altmann (Ed.), *Social communication among primates*. Chicago: University of Chicago Press.
- Lane, D., & Maxfield, R. (2005). Ontological uncertainty and innovation. *Journal of Evolutionary Economics*, *15*, 3–50.
- Lee, R. (1990). *The Dobe Ju'hoansi* (3rd ed.). Toronto: Wadsworth.
- Marshall, L. (1976). *The !Kung of Nyae Nyae*. Cambridge: Harvard University Press.
- McGrew, W. C. (1992). *Chimpanzee material culture: Implications for human evolution*. Cambridge: Cambridge University Press.
- Mitani, J., Merriwether, D. A., & Zhang, C. (2000). Male affiliation, cooperation, and kinship in wild chimpanzees. *Animal Behaviour*, *59*, 885–893.
- Mitani, J. C., Watts, D. P., Pepper, J. W., & Merriwether, D. A. (2002). Demographic and social constraints on male chimpanzee behaviour. *Animal Behaviour*, *64*, 727–737.
- Morin, P., Moore, J., Chakraborty, R., Jin, L., Goodall, J., & Woodruff, D. (1994). Kin selection, social structure, gene flow, and the evolution of chimpanzees. *Science*, *265*, 1145–1332.
- Nishida, T., & Hiraiwa-Hasegawa, M. (1987). Chimpanzees and bonobos: Cooperative relationships among males. In B. B. Smuts, D. L. Cheney, R. M. Seyfarth, R. Wrangham, & T. T. Struhsaker (Eds.), *Primate societies* (pp. 165–178). Chicago: University of Chicago Press.
- Parkin, R. (1997). *Kinship: An introduction to basic concepts*. London: Blackwell.
- Parsons, T. (1964). *The social system*. New York: Free Press of Glencoe.
- Pigeot, N. (1991). Reflexions sur l'histoire technique de l'homme: De l'évolution cognitive à l'évolution culturelle. *Paléo*, *3*, 167–200.
- Povinelli, D. J., Nelson, K. E., & Boysen, S. T. (1990). Inferences about guessing and knowing by chimpanzees (*Pan troglodytes*). *Journal of Comparative Psychology*, *104*, 203–210.

- Premack, D., & Woodruff, G. (1978). Does the chimpanzee have a theory of mind? *Behavioral and Brain Sciences*, 4, 515–526.
- Pusey, A. E., & Packer, C. (1987). Dispersal and philopatry. In B. B. Smuts, D. L. Cheney, R. M. Seyfarth, & T. T. Struhsaker (Eds.), *Primate societies*. Chicago: University of Chicago Press.
- Read, D. W. (1984). An algebraic account of the American Kinship terminology. *Current Anthropology*, 25, 417–440.
- Read, D. W. (1987). Foraging society organization: A simple model of a complex transition. *European Journal of Operational Research*, 30, 230–236.
- Read, D. W. (2001). What is kinship? In R. Feinberg & M. Ottenheimer, *The cultural analysis of kinship: The legacy of David Schneider and its implications for anthropological relativism* (pp. 78–117). Urbana: University of Illinois Press.
- Read, D. W. (2004). The emergence of order from disorder as a form of self organization. *Computational & Mathematical Organization Theory*, 9, 195–225.
- Read, D. W. (2005). Change in the form of evolution: Transition from primate to hominid forms of social organization. *Journal of Mathematical Sociology*, 29, 91–114.
- Read, D. W. (2006). Working memory: A cognitive limit to non-human primate recursive thinking prior to hominid evolution? *Cognitive Science Journal Archive*, 2674–2679.
- Read, D. W., & van der Leeuw, S. (2008). Biology is only part of the story. *Philosophical Transactions of the Royal Society B*, 363, 1959–1968.
- Reader, S. M., & Laland, K. N. (2002). Social intelligence, innovation, and enhanced brain size in primates. *Proceedings of the National Academy of Sciences (USA)*, 99, 4436–4441.
- Rodseth, L., Wrangham, R. W., Harrigan, A. M., & Smuts, B. B. (1991). The human community as a primate society. *Current Anthropology*, 32, 221–254.
- Sahlins, M. (1962). *Moala: Culture and nature on a Fijian Island*. Englewood Cliffs: Prentice-Hall.
- Savage-Rumbaugh, E. S., Rumbaugh, D. M., & Boysen, S. T. (1978). Sarah's problems in comprehension. *Behavioral and Brain Sciences*, 1, 555–557.
- Schwartz, T. (1981). The acquisition of culture. *Ethos*, 9, 4–17.
- Shettleworth, S. (1998). *Cognition, evolution and behavior*. Oxford: Oxford University Press.
- Silk, J. B. (2003). Cooperation without counting: The puzzle of friendship. In P. Hammerstein (Ed.), *Genetic and cultural evolution of cooperation* (pp. 37–54). Cambridge: MIT Press.
- Spinozzi, G., Natale, F., Langer, J., & Brakke, K. E. (1999). Spontaneous class grouping behavior by bonobos (*Pan paniscus*) and common chimpanzees (*P. troglodytes*). *Animal Cognition*, 2, 157–170.
- Tomasello, M. (1998). Uniquely primate, uniquely human. *Developmental Science*, 1, 1–30.
- Tomasello, M. (2000). *The cultural origins of human cognition*. Cambridge: Harvard University Press.
- van der Leeuw, S. E. (2000). Making tools from stone and clay. In P. Anderson and T. Murray (Eds.), *Australian archaeologist: Collected papers in honour of Jim Allen* (pp. 69–88), Coombs, Australia: Coombs Academic Publishing.
- Vigilant, L., Hofreiter, M., Siedel, H., & Boesch, C. (2001). Paternity and relatedness in wild chimpanzee communities. *Proceedings of the National Academy of Sciences (USA)*, 98, 12890–12895.
- Vinden, P. G. (2004). In defense of enculturation. *Behavioral and Brain Sciences*, 27, 79–151.
- Whiten, A., & Byrne, R. W. (Eds.). (1988). *Machiavellian intelligence: Social expertise and the evolution of intellect in monkeys, apes, and humans*. Oxford: Clarendon Press.
- Woodruff, G., & Premack, D. (1979). Intentional communication in the chimpanzee: The development of deception. *Cognition*, 7, 333–362.

Chapter 3

The Long-Term Evolution of Social Organization

Sander van der Leeuw, David Lane and Dwight Read

3.1 Introduction: Linking the Dynamics of Innovation With Urban Dynamics

In the first chapter of this book, David Lane et al. point out that the Darwinian approach to biological evolution is insufficient for the description and explanation of the cultural and social transmission of many, if not most, of society's characteristics. Instead, the chapter proposes that we shift from 'population thinking' to 'organization thinking' to understand socio-cultural phenomena. In essence, such thinking focuses on the role of information in shaping institutions and societies. In the second chapter, Dwight Read et al. outline a crucial stage in the evolution of human societies, which they call 'The Innovation Innovation'. It concerns the beginnings of information processing by (small-scale) societies about societies. They outline, in a few steps, how human beings may have developed a conceptual apparatus to describe and to manage their own bio-social (kinship) relations. The main innovation in this story is the capacity to abstract from substantive observations of such relationships to concepts that encapsulate the underlying structure of these relationships.

The current chapter continues the story, outlining how the innovation innovation transformed the world of our distant ancestors into that in which we live today. It focuses on the relationship between people and the material world, as it is the material world that has been most drastically, and measurably, transformed over the last several tens of thousands of years. In view of what we know about such distant periods, and in view of the space allotted to us here, it will not surprise the reader that we do so in the form of a narrative that is only partly underpinned by substantive data.¹ We emphasize this because we do not want to hide from the reader

S.E. van der Leeuw (✉)

School of Human Evolution and Social Change, Arizona State University, P.O. Box 872402, Tempe, AZ 85287-2402, USA
e-mail: vanderle@asu.edu

¹ Nevertheless, we have referred the reader to other literature where that seemed appropriate.

the speculative nature of the story that follows. Yet we firmly believe that, in very general terms, this scenario is correct, and that further research will vindicate us.

We first give examples of the kinds of abstractions, and the hierarchy of conceptual dimensions necessary for prehistoric human beings and their ancestors, to conquer matter, i.e. to conceptually understand, transmit and apply the operations needed to master the making of a range of objects made out of stone, bone, wood, clay and other materials. Some of the abstractions that had to be conceived in this domain resemble those that Read et al. refer to (and may therefore have been transposed from one domain to another), while others apply to this domain alone, and had to be truly ‘invented’. It is then argued that such ‘identification of conceptual dimensions’ is a process that underlies all human activity, and we look a little closer at how that process relates to invention and innovation.

Lastly, we shift our attention to the role of innovation, information processing and communication in the emergence of social institutions, and in the structural transformation of human societies as they grow in size and complexity. In particular, we look at the role that problem solving and invention play in creating more and more complex societies, encompassing increasing numbers of people, more and more diverse institutions, and an – ultimately seemingly all-encompassing – appropriation of the natural environment. To illustrate this development we will focus on the origins and growth of urban systems, as we have done for the ISCOM project as a whole.

3.2 Why did it Take Us So Long to Become Inventive?

Human beings and their precursors have lived on this Earth for several million years. In their current guise (*Homo sapiens*), they have roamed over its surface for around 195–160,000 years. For the great majority of that long time-span, our species moved around in small groups, using very basic tool-making techniques to ensure modest success in defending itself and obtaining the necessary foods for its subsistence. Although there are clear signs of innovations from at least 50,000 years ago, the rate of innovation increased dramatically from about 10,000 years ago. From that moment on, in a relatively short time, human beings have managed to upend the balance between themselves and their natural environment, and gain something approaching control over many environmental processes. The history of our species therefore raises three interesting questions:

1. How did the species survive for so long with a minimal toolkit to defend and nourish itself, and under a wide range of natural circumstances?
2. Why did it take so long to ‘invent’ and accelerate innovation?
3. Why did innovation increase so rapidly, once that point was reached?

Clearly, the long-term survival of the human species depended, and depends, on the adequacy of human behavior with respect to the environment. Most human behavior was of course routine, and well-adapted to known circumstances, but whenever people encountered unknown phenomena, they initially suffered until they

learned how to deal with them. If that took too much time, the consequences were dire. Learning and adaptation constituted the key to survival (as they do now!).

For much of human history adaptation was favored by the fact that people lived in small groups, had an essentially mobile lifestyle, and roamed over large territories. Their capacity to extract the necessary resources for survival may have been lower than at present because highly efficient extraction technologies were not available, but the small size of the groups and the large size of groups' territories made it possible to gather sufficient food to survive on, and the mobility ensured that if a group could not do so in one place, it had a good chance of finding another, better, location before it was too late.

All in all, for a very long time, human populations thus lived within fairly narrow constraints, surviving on whatever the land offered and moving on when that was not enough. The major areas of invention and innovation concerned the immediate interface between people and their environment: tools and procedures to enhance the impact of human actions, and to extend the range of resources that could be used. Examples of such innovations are the control over fire, the manufacture of clothing and weapons, the construction of shelter, etc.

However, that does not explain why this way of life persisted for so long. If, in the last 10,000 years, the species managed to increase so dramatically the rate of innovation, why did it not do so before? One could argue that external perturbations forced the species to accelerate learning, invention and innovation. But, if external perturbations are at the origins of the acceleration of invention, why did this not happen during earlier or later periods of (often much more drastic) transformations of the external circumstances under which human societies lived? And why, once the threshold had been crossed, could invention and innovation accelerate exponentially, independent of external circumstances? There must have been other factors at play. . . . It seems to us that the answer to both these questions lies in the fact that achieving a certain level of development of the human conceptual apparatus was the necessary and sufficient condition for the acceleration of invention and innovation. The next section of the current chapter will be devoted to answering the second of the three questions: "Why did it all take such a long time?" (cf. Wobst, 1978).

3.3 The Earliest Tool-Making and the Conceptualization of Three Spatial Dimensions

In the last chapter, Read et al. called this threshold "the innovation of innovation," and argued that, in the realm of kinship relations, crossing it essentially entails attaining the capacity to abstract and generalize locally-made observations to a much wider realm that includes unobserved situations. Here, we are going to extend that argument to the domain of the relations between people and matter, because *in this realm, we can assign dates to some of the steps involved, and we can thus substantiate the claim that it did indeed take a very long time to innovate innovation*. In the process, we will point to the invention of other conceptual dimensions and operators that were needed to 'conquer the material world'.

For most of human history, there is only one category of artifacts that allows us to monitor the development of human cognition: flaked stone tools. In what follows, we will base ourselves on Pigeot (1992), as well as on a more extensive paper one of us has written on this topic some years ago (van der Leeuw, 2000). Pigeot has presented us with an outline of the development of human cognition as reflected in the ways and means for knapping stone tools. In essence, she argues that one may distinguish five stages in the development of the techniques to make stone tools: the conceptualization of (1) the point, (2) the angle, (3) the edge, (4) the surface and finally (5) the volume (Figs. 3.1–3.4).²

Dimension 0: the point

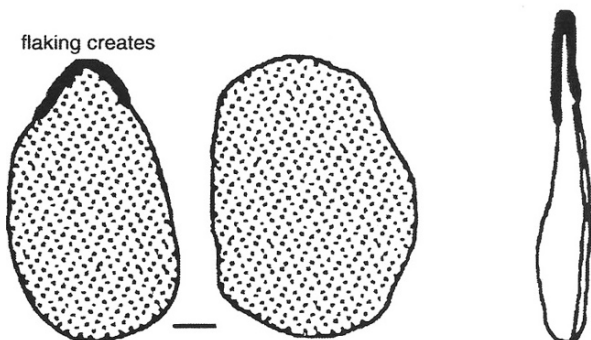


Fig. 3.1 The cognitive capacities of the preliminary stage (after Pigeot 1991)

First Dimension: the line

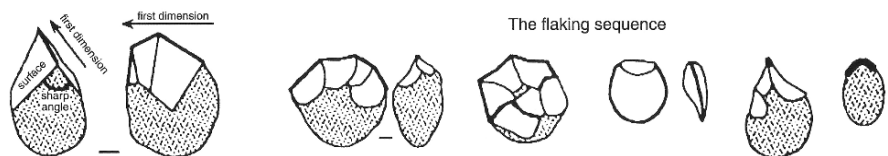


Fig. 3.2 The cognitive capacities of the first stage (after Pigeot 1991)

² Inevitably, since the publication of this paper, colleagues have disputed some of its finer points of technology, as well as the relatively simple chronological sequence in which Pigeot cast the ideas. But for the purposes of our argument here, these points are less relevant than her final conclusion – that by the end of the Paleolithic, flint-knappers had mastered the conceptualization of objects in three dimensions.

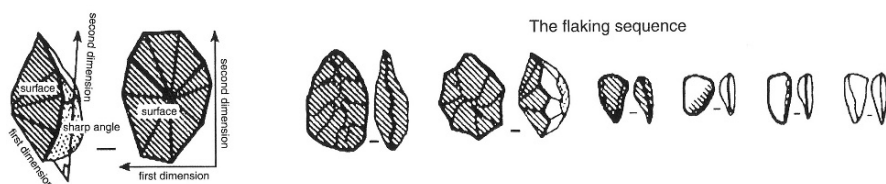
Second Dimension: the surface

Fig. 3.3 The cognitive capacities of the second stage (after Pigeot 1991)

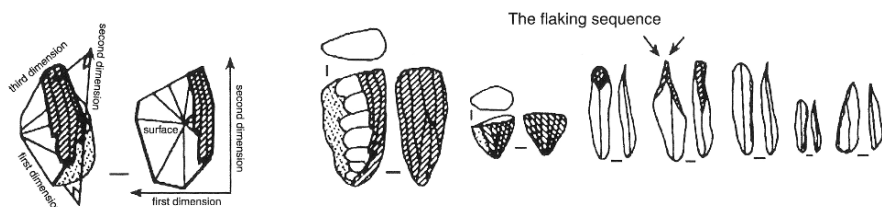
Third Dimension: the volume

Fig. 3.4 The cognitive capacities of the third stage (after Pigeot 1991)

3.3.1 The Concept of Points: The Very Earliest Tools of Primates and Humans

Fundamental for all deliberate tool making is the capacity to invert causal sequences in the mind (“A causes B, therefore if I want B to happen, I have to do A”), transforming an observation into a deliberate action (e.g., Atlan, 1992). Humans share this capacity with many primates. Both are therefore theoretically able to manipulate (aspects of) the material world.

Primates and early humans use stone tools both actively (as hammers, or as objects to throw) and passively (as an anvil, for example). These uses may be combined, for example, in cracking nuts or bones. Like human beings, primates are therefore aware of certain properties of their tools, such as their weight, their hardness and their resistance to shock, as well as of the motions they can execute with them.

Finally, primate tools and the earliest human tools are very much alike. Both show the impact of blows, because flakes have been struck off at that point, which in turn indicates that the angle at which these stones hit each other was smaller than 90° . But there the similarities end. Primates never learned to conceptualize either the *point as an abstract object*, or the *angle as a relational concept linking two lines or planes*. As a consequence, they are not able to deliberately shape the tools they use. In this, they are different from (proto) humans.

3.3.2 Thinking in One Dimension: Edge Tools

As early as two million years ago, the first collections of broken pieces of stone are identified as tools. Although it is not sure that they have actually been deliberately shaped, they share one feature: that of a (cutting) edge. To make such objects deliberately is presumably a sign that some conceptual capacities have been developed. The first such tools, the so-called choppers and chopping tools, show that several adjacent flakes have been removed deliberately by hitting the stone with another one. In each case, the toolmaker focused on sharpening only (part of) the tool's edge. In so doing, he followed the original shape of the stone; there is no attempt to achieve control over the whole shape. Nevertheless, the alignment of blows shows a degree of deliberation and choice in the way stones were hit, and shows that tool preparation has moved from a point-based conception towards a line-based conception (Fig. 3.1). Moreover, the fact that, at this point in time, both one-sided and two-sided flaking occurs, confirms that the angle has been conceptualized. The edge occurs where it best suits the natural form of the pebble.

Finally, although both the object from which a flake was removed and the flake itself may be used in further actions, from the toolmaker's perspective they are viewed differently. The former may be further modified by flaking, and is therefore the true object of the toolmaker's attention, while the latter is not – it is a by-product rather than an object.

3.3.3 Thinking in Two Dimensions: Surface Tools

At some point in time, our ancestors extended their linear conception of tool-making by removing flakes all around the edge of a pebble or stone, making what Pigeot calls 'discoïdal tools.' By thus closing the knapped line onto itself, they implicitly defined a plane or surface. The next conceptual step concerned the transition from defining a surface by instantiating its perimeter to defining an edge by instantiating the surface within it, transforming such tools from objects consisting of an edge flanked by two planes into objects consisting of two planes between which one finds an edge. Identifying the transition is a question of reconstructing the sequence in which the tool is made: has the maker first sharpened the edge of the tool and then taken large flakes off both surfaces, or has he done the reverse? In practical terms, this does not change the shape much, or the function or the effectiveness of the resulting tool, but the conceptual step is a major one for the toolmaker. It implies a move from a tool conceived in terms of one-dimensional conceptual objects (edges) to one conceived in terms of two-dimensional conceptual objects (surfaces). Once that step has been taken, tool making becomes inherently a matter of dealing with surfaces. Around 250–300,000 years ago, this leads to the development of a special technique (called 'Levallois,' see Pigeot, 1992, p. 184–186, and Fig. 3.2) in which the removal of one flake is at the same time the preparation for the removal of the next one. In the process, the makers substantively increased control over the shape of their products.

3.3.4 *Thinking in Three-Dimensions: Tools Conceived as Volumes*

Whereas the Levallois technique exploits the stone by knapping on two crosscutting planes and creating flake tools with one edge where two planes meet, the Upper Paleolithic knappers work at the intersection of *three* planes. They create long blade tools with three edges, at each of which two out of the blade's three planes meet (Fig. 3.3). The Upper Paleolithic nucleus is thus a volume in the true sense; it is prepared in different ways, but it always consists of three crests that guarantee an optimal exploitation of the nucleus because the flaking reduces the volume everywhere in turn (Boëda, 1990). Pigeot adds that the volume that is thus defined is exploited on the smaller of its surfaces, so that the volume literally is more important than the surface, and the management of the core volumetric rather than planar. A volumetric conception of tool manufacture is attested in the Gravettian and Magdalenian traditions (c. 27,000 and 13,000 years ago, respectively). It is conceptually and economically very efficient, and it involves the simultaneous mastering of new knapping techniques, such as soft percussion flaking, which increase control over the way flakes and blades are removed from the core. As a result, a completely new range of small blade-based stone tools is invented (see below).

3.3.5 *What are the Conceptual Advances So Far?*

Depending on when one assumes the above collective learning process to have started, it seems to have taken human beings and their ancestors one to two million years (or more) to achieve true volumetric conceptualization of the manufacture of stone tools. At that point in time (say, 20,000 years ago), they acquired the capacity to conceive of *each volume as constituted of an infinity of surfaces, each surface of an infinity of lines and each line of an infinity of points, and vice versa*. This has important implications. As long as the object (tool or flake) is conceptualized in fewer than three dimensions, but in reality exists in all three, full manipulation of matter is impossible. It is only when all three dimensions of material objects are conceptualized that full control can be achieved over the making of stone tools. A threshold is thus crossed, which, from this time onwards, allows for a major acceleration in the development of human control over matter.³

But in the process of acquiring the capability to conceive of and to make stone tools in three dimensions, our ancestors had also acquired a number of other conceptual tools. One of these is the *capacity to plan and execute complex sequences of actions*. As long as individual flakes are being removed from a tool to create an edge that follows no predetermined pattern (as is the case for many of the earliest tools), there is little or no anticipation. The controlling feed-back loop at most relates the removal of a flake *to the past*, to the removal of the last one. And even when the edge comes full circle, the result of blow $n+10$ is rather under-determined by the result of blows n to $n+5$ (Pigeot, 1992, p. 182). When the Levallois technique is

³ For an extended discussion of the above developments, as well as those pointed to in the next few pages, see van der Leeuw (2000).

introduced, this approach changes radically. Now, the knapper has to look several steps ahead *into the future* in evaluating the results of past actions, so that each removal prepares the way for future steps. Initially, the span of such control loops⁴ is relatively short. However, in the case of the industries that use a truly volume-based approach, the preparatory phase is quite elaborate, involving the creation of a strike platform, as well as the preparation of the surfaces from which the blades are to be removed. Once all that is done, the blades are removed one after the other without intervening reshaping. From the sequencing point of view, this implies a stringent separation between ‘preparation’ and ‘exploitation’ of the raw material. The Levallois technique is thus an early instance of a *manufacturing sequence in which the total process is divided in different phases that are not interleaved*. Being able to conceive and manage such sequences in turn testifies to the fact that the toolmakers developed the conceptual capacity to link different steps in a process together in such a way that one might speak of a plan. Their repeated use of the different knapping sequences has made it possible to identify the steps involved, and, to some extent, the cause-and-effect relationships between them. But what does this require at the conceptual level?

To *identify* cause-and-effect one must have inserted a control loop between observations and conceptualizations. To *(re)create* a transformation at will, the individual must be able to retrieve the whole sequence associated with the desired result, to ‘wind it back’ to its beginning, and to ‘replay’ it in the appropriate order. To *understand and manipulate* the dynamics of cause and effect, the individual must be able to play such sequences backwards and forwards in an interactive way and to retrieve (parts of) sequences that are associated with any of the stages of transformation or actions concerned. Moreover, such understanding also requires the capability to observe differences between manufacturing sequences, and to generate variations by mentally or physically inserting operations or modifying them. *That in turn requires the capacity to mentally associate different strings of events, for example, on the basis of an assessment of similarities and differences between the transformations and/or the products at certain stages*. It requires full conceptualization of all the steps and their interconnections, hidden and apparent, so as to be able *to anticipate all different parts of the manufacturing sequence*, and to create the right conditions for them to be implemented and controlled.

3.4 Creating a New Material World

From around 12,000 years ago, we clearly observe a drastic acceleration in the speed of invention and innovation. Many new categories of artifacts emerge, new materials are used, new techniques are introduced, and new ways to deal with aspects of the material world are ‘discovered’ in a comparatively short time – a span of a few

⁴ We will here use the term “control loop” instead of the more common “feed-back” or “feed-forward” loop because it seems to us that feed-back and feed-forward are always combined in monitoring the actions undertaken as part of a manufacturing sequence.

thousands of years. The acceleration is so overwhelming, that in that time-span, the whole way of life of many humans on earth changes: rather than live in small groups that roam around the Earth, people concentrate their activities in smaller territories, they invent different subsistence strategies, and in some cases, they literally settle down in small villages.

If we look forward in time from that point, the change is even more dramatic: within a couple of thousands of years more, people congregate in much larger groups, all kinds of new social institutions are instantiated, technological inventions explode and reach very different realms when people start building cities, communicate in writing, etc. All in all, it is clear that just before the beginning of the Neolithic, a threshold of innovativeness had been crossed. This section is in part devoted to the description of some of the conceptual tools acquired, but its aim is to answer the question what made this acceleration possible.

3.4.1 *New Kinds of Tools in Stone, Bone, Wood, Etc.*

During the tail end of the Upper Paleolithic and the Mesolithic, we see a rather large number of important new techniques emerge almost simultaneously.⁵ One such new development is the manufacture of a wide spectrum of smaller and smaller stone tools. These so-called ‘*microliths*’ are very finely made, and show that the makers are extending their control over manufacture to finer and finer details, something that would not have been possible if these objects had not been conceived in three dimensions, and that the manufacturing sequences were planned in detail. It testifies to the *extension of the range of orders of magnitude of volume* manipulated by toolmakers.

As part of this process, we see an increasingly wide range of differently shaped objects, which implies that there is a increasingly close match between individual objects and the functions that they are meant to fulfill; that in turn suggests that tool-makers have acquired a *more versatile spatial topology*, and an improved capability to analyze the requirements their artifacts should meet in order to fulfill their intended function most effectively.

A closely related innovation is the introduction of *composite tools*, consisting of a number of microliths hafted together in objects of wood or bone. This phenomenon is interesting from the conceptual point of view, as it implies a certain *reversibility of scalar hierarchies*: not only are tools made by reducing a larger piece of stone into one or more small flakes, which are then retouched to give them the required shape, but such small pieces are assembled into something bigger.

Although it is unlikely that this moment represents the first use of non-stone materials as tools, one now observes the use of other materials (wood, bone, antler, etc.) alongside stone in new, composite, tools. Some of these tools fit the new, small

⁵ In archaeology, it is often very difficult to determine the sequence in which phenomena appear, in part because of either a lack of dates or because dates have a wide margin of error, but also because our record is often so fragmentary that it is very easy to miss the first manifestation of a phenomenon.

stone blades in larger objects made from bone, for example. The fact that such a wide range of new materials was used testifies in yet another way to the increasing innovativeness of human beings around this time, as it meant developing a wide new range of (motor and other) skills and tools to work all these materials.

3.4.2 The Introduction of Ground Stone Objects

Alongside these newly emerging techniques, knapping of larger stone tools clearly continues. From about 10,000–7,000 years ago, however, these tools are transformed beyond recognition as toolmakers discover grinding. Grinding tools in the last stage of manufacture achieves much better control over the final shape and enables toolmakers to create shapes that until then had been beyond their capability. In many ways, this development caps and completes the mastery of stone-working at all scales. Objects such as Neolithic stone axes and adzes are first roughly flaked out of appropriately fine-grained blocks of stone. Next, they are refined, by removing smaller and smaller flakes. Finally, the toolmaker removes particles of infinitely small size by pecking or grinding. The resulting objects have a completely smooth surface, which can be as flat, rounded, or irregular as desired. *Control over the final shape is complete, as is the use of different scales of removal from the initial stone block – from very large flakes to individual grains.* That control leads, ultimately (in the British Late Neolithic, for example), to very highly standardized production of very large numbers of polished stone axes (cf. Bradley & Edmonds, 1993).

3.4.3 The Introduction of Containers

The making of containers was invented anywhere between 12,000 and 9,000 years ago (depending on the material and the part of the world one looks at). Such containers occur in wood, leather, stone and pottery. In each case, the actual manufacturing technique is of course different, but the conceptual underpinnings are the same. They combine different innovations that emerge subsequent to three-dimensional conceptualization of artifact manufacture (van der Leeuw, 2000):

- The introduction of topologically fundamentally different objects, consisting of *solids around a void*. This requires the conceptual separation of the surface of an object from its volume, and making the distinction between outside and inside surfaces. Neither is conceivable in the absence of a true 3-D conception of objects.⁶ In some parts of the world, gourds may have provided an example, and other natural containers may have served that function elsewhere. Nevertheless,

⁶ One craft in which the ‘discovery’ of such hollow shapes may have occurred is leather-working, where skins are removed from one kind of object (the animal), only to be transformed into another, differently shaped, object (the container).

the step from observing a natural object to re-creating it conceptually was an important innovation.

- The inversion of the sequence of manufacturing, *beginning with the smallest particles and assembling them into larger objects*. Basketry and weaving are examples of this. In each case, one assembles small linear objects (animal hairs or vegetal fibers) into longer and thicker ones that are subsequently assembled into two- or three-dimensional surfaces (cloth, baskets). Although in the case of pottery the smallest particles (the clay platelets) are found in nature, the earliest shaping techniques (coiling; shaping in a net) are closely related to basketry.
- One of the main advantages of such additive techniques over subtractive ones is the fact that one can *correct errors* by going one or more steps back in the procedure, making the correction, and then proceeding again. This presumes not only that control loops link the past with the present and the future, but also that, in this particular domain, *actions are conceived as being reversible*.
- In the case of pottery, the *separation between different stages of production* is also pushed a step further. Resource procurement, clay preparation, shaping, decoration, and firing occur one after another, and, during the whole manufacturing process, the maker has to keep all later stages in mind. The choice of raw materials, for example, is intimately linked to a pottery vessel's shape, function and firing conditions. Manufacture involves a large number of embedded control loops, and small variations early in the process do have major consequences later. It is therefore important that errors are easily corrected.

3.4.4 *Inventing the Conceptual Tools to Conquer the Landscape*

These conceptual advances also opened up completely new realms of problem-solving and invention, including the transformation of subsistence risks from a daily concern over which people had little control, to a seasonal or pluri-annual concern over which they had a little more control. That transformation was achieved by (1) a mobile lifestyle with the breeding and herding of domesticated animals, or the seasonal cultivation of wild plants, as principal subsistence strategy, or by (2) settling in one place, building houses, clearing fields, and cultivating domesticated plants.⁷ In both cases, the result was *long-term human investment in certain aspects of the environment*, which cut off some hitherto existing options for flexible interaction with the environment, even as it opened up new opportunities for social and cultural development. As we will see, the concomitant 'inventions' are difficult to imagine without the conceptual changes in artifact manufacture that we discussed.

Gathering and hunting are essentially intermittent, almost instantaneous (albeit periodic) interactions between the various temporalities that rule the natural

⁷ The fact that different crops and different animals became domesticated in different parts of the world seems to argue against the spread of domestication as a 'technique,' but not necessarily in favor of independent invention. The invention of different subsistence activities may simply have been enabled by the changes in the conceptualization of space and time.

environment and the rhythms of human subsistence needs. As the human beings concerned need only to be at the right time in the right place to feed themselves, such 'culling' only requires descriptive knowledge in which space-time can be represented as strings of (favorable or unfavorable) encounters between people and the landscape. It suffices to include the season in such descriptions to make them effective as subsistence 'manuals.' Herding, cultivation and domestication, on the other hand, involve a longer-term, intimate symbiosis between humans and their food sources, in which people influence the natural processes. That presupposes the existence of certain conceptual functions at the spatial scale of the landscape and the temporal scale of years.

Spatially speaking, this requires a *two-dimensional map of the landscape* and, in the case of houses and cultivation, the conceptual distinctions between '*inside*' and '*outside*,' – marked by the walls of a house or the perimeter of gardens or fields – as well as between '*self*' and '*other*' that is acquired as part of the conceptualization of kinship systems. But it also involves the extension of two-dimensional conceptualization of space to surfaces larger than those of tools, and the distinction between inside and outside that is characteristic of pottery-making.

Temporally speaking, in the case of cultivation clearance, planting or seeding are separated by several months, if not years, from harvesting a crop, whereas in the case of herding it takes years to build up a herd of sufficient size to entrust the survival of the group to symbiosis with a particular herd. We thus see the same mechanisms at work as in artifact manufacture: stretching of temporal sequences and temporal separation between different parts of a 'manufacturing' sequence.

We can infer there is also a *change in people's relationship with time and space at the scale of the landscape*. During the Mesolithic, artifact distributions point to increasing circumscription of the areas within which each human group moves around. A little later (10,000–7,000 years ago, depending on the region), this may result in a settled lifestyle. In the absence of archaeological data, we must look to ethnography in order to understand the implications of these changes for people's perception of time and space.

The Australian Aborigine use of 'song-lines' provides an example of mobile space-time perception. Song-lines are strings of chant sung while traveling through an unknown landscape. They allow the traveler to find his way by matching the song with what one sees while traveling. The landscape cannot be interpreted without the song, nor does the song make sense without the landscape. Individual song-lines are learned by rote, as people do not return to any area frequently enough to acquire the necessary knowledge by experience. The song provides a guideline through a landscape because it invokes time as an independent dimension to interpret space, stringing a series of punctual perceptions of space into a sequence.

In many sedentary cultures, on the other hand, spatial perception is encoded on a map. There are, again, two dimensions, but both are spatial. Three factors facilitate the transition. Firstly, settling down provides fixed points (settlements) around which a two-dimensional map can be organized. Secondly, frequent movement over limited distances replaces long-distance movement, so that every trajectory between any two points is observed from every direction, and the relationships between these

trajectories can be memorized. Thirdly, settling down provides the temporal continuity of observation needed to unravel the respective roles of the spatial and temporal dimensions in observed changes. Together, these factors are necessary to enable people to develop two- and three-dimensional conceptions of the landscape.

3.4.5 The Impact of the Invention Explosion

In conclusion, we would argue that the ‘invention explosion’ of the Neolithic is the result of the fact that human beings have internalized the conceptual apparatus necessary to conceive of space in four nested dimensions (0, 1, 2, 3) across a wide range of spatial scales (from the individual fiber or grain to the landscape), to separate a surface from the volume it encloses, to use different topologies, to distinguish and relate time and space, to distinguish between different sequences of cause and effect, and to plan, etc.

Together, these conceptual advances greatly increased the number of ways available to tackle the challenges posed by the material environment. That allowed them to meet increasingly complex challenges in shorter timeframes. Hence, it triggered a rapid increase in our species’ capability to invent and innovate in many different domains, substantively increasing humans’ adaptive capacity. It is as if, rather suddenly, human beings had achieved an exponential increase in the dimensionality of the conceptual hyperspace (‘possibility space’) that governed their relationship with the external world. This afforded them a quantum leap in the number of degrees of freedom of choice they had in dealing with their material and ideational environment.

But the other side of the coin was that these solutions, by engaging people in the manipulation of a material world that they only partly controlled, ultimately led to new, often unexpected, challenges that required the mobilization of great effort to be overcome in due time. The fact that, in the process, human societies invested more and more in control over their environment (such as by building infrastructure), anchored them more and more closely to the territory in which they had chosen to live. The symbiosis that thus emerged between different landscapes and the life-ways invented by human groups to deal with them eventually irrevocably channeled the ways in which the societies concerned could interact with their physical environment, driving them to devise increasingly complex solutions, with more unexpected consequences resulting. Overall, therefore, increasing control over the material and natural environment was balanced by increasing societal complexity, which was not always simple to keep under control.

3.5 Invention, Innovation and Collective Problem Solving

In the remainder of this paper, rather than attempt to outline, necessarily very poorly, the innumerable individual conceptual inventions made by humankind since the innovation threshold has been crossed, we will assume that the process of invention

and problem generation accelerates from the Neolithic onwards, and focus our attention at the aggregate level, investigating the processes that were triggered at the level of whole societies by the crossing of the invention threshold.

Returning with the benefit of hindsight to the questions asked at the beginning of the paper, we can now answer them by describing the emergence of the human species as a dominant player on Earth as a bootstrapping process in which humans slowly gained an edge over other species and over their physical environment by using the faculty that distinguished them from all other species: the capacity to learn and to learn how to learn.⁸ That capacity allowed them to categorize, make abstractions and hierarchically organize these abstractions, and, in so doing, to develop their capacity to identify and solve problems by inventing suitable conceptual tools. They learned various kinds of (symbolic and other) means to communicate among themselves, and they increased their capacity to transform their natural and material environment in many different ways and at many spatial and temporal scales.

A ‘shorthand’ description of this bootstrapping process as it occurs in any individual looks more or less like this:

1. A trial-and-error process identifies conceptual dimensions that summarize observations and experiences in a particular domain, so that these can be stored and transmitted in an economic and efficient manner;
2. The more such dimensions are available, the more questions can be asked, and the more answers found, further increasing the available know-how to solve emergent problems;
3. The human capacity for abstraction allows increasing numbers of conceptual dimensions, questions, and functional domains to be conceptually and hierarchically linked, thus structuring and increasing the connectivity between different domains of knowledge and understanding;
4. This leads to a continual increase in the density of identified conceptual dimensions in the cognized ‘problem space’ of the individuals involved, and thus gives those individuals an immediate edge over others, as well as over their non-human environment.
5. In the longer term, each solution brings with it its own unexpected challenges, requiring more problem-solving, and a more costly conceptual and material infrastructure to survive.

Ultimately, human survival in the early stages was because no dependencies emerged between the human species as a whole on the one hand, and any specific set of environmental conditions on the other (even though individual groups did depend on particular environments). Human beings were, in the true sense, omnivorous,

⁸ Their capacity to process information is genetically encoded, but the information they process, and the ways in which they do so, is not. It is socio-culturally and self-referentially developed and maintained.

living in the widest possible range of environments, and using in each environment the widest range of available resources. They invested little or nothing in their environment, and took from it whatever they could use. Their social organization, in small groups, put minimal pressure on the environment, allowed a huge range of different ways to survive, and operated with minimal overhead and maximum spread of risk. In short, and in systems terms, the species survived because the coupling between human groups and their environment was extremely loose.

Both the slow start and the subsequent acceleration in the innovation process are, in our opinion, best explained by looking at invention and innovation as a process occurring at the level of the group, rather than the individual. It seems reasonable to assume that there are in every group a number of inventive individuals, and that there is, therefore, an average 'invention rate.' For much of prehistory, however, population densities were very low, and encounters incidental, so that the transmission of inventions was irregular and its rate of loss was high. For innovation to take off, the interactivity among individuals had to exceed a certain threshold, attained by bringing more people together in the same place for longer periods. The circumscription of mobility and/or an overall increase in population were therefore a necessary part of the process leading to increased invention and innovation levels. To enable either, it was necessary to have sufficiently dependable and storable foodstuffs, and that, in turn, required certain innovations.

Once the conquest of the material world was made possible by the invention of conceptual tools such as described above, the coupling between humans and their environments became much tighter, initiating a true co-evolution between the two. That coupling increased investment in specific environments and subsistence strategies and concomitantly increased the risks involved in any individual survival strategy. The problems that emerged prompted a search for solutions, leading to more problems, etc. A control loop emerged between innovation and population density growth that was responsible for an exponential increase in both, over the last 10,000 years. That control loop is summarized in the following box (cf. van der Leeuw & McGlade, 1993; van der Leeuw & Aschan-Leygonie, 2005):

Problem-solving structures knowledge → more knowledge increases the information processing capacity → that in turn allows the cognition of new problems → creates new knowledge → knowledge creation involves more and more people in processing information → increases the size of the group involved and its degree of aggregation → creates more problems → increases need for problem-solving → problem-solving structures more knowledge . . . etc.

The result of this loop is the continued accumulation of knowledge, and, thus, of information-processing capacity, enabling a concomitant increase in matter, energy and information flows through the society, and, therefore, the growth of the number

of people participating in that society. We will discuss the relationship between these flows in the next part of this paper.

3.6 The Emergence of Complex Societies

In archaeology and anthropology, although the boundaries between them are not very sharply defined, we distinguish between the small-scale ‘gatherer-hunter-fisher societies;’ the larger, but still relatively homogeneous ‘tribal societies’ in which there are few structuring institutions based on any form of power or control, simply because the size of the societies (hundreds to a few thousand people) does not require such institutions; and the so-called ‘complex societies,’ which encompass tens of thousands of people or more, and in which control becomes a problem that is solved by the creation of ad-hoc institutions.⁹

On this topic, our ideas run somewhat counter to established, energy-based, theories that argue that in order to establish such societies, the first requirement was to institute subsistence strategies that could yield a food surplus, so that those ‘in power’ would not have to provide their own subsistence, and could harness some of the population at least part of the time to invest in collective works, etc. Instead, in our opinion, the emergence of such societies is linked closely to the problem-solving control loop we have just discussed.

We would argue that, because matter and energy are subject to the laws of conservation, they could never have played the role of driver in the emergence of complex societies. Matter and energy can be transmitted and stored, they may feed people and provide them with other necessary means of survival, but they cannot be shared. They can only be passed on from person to person, and any fortuitous constellation of people that only processes energy and matter together will therefore immediately lose structure. In other words, flows of matter and energy alone could never have created durable human social institutions, let alone complex societies. Information, on the other hand, is not subject to the laws of conservation, and *can* be shared. It follows that the coherence of human societies is due to the exchange of information, which, by spreading similar ideas, links more and more individuals into a network of shared meanings. In effect, human societies are held together by expectations, by institutions, by world-views, by ideas, by technical expertise, by a shared culture! The larger and more complex the society, the more information is processed by its members, requiring an ever more sophisticated information-processing apparatus if the members of the society are to act in concert.

⁹ In practice, it turns out to be difficult to assign precise limits to these different categories, and there is, therefore, an important debate in the discipline about whether these categories are ‘real’ (e.g., Feinman & Neitzel, 1984). We have chosen to maintain them for simplicity’s sake because (a) we are not involved with the detailed distinctions between these categories, and, (b) we are not, here, trying to determine the position on this scale of any individual society. We use the terms in an indicative manner, to point out different parts of a continuum.

This does of course not mean that there was no need for any surplus to provide food and other resources for those people who were executing tasks for the group as a whole. Surplus was necessary, but the ability to procure a surplus was not sufficient for the emergence of complex societies. For such societies to emerge, the people involved must have been sufficiently interactive over considerable periods of time to understand each other, divide tasks and duties, and in general terms organize their society in such a way that everyone has a role, that the society has a stable subsistence and resource base, that there are institutions in place to deal with potential or actual, internal and external conflict, etc.

The multiplication of the number of people interacting together, and of the different kinds of tasks to be fulfilled in various contexts, very quickly did place considerable strain on the communication channels of such groups because the information load increased somewhere near geometrically, or even exponentially, when the number of people increased arithmetically (Mayhew & Levinger, 1976, 1977; Johnson, 1978, 1982, 1983).

With the technical means available, the only way to solve such problems involved reducing the time necessary for communication, making communications more efficient by eliminating error and noise, and inventing ways to communicate by other than oral means. That process must therefore have been going on almost everywhere in complex societies. We see the tangible results in the emergence of increasingly large interactive groups, which occurred in all complex societies in one form or another, once certain thresholds of population size were reached (van der Leeuw 1981). The emergence of towns and cities is among the most prominent consequences of this trend.

3.7 The Emergence of Towns

How did towns and cities emerge? The topic has been studied in many different archaeological contexts (China, the Indus valley, Mesopotamia, Crete, Etruria, Yucatan, etc.). The principal conclusions of such studies are that the emergence of towns is not related to any particular environment, as they emerge in many different ones, deserts, temperate plains, and tropical jungles among them. All one can say is that in the earliest cases, these environments did not facilitate the communication between large numbers of dispersed villages.

It is noteworthy in all the different archaeological contexts in which towns and cities emerged, – China, the Indus valley, Mesopotamia, Crete, Etruria, Yucatan – they never emerged singly. Rather, they always constituted clusters that differentiated themselves from their rural environments at around the same time. This reflects the fact that they are, in effect, always part of a network of exchanges and communications that links their immediate (rural) hinterland with other towns (and their hinterlands) farther away. It is thus not surprising that archaeologists always find trade goods in early cities – goods that sometimes come from hundreds or thousands of kilometers away. Moreover, if one corrects for the particularities of local geographical circumstances, such networks have a very specific spatial structure

that is linked to spatiotemporal constraints in the transport of goods, energy and information both locally (from town to hinterland and *vice-versa*) and among the towns themselves (Reynaud, 1841; Christaller, 1933; Berry, 1967 and many others; for a summary, see Abler, Adams, & Gould, 1997).

Similarly, there appear to be recurrent regularities in the relation between the position of a town in the urban hierarchy and its size (Crumley, 1976; Johnson, 1981; Paynter, 1982, and many others). These regularities generally are explained by the fact that for a complex society to operate coherently, a number of administrative, commercial and other functions need to be fulfilled for all members of the society. Some of these functions are invoked only rarely, but others much more frequently. In order to optimize the overall effort involved in meeting these needs, the frequently invoked functions are present in every town, the rarely invoked ones only in one town, and, for those in between, the number of towns where they are present is dependent on the frequency with which they are needed. As a result, it is argued, town size and town rank (in the hierarchy) scale according to a Pareto or Zipf distribution.

There do not seem to be any external causes for the emergence of towns and cities that one could point at, as they emerge at different times in different regions, and do not emerge at times of particular climatic or other environmental stresses. That has led many urbanists, as well as some archaeologists, to hypothesize that towns and cities emerge spontaneously, due to a process of auto-organization (e.g., Pumain, Sanders, & Saint-Julien 1988; Durand-Dastès et al., 1998). One possible scenario has been elaborated by van der Leeuw and McGlade (1993, 1997).

How did towns and cities affect communication? Van der Leeuw and McGlade (1993, 1997) present us with a detailed discussion of this question in abstract, dynamic terms. We will here summarize it very briefly, and in terms that are more accessible. First, towns concentrated people in space, thereby reducing the time needed to access most information, especially when, within such towns, people were exercising similar activities, and, therefore were most closely involved with each other in everyday life. Secondly, the towns soon became foci of attention for those in the society who were not living there – as marketplaces, they became an important source of goods and information for the surrounding countryside. In the process, trading tokens and, eventually, money were invented as means to communicate and store value and to facilitate material exchanges. Thirdly, in all urban societies we see the development of writing (or some similar means of accounting and communication, such as *quipu's* in Peru), which, on the one hand reduced the error rate in communication, and, on the other, enabled communication by non-oral means. Fourth, in such early towns we see the development of an administration – i.e. institutionalized channels of communication and conflict resolution (e.g., Wright, 1969).

Ultimately, the conjunction between the absence of external drivers towards urbanization on the one hand, and the fact that towns and cities facilitate communication in a major way on the other, convinced us to ascribe the emergence of urban systems to yet another nexus in the development of information processing and communication networks in human societies. But as we have seen above, this interpretation challenges a considerable body of extant theory that ascribes the

emergence of towns in terms of economies of scale in providing subsistence and other resources for the populations concerned. We would therefore like to devote some space to countering these arguments.

3.7.1 The Role of Energy in the Dynamics of Complex, Urban Societies

As biological organisms, individual human beings require that they continually dispose of sufficient energy and matter to stay alive. According to biologists' calculations, that takes about 100 Watt per person. Yet in the developed world today, the average energy consumption per person is of the order of 10,000 Watt, two orders of magnitude higher (IEA 2006). With what we know about the subsistence and lifestyle of most hunter-gatherer people, it seems highly improbable that this increase occurred before the Neolithic. If it began at that time, as it may in some societies that were blessed with sure and plentiful resources, it must still have begun very, very slowly, because there was nowhere to spend or invest that energy. Human exertion may increase the total energy intake of a society somewhat, but could not possibly be responsible alone for a hundredfold increase. Such an increase is only imaginable in the context of a substantial increase in infrastructure of the kind that uses large amounts of wood or fossil fuel, animal energy, or water or wind energy. Those conditions did not come into existence until the emergence of complex, urban societies, several millennia later (around 7,000 years ago). It has therefore been argued that the emergence of towns is in fact driven by economies of scale in energy procurement, transport, and use (Bettencourt, Lobo, Helbing, Kühnert and West, 2007).

In our opinion, this is not so. Assuming an average yield in energetic resources per unit of surface, growth of the urban population would have meant that foodstuffs and other resources would have had to come from further and further away, rapidly leading to important increases in the cost of transportation of these resources. We would therefore argue that the transition towards urbanization is not driven by economies of scale in matter and energy provision, but, that said transition is very costly in energy terms. *Rather than a driver, energy usually must have been a constraint that limited urbanization and the growth of complex societies.* The need to ensure that enough energy and matter reached every member of an urban society must have pushed people towards attempts to solve problems of energy acquisition, distribution, and use.

Not much could be done in the domain of energy *acquisition*. Until the agricultural revolution, the available subsistence resources do not change much: essentially, they are products of agriculture, forestry, hunting and fishing, the breeding of animals and the collection of plants. Until the introduction of fossil or artificial fertilizer in the 19th century, all are essentially dependent on solar energy, and their maximum yield per acre is therefore limited. The available forms of energy also remained essentially the same until the industrial revolution. They included human and animal labor, hydraulic, wind and solar energy, wood, and (to a very limited

degree), coal. The cost of acquisition and the yield of these energy sources did not change enough for any society to dramatically increase its *per capita* food or energy resources (except by appropriating resources accumulated over time by others, such as in the case of conquest).

In the *distribution* of these resources, on the other hand, there was room for energy savings. In contrast with biological organisms, the channels for the flows of energy and matter are not ‘hard-wired’ in human societies. Therefore, sharing information creates the channels through which matter, energy and information are processed and distributed, whether these channels are material (e.g., roads, cables etc.) or remain virtual (e.g., exchange networks such as the *kula* (Malinowski, 1922)). There is great flexibility both in the organization of the networks and the forms in which matter and energy are distributed. Therefore, between the Neolithic and the Industrial Revolution, many improvements in the efficiency of matter and energy transport (such as the introduction of slaves, beasts of burden, wheeled carts, boats, roads, etc.) helped alleviate the energy constraint on urbanization.

Other savings were made by *reducing energy use*, for example by improving crops, reducing crop losses and ameliorating agricultural techniques, redesigning clothing, inserting glass in windows, improving pottery kilns and fireplaces, inventing more efficient tools (levers, pulleys, etc.) to assist in building, and so forth. Finally, as we have seen above, major energy savings were achieved by minimizing the cost of information acquisition and transfer. The organization of regular markets, for example, reduced the time spent in finding appropriate items or information; the introduction of coins and money facilitated the exchange and transmission of value; and the invention of bookkeeping and writing reduced the cost of long-distance communication of information (and increased its efficiency).

All of these savings, though, were not enough to facilitate the growth of truly large cities, such as Rome. This growth required the invention of new techniques to *harness more energy*. Many of these techniques were essentially of a social nature, serving to enhance the control of few over many: feudalism, slavery, serfdom, wage labor, taxation, administration, and so forth. However, as Tainter argues (1988), all the complex societies based on these kinds of harnessing techniques were in the end not sustainable in the absence of fossil energy. The counterpart of this argument is visible with our contemporary eyes: since the industrial revolution, and in particular since the introduction of fossil fuel, cities seem to grow almost exponentially, and no limit seems in sight.¹⁰

In summary, the observation that the growth of towns goes hand in hand with economies of scale in energy use is correct, but the conclusion drawn from it is not: *these economies enabled urban growth by alleviating the energy constraints, but they did not drive it. Quite the reverse, energy savings were forced on urban societies to meet the growth of societies that was driven by economies in communication and information processing.* Towns emerged and grew as a way to deal with the fact

¹⁰ Half the world’s population of c. 6 billion people currently lives in towns and cities, and 80% is expected to live there in twenty or thirty years!

that more and more people were involved in a society's problem solving, leading to increasing diversification and specialization, and therefore to an increasing dependency on frequent interaction and communication. Because town size fundamentally was limited by the need to distribute energy to all inhabitants, much inventiveness was invested in finding ways to do more with less energy. Once fossil energy was domesticated, social systems could, and did, grow without constraint, to the point that we now use 10,000 watts per person on average. Of these, 9,900 watts are invested in our society's infrastructures, and only 100 in our own survival! Finally, it is noteworthy that where it took human beings 200,000 years or more to remove the conceptual constraints on dealing with matter, they removed the energetic constraints in less than 8,000 years. That is in itself a remarkable acceleration, to which we will return.

3.7.2 Complex Societies as Webs of Networks

In what follows, we have chosen to represent human societies as organizations whose existence is dependent upon flows of matter, energy and information that meet the needs of the individual participants by distributing resources throughout the society. Material and energetic resources are identified in the natural environment, transformed by human knowledge such that they are suitable for use in the society, and again transformed during use into forms with higher entropy. These forms can then be recycled, or excreted by the society. The first kind of transformation increases the information content of the resources, whereas the second one reduces their information content. Indirectly, therefore, the information content (or information value) is a measure of the extent to which the resource has been made compatible with the role it fulfils in the society.

Channels for the distribution of energy, matter and information link all individuals in a society through one or more networks. In the smallest of societies, there is essentially one, kinship-based network. Kin relations determine social contexts and exchanges of genes, information, food and other commodities, etc. In complex societies, on the other hand, the networks are many, and are functionally differentiated: kinship, friendship, business relations, and clubs do all constitute social networks. But in such societies, we also have networks of different kinds, such as distribution networks for communication, gas, electricity, fuel, ice, food, etc. Sometimes the channels for information and matter and/or energy are the same, but that is not necessarily so: electricity, petroleum and coal are transported, processed and delivered in different ways. The same is true of virtually all goods in everyday life that we do not collect or process ourselves. We conclude that the 'fabric of society' consists of flows through multiple networks, held together by different kinds of (information) relations, and transmitting different combinations of the three basic commodities (energy, matter and information). In Chapter 5, White introduces the concept of 'multi-net' to refer to the sum of these networks.

How did these all-important networks emerge? It follows from our basic premise that they emerged through a continued exchange of information, matter and energy

that eventually allowed different people at least partially to share perspectives, ways of doing things, beliefs, material culture, etc. The recursive communication underlying this process both facilitated shared understanding among individuals, and drew more and more individuals into a network in which they could communicate and share more easily, and with less risk of misunderstanding, than they would experience with non-members of the network. There was thus a decided adaptive advantage to being part of such a network. Moreover, when the recursive communication remained below a certain threshold, it kept people out of the network because they could not sufficiently maintain their alignment.

What is the nature of such networks? In a network, nodes (actors, towns, hubs) are linked by edges (links) of different kinds (reciprocal, non-reciprocal, symmetric or asymmetric, etc.). We can define a network between any number of components, at any scale, linking any two phenomena (such as people to objects, to functions, to ideas) or serving the transmission of any conceivable commodity. Thus, there are networks of scientists, networks of administrators, networks of pipelines or cables, or roads, etc., but also networks of ideas, principles, artifacts, etc. For each network, both the relations between nodes and the nodes themselves must be defined *ad hoc*. In the case of self-structuring networks such as most social networks in society, one would also have to elicit or define what the thresholds are for sufficiently intensive participation that nodes may be included among the membership.

Like everything else, the nature of these networks can, and does, change. In most complex societies, not only do the actors in the network change, but so does the function of the network. Let us illustrate this with reference to the link between a potential resource, such as a vein of a particular kind of ore, and a society. As soon as a prospector (who has a certain kind of knowledge) identifies a promising geological formation, he or she creates a link between herself or himself and the potential resource. That link is purely informational. Once (s)he gets his (her) claim registered, and searches for funds to start exploiting the vein, there comes into being an embryonic informational (financial, legal, institutional) network that links the vein to other members of the society. In the next phase, that of preparing and beginning mining, a material network is instantiated by linking the mine to an existing road (or rail) network, and, subsequently, to electricity and/or telecommunication networks. From the start, however, the operator will hire mining personnel, tapping into a whole set of new networks (kin and business relations between the workers). Provided the mine is successful and the product can be sold, the mine will be linked into an industrial network that transforms the raw material into a number of finished products. That, however, presupposes that the society itself has not only identified the ore as a potential resource, but has been structured to use it as such, by the emergence of an industrial chain that can transform the resource into something considered valuable to the society. In a short time, the mine has become a node in a large number of functionally different networks that integrate it into the existing society.

In the case of an invention, the process is very similar. It also begins with a single person and a potential commodity or artifact, and links both into the society. But the process is much slower, as it entails, first, the spread of the underlying ideas

into the society (“this box serves to telephone . . . it seems useful . . . I wish to use one”), so that the potential resource is recognized as an actual resource, and then the creation of the appropriate support networks (the communications towers, the dealers, the salesmen, the clients), etc. In either (and any other) case, though, the network is initially one of ideas, and, subsequently, may become material, energetic, communicative, or all three

The configuration of the network is closely related to its dynamics. In recent years, there have been a number of impressive studies looking at the relationship between these two aspects of the networks. In Chapter 5, White presents ways to formalize the interactions between nodes, and, on that basis, show how to define the structure and dynamics of large-scale urban networks.

3.8 Society-Environment Dynamics

We have seen that a society uses information processing to ensure that the necessary matter and energy reach all its members.¹¹ The matter and energy are found in the environment, while the information processing is found in the society. There is thus a control loop with matter and energy as input into the society, information as output (van der Leeuw, 2007). Maintaining a society’s growth requires a continued increase in the quantity of energy and matter flowing through the society. Such growth is achieved through the identification, appropriation and exploitation of more and more resources. At the most abstract level, therefore, the flow of information (structuration) into the environment enables the society to extract from that environment the matter and energy it needs to ensure the survival of its members. The dynamic is driven by the information-processing control loop that aligns more and more people into a connected set of social networks, thus at once increasing the degree of structuration of the society and the number of people involved.

In order for the whole to function correctly, the rates of information processing and those of processing energy and matter need to be commensurate. If not enough information is processed collectively, the society loses coherence; people will act in their own immediate interests and the synergies inherent in collaboration will be lost. Thus, for a center to maintain its power, it must ensure that there is an information-processing gradient outward from itself to the periphery of its territory. At every point in the trajectory between the center and the periphery, it must be more advantageous for the people involved to align themselves with what is happening closer to the center than with what is happening locally or further away.

Over time, such a gradient can only be driven by a continual stream of innovations emanating from the center towards the periphery. Such innovation is facilitated by the fact that, the closer one is to the center, the higher the density of aligned individuals, and thus the more rapid is the information processing. One could say that the innovation density of such a system is thus always higher nearer the center.

¹¹ This is, of course, an idealized situation. Societies may have pathologies where these resources do not reach all of its members, unintentionally or deliberately (Rappaport, 1971).

Innovations create value for those for whom they represent something desirable, but unattainable. The farther one is from the center of the system, the more unattainable the innovations are (because one is farther from the know-how that created the value). In general, therefore, the value gradient is inversely proportional to the information–processing gradient. Value in turn attracts raw materials and resources from the periphery. These raw materials and resources are transformed into (objects of) value wherever the (innovative) know-how to do so has spread, thus closing the loop between the two gradients. The objects are then exchanged with whoever considers them of value, i.e. whoever cannot make them (or make them as well, or as efficiently).

3.9 An Example: The Expansion of Ancient Rome

To illustrate how this works in practice, we could look at the history of the Roman Empire (van der Leeuw & De Vries, 2002). The expansion of the Roman republic was enabled by the fact that, for centuries, Greco-Roman culture had spread around the Mediterranean coasts. It had, in effect, structured the societies in (modern) Italy, France, Spain and elsewhere in a major way, leading to inventions (such as money, the use of new crops, the plough), the building of infrastructure (towns, roads, aqueducts), the creation of administrative institutions, and the collection of wealth. Through an ingenious policy, the Romans aligned all these and made them subservient to their needs, i.e. to the uninterrupted growth of the flows of matter (wealth, raw materials, foodstuffs) and energy (slaves) throughout their territories, toward their urban centers and ultimately toward Rome itself.

The Roman Republic and the Empire could expand as long as there were more pre-organized societies that could be conquered. Once their armies reached beyond the (pre-structured) Mediterranean sphere of acculturation (i.e. when they came to the Rhine, the Danube, the North African deserts and the Middle Eastern Empires), that was no longer the case, and conquests stopped. Then they began major investment in the territory thus conquered. This investment consisted of expanding the infrastructure (highways, *villae*, industries) and the trade sphere (Baltic, Scotland, but Roman trade goods have been found as far a-field as India and Indonesia) to harness more resources. As large territories were thus ‘Romanized,’ these territories became less dependent on Rome’s innovations for their wealth, and thus expected less from the Empire. One might say that the ‘information gradient’ between the center and the periphery leveled out, and this made it increasingly difficult to ensure that the necessary flows of matter and energy reached the core of the Empire. As the cost of maintaining these flows grew (in terms of maintaining a military and an administrative establishment, for example: see Tainter, 2000; Crumley & Tainter, 2007), the coherence of the Empire decreased to such an extent that it ceased, for all intents and purposes, to exist. People began to focus on their own interest and local environment rather than their interest in maintaining a central system. Other, smaller, structures emerged at its edges, and there the same process of extension from a core began anew, at a much smaller scale, and based on very different kinds

of information. In other words, the alignment between different parts of the overall system broke down, and new alignments emerged that were only relevant locally or regionally.

3.10 Networks of Cities Constitute the Skeleton of Large-Scale Societies

In such an overall flow-structure dynamic, cities play a major part. They are demographic centers, administrative centers, foci of road systems, but above all, they are the nodes in the network, the locations where most information processing goes on. As such, they are the backbone of any large-scale human social system. They operate in network-based 'urban systems' which link all of them within a particular sphere of influence. Such systems have structural properties that derive from the relative position the cities occupy on the information-processing gradient, and in the communications and exchange networks that link them to each other (Chapter 5). Although the role of individual towns in such systems may change (relatively) rapidly (Guerin-Pace, 1993), the overall dynamic structures are rather stable over long periods.

Because people congregate in cities, the latter harness the densest and the most diverse information processing capacity. Not only does this relatively high information-processing capacity ensure that they are able to maintain control over the channels through which goods and people flow on a daily basis, but their cultural (and, thus, information-processing) diversity also makes them into preferred loci of invention and innovation. Innovation (as represented by the number of people involved in research, the number of research organizations, the number of patents submitted, etc.) scales super-linearly with the size of urban agglomerations (see Chapters 7 and 8).

The super-linear scaling of innovation with city size enables cities to ensure the long-term maintenance of the information gradient that structures the whole system. This is due to a positive control loop between two of any city's roles. On the one hand, most flows of goods and people go through towns and cities. That supplies them with information about what is happening elsewhere and this again enhances their potential for invention and innovation. But the same connections enable them to export these innovations most effectively – exchanging some part of their information processing superiority for material wealth.

3.11 The Role of Cities in Invention and Innovation

In the last section, we alluded to the fact that cities are preferred loci of invention and innovation, but did not really elaborate. Although this topic will be dealt with in a more extensive way in other places in this volume (Chapters 6–8, 12), in order to round off our outline of the long-term evolution and role of human conceptual

systems from Early Man to the present, we will devote a couple of pages to the role of towns and cities in generating innovations.

First, we must clear the terminological ground a bit, and look somewhat closer at how inventions become innovations. *Inventions* are essentially local phenomena in the social and information-processing network that constitutes a society: they involve one or a few people, and one or a few ideas. *Innovations*, on the other hand, involve the network as a whole, they are global phenomena, as they imply the spread of the invention to all potentially relevant parts of the network, and, while spreading, they expose the invention to many other forms of knowledge and ideas (which we will here call conceptual dimensions). These conceptual dimensions are themselves linked to possibilities and challenges that the original inventors were not aware of, or did not connect in any way to their invention.

3.11.1 When are Inventions Transformed into Innovations?

The question we wish to address in this section is, therefore: *When are inventions transformed into innovations?* To investigate it, we will contrast unsuccessful inventions with successful ones. There are a number of ‘classical’ cases of inventions that were instantiated, but never transformed into innovations until very much later. Iron-working, for example, was first invented about a thousand years before it appeared as a common technique to make a wide range of weapons and tools (Stig Sorensen & Thomas, 1989; Collis, 1997), and Hero of Alexandria’s steam engine in the first century BC was not used widely until it was reinvented 1,600 years later. In both cases, the conceptual and material tools and techniques were available to instantiate these inventions, but the societal and/or technological context was not such that the invention could spread. In the former case, the society was very hierarchically organized, and those in power (who controlled the bronze industry by controlling the sources of bronze) were not ready to allow a technology to emerge that they could not control (because iron was found, literally, in every bog or riverbed). In the latter case, in a society based on slavery, there was no demand for steam power . . .

The contrast with the present is striking. The closer we come to modern times, the more clearly we can observe major innovations coming in waves that rapidly succeed each other. In such a wave, once an initial invention is transformed into an innovation (i.e. when the invention has become popular, changing the way people do things and think about them), this triggers a cascade of other, related, inventions/innovations so that, together, these innovations completely change one or more domains of daily life, trade and/or industry. Think of the introduction of printing, or, more recently, the development of the computing and biotechnology industries – with nanotechnology already on the horizon.

What makes the difference? Spratt (1989) summarized what it takes to get a relatively complex innovation, such as a car, ‘up and running.’ He outlines some of the many inventions and innovations that had to exist, and to be linked together in a single conception, before the first car came on the road. These go back several centuries, and his example shows very clearly that the emergence of an invention

is highly context-dependent. Many of the contributing innovations were made in domains that had nothing to do with transport (such as the discovery of rubber and the invention of new manufacturing techniques for steel), and clearly were not driven by demand in the transport sector. Others, however, were indeed triggered when problems emerged that had to be solved for the car to work. In fact, his example shows beautifully that an innovation of such complexity as a car is not possible until a certain (high) density in cognized and conceptualized problem-solving tools has been introduced and instantiated in a society as a whole.

In other words, an invention can successfully be transformed into an innovation, and then trigger other inventions and innovations when there is a sufficient density of relevant conceptual dimensions in the ‘global’ network. These dimensions can be conscious or latent (waiting to be discovered), on the one hand to instantiate the invention (e.g., the necessary raw materials, tools, techniques to make it) and on the other to link the invention to a range of new functional domains to be explored and/or exploited. In the absence of such a sufficient density, for whatever reasons, the invention may not take off at all, or may remain alone without triggering a cascade. Both densities are, of course, closely related to the available density of connected ‘grey matter’ involved in information–processing.

3.11.2 How are Inventions Transformed into Innovations?

Now let us tackle a more difficult question: *how is an invention transformed into an innovation?* According to current cost-benefit theory, that depends on whether there is (latent or conscious) demand for the invention, or, if there is not, whether demand can be generated within a relatively short time-span. In that theoretical context, there is an important distinction between the distant past and the recent period. We argue that in the distant past, once an invention was available, its transformation (or not) into an innovation was demand-led, but that in the present, many inventions are made into innovations at great cost, by advertising, by the creation of scaffolding structures, etc. Innovation in the modern world is deemed increasingly supply-led.

Although that distinction makes an important point, we wish to make clear that we do not believe that innovations are either completely demand-led, or completely supply-driven. After all, what is at stake is a match between an offer (an invention) and a demand (a need in the society), or between a solution and a problem: does the problem trigger the solution or does the solution make one aware of the problem? There is always a bit of both in the transformation of an invention into an innovation. The differences are, in our opinion, a matter of proportions. The important thing is that a match is made that is of sufficient relevance to the society to adopt the invention and generalize it, so that it may pose new problems and trigger new solutions (and *vice versa*).

However, there are some consequences for the frequency and structure of innovating. If demand were the limiting factor, generally, there would first seem to be an important percentage of inventions that never are transformed into innovations, simply because they are forgotten before anyone outside the circle of the inventor(s)

notices them. The loss of inventions is probably more important than is the case when there is a deliberate policy of innovation (though it would be difficult to demonstrate for lack of information on demand-led innovations). Secondly, it would seem that in the supply-led case there is, at least to some extent, a 'logic of innovation' that it may be possible to retrace, whereas in the demand-led case that would seem absent because the process has too many degrees of freedom and resembles random walk. As a result, when the emphasis is on demand, it would seem that the overall rate of innovation is slower than when there is a strong supply-driven effort at innovation. That may also explain some of the acceleration in innovation that we have seen over the last couple of centuries.

But the very tentative way in which we have put these ideas testifies to the fact that, in our opinion, demand-led innovation merits much more study from the kind of perspective we have tried to open here. There is, as of yet, too much study of *either* innovation *or* the converse, tradition, and not enough of the *relationship between the two* (van der Leeuw, 1994). Yet it is in the interaction between what exists at any time and what is invented, that the emergence of innovations has to be explained. It is for that reason that we (Ferrari, Read, & van der Leeuw, Chapter 14, this volume) have tried to create a simulation model of the interaction between invention and the formation of a consensus about it, which would be a first step towards transformation of the invention into innovation.

3.11.3 What Is the Role of Cities In the Transformation of Inventions Into Innovations?

Getting back now to the relationship between cities and innovations, why, then, are cities essential to invention and innovation and *vice-versa*? In recent years, several important characteristics of cities have emerged that are of relevance here: high population density, demographic diversity, above-average interactivity among the city's inhabitants, and accessibility from outside. We will briefly deal with each in turn.

The first of these, high population density, is critical, as described in detail in Chapters 6–8. One implication of high population density is a rich interaction structure, which increases the speed with which innovations can be propagated.

The high interaction level has a counterpart in the diversity characteristics of urban populations, particularly with respect to specialized competences and functionality. A study undertaken by ISCOM team members for 5,500 settlements in the southern half of France demonstrated very clearly that the resilience and growth of individual cities are, among other things, directly related to diversity in age groups, diversity in economic activities, and diversity in level of education of the population (ARCHAEOMEDES, 1999). We explain this by pointing out that the more diverse a population, the more encompassing is its potential 'possibility space,' i.e., the total set of potential domains and directions in which the town can grow through innovation. Thus, if an invention is made or brought into an urban context (either

because the inventor, as often happens, moves to a city or because someone picks up on an idea), suddenly the problem space with which the invention is confronted is much wider, and the chances are increased that the invention triggers, or meets, more problems for which it can provide a solution . . . and that in turn enhances the possibility that the invention triggers a cascade.

In addition, high levels of interurban interaction within an urban system increase the innovative capabilities of these systems (Chapter 6). This increase points clearly to the fact that both the increased access to ideas and resources, and the capacity to favor the spread of innovations that come with these interactions are also important elements in the equation. It is that capacity which has, initially slowly but increasingly with explosive rapidity, led to what Ingold (1987) calls '*The Appropriation of Nature:*' the reduction of the complexity and diversity (biodiversity, spatial diversity) of our natural environment and the increase of human control over it.

3.12 Do We Need a Lesson From the Past?

We are aware that this paper has only begun to identify where we can, and must, scratch the surface on invention, innovation, and the history of our species that led us from life in tiny groups without man-made shelter and only a few stone tools, to a worldwide society of more than six billion people who live in an infinitely complex social and material culture, and in an environment that borders on the artificial. Moreover, our map of the places to be scratched is very fragmentary in coverage, and only has touched on a few of the scales at which it approaches different problems. Yet, we do believe that we can sketch some of the implications of our work for the future, and we do not want to end the paper without sharing some of them with the reader.

It took our ancestors of different subspecies hundreds of thousands of years to establish the conceptual tools to deal with matter, and thousands of years to do the same with energy. We currently are in the first years of the third of these revolutions, the information revolution, which will, by extrapolation, last a number of decades or even centuries.

The driving force behind these developments has been the interaction between problems and solutions, in which known problems beget solutions beget unknown problems. With each new invention, new conceptual dimensions were added to the existing arsenal, and the total problem/possibility space now counts an almost infinite number of dimensions. That in turn has exponentially increased the potential problems (or to use the softer term 'unforeseen consequences') that any innovation can trigger. As in the case of Rome, when the possibility space expands through inventions and innovations, the problem space expands even faster. In the end, therefore, the rate at which problems emerge overcomes the rate at which people can innovate to solve them, which causes crises and re-structurations (van der Leeuw, 2007).

Thus far, these crises and the need to restructure from the bottom up (as in the case of Rome, but also the Chinese, the Maya, the Indus, and other civilizations)

have limited the overall speed of change in human societies. However, at present, we are on the threshold, for the first time, of innovations that depend on, and enable, reflexive intervention in our own systems. These innovations, moreover, occur at a scale and speed, and are of such a complexity, that the intuitive human apparatus to deal with new problems by using models from other domains may quickly become obsolete.

On the one hand, that may pose enormous dangers, for societies that completely get cut off from the traditions that have enabled them to maintain a degree of coherence. On the other hand, that may offer, for the first time, the opportunity to move away from the means by which most societies have survived thus far (i.e., by gaining control over, and destroying large parts of our environment). The challenge is unimaginable in scope.

We had better get used to and begin to deal with this challenge. One starting point might be to gain better knowledge of what drives innovations in our societies, so that rather than deal with the consequences of our innovation drive, we can begin to deal with that drive itself. This book aims to make a beginning with that task, by juxtaposing the acceleration of innovation and urbanization, trying to improve and formalize our descriptions of these twin phenomena, modeling them, and, thus, gaining a deeper insight in what drives them, what constrains them, and what might help us control them.

Acknowledgments The authors wish to acknowledge the pleasure of close collaboration with Denise Pumain, Geoffrey West, Douglas White, José Lobo and all the others participants of the ISCOM project over a period of four years preceding the publication of this paper. They also gratefully acknowledge the funding of the project by the Research Directorate of the European Union as part of its ICT (FET) program, under contract n° ICT-2001-35505.

References

- Abler, R., Adams, J. S., & Gould, P. (1997). *Spatial organization, the geographer's view of the world*. Englewood Cliffs, NJ: Prentice Hall.
- ARCHAEOMEDES. (1999). *ARCHAEOMEDES II: Policy-relevant models of the natural and anthropogenic dynamics of degradation and desertification and their spatiotemporal manifestations*. In S. E. van der Leeuw (Ed.), Unpublished Report to DG XII of the European Union, Vol. 5: A multiscale investigation into the dynamics of land abandonment in Southern France (S. Lardon, L. Sanders, et al.).
- Atlan, H. (1992). Self-organizing networks: weak, strong and intentional. The role of their under-determination. *La Nuova Critica*, 19–20, 51–70.
- Berry, B. J. L. (1967). *Geography of market centers and retail distribution*. Englewood Cliffs NJ: Prentice-Hall.
- Bettencourt, L. M. A., Lobo, J., & Strumsky, D. (2007). Invention in the city: Increasing returns to patenting as a scaling function of metropolitan size. *Research Policy*, 36, 107–120.
- Bettencourt, L. M. A., Lobo, J., Helbing, D., Kühnert, C., & West, G. B. (2007). Growth, innovation, scaling, and the pace of life in cities. *Proceedings of the National Academy of Sciences*, 104(17), 7301–7306.

- Boëda, E. (1990). De la surface au volume: analyse des conceptions des débitages Levallois et laminaires. In C. Farizy (Ed.), *Paléolithique moyen récent et Paléolithique supérieur ancien en Europe. Ruptures et transitions: examen critique des documents archéologiques - Actes du Colloque international de Nemours, 1988. Nemours: Mémoires du Musée de Préhistoire d'Ile de France*, 3 (pp. 63–68).
- Bradley, R., & Edmonds, M. (1993). *Interpreting the axe trade; production and exchange in Neolithic Britain*. Cambridge, UK: Cambridge University Press.
- Christaller, W. (1933). *Die Zentralen Orte in Süddeutschland: eine Ökonomisch-Geographische Untersuchung über die Gesetzmässigkeit der Verbreitung und Entwicklung der Siedlungen mit Städtischen Funktionen*. Jena, Germany: Fischer Verlag.
- Crumley, C. A. (1976). Toward a locational definition of state systems of settlement. *American Anthropologist*, 78, 59–73.
- Crumley, C. A., & Tainter, J. A. (2007). Climate, complexity, and problem solving in the roman empire. In R. Costanza, L. J. Graumlich, & W. Steffen (Eds.), *Sustainability or collapse? An integrated history and future of people on earth* (pp. 61–77), Cambridge, MA: MIT Press.
- Collis, J. (1997). *Iron age Europe* (2nd ed.). London, UK: Routledge.
- Durand-Dastès, F., Favory, F., Fiches, J.-L., Mathian, H., Pumain, D., Raynaud, C., et al. (1998). *Des Oppida aux Métropoles*. Paris, France: Anthropos.
- Feinman, G. A., & Neitzel, J. (1984). Too many types: An overview of sedentary pre-state societies in the Americas. *Advances in Archaeological Method and Theory*, 7, 39–101.
- Guerin-Pace, F. (1993). *Deux siècles de croissance urbaine*. Paris, France: Anthropos.
- Ingold, T. (1987). *The appropriation of nature*. Manchester, UK: Manchester University Press.
- International Energy Agency (IEA). (2006). *Energy Balances of OECD Countries (2006 edition) and Energy Balances of Non-OECD Countries (2006 edition)*. Paris, France: IEA. Available at <http://data.iaea.org/ieastore/default.asp>.
- Johnson, G. A. (1978). Information sources and the development of decision-making organizations. In C. L. Redman, M. J. Berman, E. V. Curtin, W. T. Langhorne Jr., N. M. Versaggi, & J. A. Wander (Eds.), *Social archaeology: Beyond subsistence and dating* (pp. 87–122). New York: Academic Press.
- Johnson, G. A. (1981). Monitoring complex systems integration and boundary phenomena with settlement size data. In S. E. van der Leeuw (Ed.), *Archaeological approaches to the study of complexity* (pp. 143–188), Amsterdam, The Netherlands: Institute for Pre and Protohistory (Cingula VI).
- Johnson, G.A. (1982). Organizational structure and scalar stress. In A. C. Renfrew, M. J. Rowlands, B. Segraves (Eds.), *Theory and explanation in archaeology: The southampton conference* (pp. 389–421), New York: Academic Press.
- Johnson, G. A. (1983). Decision-making organizations and pastoral nomad camp size. *Human Ecology*, 11, 175–200.
- Malinowski, B. (1922). *Argonauts of the Western Pacific: An account of native enterprise and adventure in the archipelagoes of Melanesian New Guinea*. London, UK: Routledge.
- Mayhew, B. H., & Levinger, R. L. (1976). On the emergence of oligarchy in human interactions. *American Journal of Sociology*, 81, 1017–1049.
- Mayhew, B. H., & Levinger, R. L. (1977). Size and density of interaction in human aggregates. *American Journal of Sociology*, 82, 86–110.
- Paynter, R.W., (1982). *Models of spatial inequality: Settlement patterns in historical archaeology*. New York: Academic Press.
- Pigeot, N. (1992). Reflexions sur l'histoire technique de l'Homme: de l'évolution cognitive à l'évolution culturelle. *Paléo*, 3, 167–200.
- Pumain, D., Sanders, L., & Saint-Julien, Th. (1988). *Villes et Auto-Organisation*. Paris, France: Economica.
- Rappaport, R. A. (1971). The sacred in human evolution. *Annual Review of Ecology and Systematics*, 2, 23–44.
- Reynaud, J. (1841). Villes. In *Encyclopédie nouvelle* (Vol. VIII, pp. 670–687). Paris, France: Gosselin.

- Spratt, D. (1989). Innovation theory made plain. In S. E. van der Leeuw & R. Torrence (Eds.), *What's new? A closer look at the process of innovation*. London, UK: Hyman and Unwin.
- Stig Sorensen, M. L., & Thomas, R. (1989). *The Bronze-age iron-age transition*. Oxford, UK: British Archaeological Reports (International Series, 483), 2 vols.
- Tainter, J. A. (1988). *The collapse of complex societies*. Cambridge, UK: Cambridge University Press.
- Tainter, J. A. (2000). Problem solving: complexity, history, sustainability. *Population and Environment*, 22, 3–41.
- van der Leeuw, S. E. (1981). Information flows, flow structures and the explanation of change in human institutions. In S. E. van der Leeuw (Ed.), *Archaeological approaches to the study of complexity* (pp. 230–329), Amsterdam, The Netherlands: University Printing Office (Cingula VI).
- van der Leeuw, S. E. (1994) La dynamique des innovations. *Alliages 21* (Penser la technique): 28–42.
- van der Leeuw, S. E. (2000). Making tools from stone and clay. In T. Murray & A. Anderson (Eds.), *Australian archaeologist. Collected papers in honour of J. Allen* (pp. 69–88), Canberra, Australia: ANU Press.
- van der Leeuw, S. E. (2007). Information processing and its role in the rise of the European world system. In R. Costanza, L. J. Graumlich, & W. Steffen (Eds.), *Sustainability or collapse? An integrated history and future of people on earth* (pp. 213–241). Cambridge, MA: MIT Press (Dahlem Workshop Reports).
- van der Leeuw, S. E., & Aschan-Leygonie, C. (2005). A long-term perspective on resilience in socio-natural systems. In U. Svedin & H. Liliénstrom (Eds.), *Micro–Meso–Macro. Addressing complex systems couplings* (pp. 227–264), London, UK: World Scientific.
- van der Leeuw, S. E., & De Vries, B. M. (2002). Empire: The Romans in the mediterranean. In B. L. de Vries & J. Goudsblom (Eds.), *Mappae Mundi: Humans and their habitats in a long-term socio-ecological perspective* (pp. 209–256). Amsterdam, The Netherlands: Amsterdam University Press.
- van der Leeuw, S. E., & McGlade, J. (1993). Information, Cohérence et Dynamiques urbaines. In B. Lepetit & D. Pumain (Eds.), *Temporalités Urbaines* (pp. 195–245), Paris, France: Anthropos/Economica.
- van der Leeuw, S. E., & McGlade, J. (1997). Structural change and bifurcation in urban evolution: a non-linear dynamical perspective. In J. McGlade & S. E. van der Leeuw (Eds.), *Archaeology: Time, process and structural transformations* (pp. 331–372). London, UK: Routledge.
- Wobst, H. M. (1978). The archaeo-ethnology of hunter-gatherers or the tyranny of the ethnographic record in archaeology. *American Antiquity*, 43(2), 303–309.
- Wright, H. T. (1969). *The administration of rural production in an early Mesopotamian town*. Anthropological Papers, Number 38, Museum of Anthropology, Ann Arbor, MI: University of Michigan.

Chapter 4

Biological Metaphors in Economics: Natural Selection and Competition

Andrea Ginzburg

Once, [Henry] George recalled a conversation a decade earlier between himself and Youmans [promoter of the publication of Spencer's works in America] on the situation of American society. "What do you propose to do about it?" George had asked. To this Youmans responded "with something like a sigh:" "Nothing! You and I can do nothing at all. It's all a matter of evolution. We can only wait for evolution. Perhaps in four or five thousand years evolution may have carried men beyond this state of affairs. But we can do nothing."

(R.C. Bannister, 1979, p. 75).

4.1 Introduction

In our time, we face stagnating production and employment but also, more generally, disappointing results in the provision of services that importantly affect our lives – in education, healthcare and social welfare, in transport and communications. To combat these ills one sole therapy is increasingly invoked: greater competition. Here, I am not referring to the idea that, in well-defined situations from the point of view of structure and of the actors involved, the stimulus from competitors, together with other forces, actions and stimuli, may revive waning energies through a spirit of emulation. A different idea has been asserting itself, with the force of common sense – all the stronger because more abstractly universal¹ and unconditional: namely, the idea that competition is the principal driving mechanism of society, and therefore steps towards increasing its presence in all fields must be the necessary and sufficient condition for economic growth and “social progress” (however defined). The

A. Ginzburg (✉)
Università di Modena e Reggio Emilia, Reggio Emilia, Italy
e-mail: ginzburg@unimore.it

¹ For a similar distinction between “universal” and “particular” determinants, aimed in that case at the limited explanatory capacity of analyses of the innovation processes starting from “diffuse and general” incentives to the reduction of costs, and which I shall recall later on, see Rosenberg (1976).

general nature of the relation between competition and growth (quantitative and qualitative) is seen as rendering unnecessary any “situated” analysis – i.e. based on knowledge of the specific situation – of the possible effects of an increase in pressure of competition. The general nature of the diagnosis thus engenders the general nature of the therapy.

At first glance, it seems rather paradoxical to invoke ever-larger doses of competition at a time when, with globalized markets, the forces of competition appear to have become incomparably more far-reaching and incisive than in any previous period in history. The current popularity of the term “competition” is partly due to its vagueness or ambiguity, hence to the wide range of its meaning. At a more superficial level, the call for greater competition stems from its identification² with the “capacity to face up to it”: thus the tautology that in order to overcome competition, greater competition is needed. At a deeper level, the ambiguity derives from the fact that the history of economic thought offers at least three different conceptions of competition: that of classical political economy (from Smith to Sraffa); the neoclassical idea of equilibrium (associated e.g. with Walras and others); and the neoclassical idea³ of competition as a process (associated with the Austrian school, particularly with Hayek and others). From these three conceptions, as I shall briefly remark later, derive very different relations between competition and growth.

There is a very close relationship between saddling competition with such demanding tasks and the ongoing assertion, as from 1980, of the culture and monetarist policies of deregulation of markets; initially associated with the figures of Ronald Reagan in the USA and Margaret Thatcher in the UK, these have gone hand in hand with the current phase of expansion of the global market. The “salvationary” role ascribed to competition stems from the, very often unconscious, acceptance of some central points in monetarist theory as expounded in the writings of Milton Friedman and, above all, Friedrich von Hayek. However different, in some respects, the ideas of these two exponents of the Chicago school may be (Friedman belonging to the Walrasian rather than the Austrian current), they share a view of the economic process in which the biological metaphor of natural selection looms large. I shall argue that making competition the main driving function of society stems from likening the competitive process to that of natural selection. To what extent is such an assimilation legitimate?

It is widely acknowledged that metaphors have a rightful place not merely in general discourse but especially in scientific discourse, and not simply by way of ornament. Metaphors enable us to use “the inferential structure of one conceptual domain [. . .] to reason about another [. . .] domain” (Lakoff & Nuñez, 2000), or, in other words, to employ what has been learnt on one subject in order to comprehend

² See, for example, Nardozzi (2004): “Where does the decline [of the Italian economy] come from? The clue to understanding this is the same as for the miracle that went before it: competition. . . But . . . unlike then . . . , the increasing international competition has not translated into increased productivity sufficient to improve or maintain the competitive position. What is lacking with respect to half a century ago is not competition itself but the ability to face up to it.”

³ Some authors would locate the Austrian school outside neoclassical theory. This view is hard to reconcile with the acknowledged influence exercised by the writings of Carl Menger and especially, as I shall note, of Hayek.

another subject in a different field. This cognitive device is indispensable if we wish to understand (or represent) something abstract in concrete terms. Moreover, as Rorty (1979) remarks, “It is pictures rather than propositions, metaphors rather than statements which determine most of our philosophical convictions.” It is in the nature of metaphors to propound analogies in different contexts, and thus to make possible creative and stimulating illumination or misleading and dangerous parallels. One may wonder whether, and to what extent, the process of natural selection put forward by Darwin actually finds its counterpart in the “evolutionist” interpretations of the competitive market propounded by the monetarists, especially by Hayek, and ultimately what role evolutionism plays in their “vision” of the social process. It will be seen that while the application of the metaphor of natural selection to the analysis of human society in general, and of the competitive market in particular, is much more problematic than is usually thought, other concepts drawn from evolutionary biology may prove to be far more useful for a “non-panglossian” analysis of the development of the social processes. I shall conclude, however, that, in the case of Friedman and Hayek, though by partly different routes, evolutionism turns out to be associated – through a rhetorical demonstration by analogy – with the attribution of characteristics of optimality (otherwise hard to justify at theoretical level, let alone empirically) to the “spontaneous” evolution of competitive markets.

4.2 Malthus, Darwin, Social Darwinism

The history of the relations between economics and biology is a tortuous one, with many hitherings and thitherings and crossing paths. Suffice it to recall here the influence exerted by Malthus’s *Essay on the Principle of Population*⁴ on Darwin in the process that led to the discovery of the principle of natural selection. In a well-known passage of his autobiography, Darwin (1958) writes:

In October 1838, that is, fifteen months after I had begun my systematic enquiry, I happened to read for amusement Malthus on *Population*, and being well prepared to appreciate the struggle for existence which everywhere goes on from long-continued observation of the habits of animals and plants, it at once struck me that under these circumstances favourable variations would tend to be preserved, and unfavourable ones to be destroyed. The result of this would be the formation of new species.

The passages in *Notebook D*,⁵ where Darwin (between 28 September and 3 October 1838) jotted down his impressions from reading Malthus, give a more precise idea of the stimuli that led Darwin to modify and rearrange his previous conceptions. Up to then, Darwin had concentrated on the struggle among the species. Reading the *Essay on the Principle of Populations*, which deals almost exclusively with human societies, prompted Darwin to focus also on conflict *within* species.

⁴ Darwin’s remarks refer to the 1826 6th edition of Malthus (1798). It is interesting to recall that also A. R. Wallace, who discovered the principle of natural selection independently of Darwin, took inspiration from Malthus.

⁵ For a close reading of these passages, upon which I have amply drawn in this section, see La Vergata (1990a) and Hodge and Kohn (1985).

Moreover, he was struck by the singularly forceful manner in which Malthus put forward “the idea of a war between the species.” Here, for the first time,⁶ the extent and intensity of the struggle is energetically associated with a conflict at once fundamental and inevitable, the fight to appropriate the means of subsistence in a situation in which the growth of population systematically outruns that of the resources. No one prior to Malthus had “seen clearly,” writes Darwin, the existence of a “great check” on the growth of population in human societies. In their relation with nature, human and non-human species are subjected to checking forces (nowadays we should say “negative feedbacks”), partly similar, partly different.

As regards the analogies, Darwin remarks that in man as in other species “even a few years plenty makes population [. . .] increase, and an ordinary crop then causes a dearth in Spring.” The existence of “repressive checks” to reproduction such as famine, cold, epidemics, etc. followed by death and/or a slowdown in reproduction retraces the re-equilibrating mechanism occurring in human societies to that existing in nature.

Malthus argued that in human societies the “repressive checks,” influencing the mortality rate, and hence *ex post*, are mainly “poverty and vice,” the latter term used to describe the illicit practices of infanticide and abortion, especially frequent among the poor. He identified the peculiarity of the human species in two circumstances: the possibility to increase the production of the means of subsistence in proportion to the increase of population, and the faculty of predicting the consequences of one’s own actions. The first of these, however, was effectively cancelled out by the effect of diminishing returns in the production of food, which underpinned the “iron law” by which the means of subsistence grew arithmetically but not geometrically. The faculty of prediction, which remained the sole feature that distinguished the human species, enabled one to choose a “preventive check”⁷ on population growth, based on moral restraint. This “conscious” check entailed delaying marriage and abstaining from the “natural propensities” to procreate until one was able to provide for children. Hence Malthus reached the conclusion that from evil (poverty and vice) good may ensue, i.e. the improvement of man’s character and the increase of his energy.

At this point in his notes, Darwin uses a “creative analogy” – in the phrase of La Vergata (1990a, p. 349) – to establish a surprising parallel between this “moral message” of Malthus and the outcome of the process of selection. The description of the struggle for survival in the presence of an imbalance between population

⁶ The subject of disproportion between demographic pressure and resources has, however, numerous precedents, some of which are recalled by Malthus himself. See La Vergata (1990a), p. 349, note 122.

⁷ The preventive check of moral constriction was added by Malthus after the first edition. Note that a further preventive check was also added, the fear of poverty, i.e. the fear of seeing one’s standard of living fall. This is not a negligible alteration, at least in principle, in Malthus’s deterministic-biological scheme: the check on population no longer involved the mere availability of the means of subsistence, but also the standard of living attained and thus expected/desired by the various social strata. In this way, Malthus sought to answer the objection regarding the lack of a restrictive check on reproduction in the more prosperous classes.

and subsistence, of the “necessary” process of adaptation, and of the outcome that ultimately ensues as a “positive” result of the demographic pressure and the struggle itself, are expressed – in this initial formulation of the theory of natural selection – in the following words:

One may say there is a force like a hundred thousand wedges trying [to] force every kind of adapted structure into the gaps in the economy of Nature, or rather forming gaps by thrusting out weaker ones. The final cause of all this wedging, must be to sort out proper structure, and adapt it to change – to do that for form, which Malthus shows is the final effect (by means however of volition) of this populousness, on the energy of Man.⁸

Malthus’s idea about the existence of universal laws of nature, valid at all times and for all living beings, takes its place in a consistent conceptual framework that we have no grounds for ascribing to Darwin (though some Malthusian suggestions do occur in Darwin’s writings). From the appeal to a law of nature, universal and thus inexorable, that man could not escape, Malthus’s doctrine had drawn “all its ideological and polemical force.”⁹ His aim had been to refute the ideas of Godwin and Condorcet on the perfectibility of man and society and the construction of a system of equality. He thus takes his place among those who opposed the ideas of progress of the French Revolution, an opposition set in motion by Burke’s *Reflections on the Revolution in France*. One corollary that Malthus had drawn from his theses was the proposal to repeal the Poor Laws,¹⁰ in order to dissuade the poor from marrying even when they lacked the means of subsistence and to spur them to a greater “sense of responsibility.” In hypothesizing “laws of nature” valid “ever since we have had any knowledge of mankind” and extending to the animal and vegetable kingdoms, Malthus (1979) also stated that such laws were inescapable, nor could the evils of humanity be in any way ascribed to human institutions. In line with a long tradition (Malthus, 1979),¹¹ he suggested that the “lash of want”, in accord with divine providence, would stir humanity to shake off its congenital idleness and bring out man’s “superior,”¹² moral faculties.

⁸ Cf. La Vergata (1990a), p. 339 and note 106.

⁹ *Ibid.*, p. 349, “Among plants and animals the view of the subject is simple,” Malthus (1979) writes “They are impelled by a powerful instinct to the increase of his species, and this instinct is interrupted by no reasoning or doubts about providing for their offspring. Wherever therefore there is liberty, the power of increase is exerted, and the superabundant effects are repressed afterwards by want of room and nourishment, which is common to animals and plants, and among animals by becoming the prey of others.”

¹⁰ Though he calls this provision “a palliative”, Malthus (cf. *Ibid.*, p. 101) claimed that it would “at any rate give liberty and freedom of action to the peasantry of England, which they can hardly be said to possess at present. They would then be able to settle without interruption, wherever there was a prospect of a greater plenty of work and a higher price for labour. The market of labour would then be free, and those obstacles removed which, as things are now, often for a considerable time prevent the price from rising according to demand”.

¹¹ On the tradition of the “lash of want” and Malthus’s natural theology, cf. La Vergata (1990b), pp. 76 et seq.

¹² It should be recalled that the progress from idleness to diligence, to self-control and responsibility under the lash of want underpins the Victorian evolutionist anthropology preceding Darwin’s theories and, as La Vergata has written (*ibid.*, p. 89) “owes very little to them.” In this literature,

In his first reading of the *Essay*, Darwin had mistakenly thought Malthus intended his theses on population to refer to the human species alone – and, to be sure, this perspective does occupy the central position in his work. Darwin’s broader aim was that the perspective should be applicable also to plants and animals other than human,¹³ and it was along this path¹⁴ that he arrived at the discovery of natural selection.

Reading Darwin’s notes, one gets the impression that he by no means takes for granted the continuity of the human world with that of other living beings but, rather, sees it as a problem. He continually confronts the one world with the other, seeking to demonstrate similarities and divergences, well beyond the lines indicated by Malthus. Their positions do however converge on one important point: both are bent on tracing out a history that dates from the “ever since we have had any knowledge of mankind;” and this is why, unlike what many interpreters were to argue, the reference to contemporary economic reality (competitive conflict in the markets) is by no means evident in Malthus’s statement of the principle of population. As Manier (1978) has remarked “The famous ‘ratios’ of geometric and arithmetic series express a mathematical thesis affirmed by both men. There is nothing distinctive of political economy or of human demography in any of this.” The idea of a constant standard of living (measured by the ratio of food production to population), with oscillations based on the hypothesis of a positive ratio of standard of living to growth of population, does not of necessity presuppose a competitive mechanism – though Ricardo would invoke such a mechanism, inspired by Malthus (and Adam Smith), to illustrate the convergence of the market wage to the long-term “natural” wage. The economic and political aspects of Malthus’s thesis mainly had to do with his forecasts on the effects of exacting poor rates (leading to a rise in the prices of subsistence goods, since he assumed their production would not be affected, then to an increased birth rate, resulting in a check on the “productive industry”) (Malthus, 1798, p. 95). These propositions, according to Manier, “applied only to the behavior of human beings; they were not general laws describing (or predicting) the behavior of all organisms” (Manier, 1978, p. 81). In the view of Malthus, the “struggle for existence” was indeed a “zero-sum competition for a scarce resource (subject to the law of diminishing returns”) (ibid., p. 82), but at the same time it acted as a brake

the creatures who are “idle par excellence” are the savages, the peoples of backward countries. The influence of this anthropology (and of the pedagogy that derives from them) can be felt in several of Hayek’s writings, which take inspiration from this literature. See, for example, the idea – in (von Hayek, 1982) discussed in Hodgson (1993) – that egalitarianism is an atavistic, tribal residue, based on primordial emotions and subsequently surmounted by social rules based on the spontaneous order of individuals.

¹³ K. Marx (cf. Marx, 1969, p.121) remarked that, “In his splendid work, Darwin did not realize that by discovering the ‘geometrical progression’ in the animal and plant kingdom, he overthrew Malthus theory. Malthus theory is based on the fact that he set Wallace’s geometrical progression of man against the chimerical ‘arithmetical’ progression of animals and plants.”

¹⁴ By taking this route, Darwin was actually led to ignore the differences between the human species and the other living species (moral checks) that Malthus had introduced after the first edition. Cf. La Vergata (1990a), pp. 354–355.

on change (not as an agent of change) since it imparted cyclical oscillations around the standard allowed by the stage of development achieved. Darwin's notion of the "struggle for existence" was, as we shall see, quite different.

We must emphasize, however, the presence of a link (arithmetic, first of all) between the indicator of the standard of living employed by Malthus and the relation that would be employed, after the "dissolution" of the Ricardian school, by the "wage fund theorists" for determining wages in a capitalist market. For Ricardo, the necessary consumption was obtained by multiplying wages at natural long-term level by the number of workers employed. The difference between production and necessary consumption, that is social surplus, divided by the capital employed was used to calculate the rate of profit and the relative prices. It followed that the main argument about distribution, for a given level of rent, turned on wages and profits. The followers of Ricardo (James Mill, McCulloch and others) transformed the identity relation from which necessary consumption was obtained (wages multiplied by employment) into a causal relation, which is also an arithmetical truism: they now took the wage fund (considered as the fund of available wage goods) as a given, by which each variation in the wage (or in the amount doled out to the non-working poor) leads to an inverse variation in employment. The argument over income distribution, *for a given wage fund*, is thus transferred within the working population.¹⁵ It is, then, the conception of the wage fund, and subsequently around 1870, the marginalist (or neoclassical) theory, concentrating on competition among the owners of the same scarce resource (though of different productivity, in the case of the marginalists), that provide the backdrop for the analogy between natural selection and competition that would be asserted by the exponents of "social Darwinism." For both these conceptions feature an inverse relation between real wages and employment that, however, finds no counterpart in the theoretical perspective of classical political economy.

Scholars bent on recruiting Darwin among the supporters of so-called "social Darwinism" have argued that he had rigorously confined the principle of evolution based on natural selection to the animal and vegetable kingdoms – though, in *The Descent of Man* (Darwin, 1871), which followed the publication of the *Origin of Species*, he did, very tentatively, apply the principle to cultural, moral and social aspects. But this argument does not stand up in the light *inter alia* of the *Notebooks* (first published in a complete edition in 1987) (Schweber, 1977, p. 232). If by the term "social Darwinism" we understand "the more general adaptation of Darwinian, and related biological concepts to social ideologies,"¹⁶ there can be little doubt that it should also embrace Darwin himself. As has been remarked,¹⁷ we must try to avoid

¹⁵ Cf. Picchio (1992), p. 34, which should also be consulted (Chapter 2) for the influence exerted by Malthus's theory of population in the transition from an analytical picture based on Ricardo's conception of the natural wage (in which differences in the growth rate of the population and in production determine oscillation in the market wage around the long-term wage) to the theory of the wage fund.

¹⁶ This is one of the definitions mentioned by Bannister (1979), p. 5.

¹⁷ Cf. Young (1985), p. 609. The quotations that follow are taken from p. 626.

the positivist illusion that in historical reconstruction we can draw a clear distinction between scientific observations and theories and the values and meanings they had for their authors. It is hard to disentangle the concepts drawn from biology and from ideology – as we have seen in the case of Malthus – so that we may conclude that “it is from society that we derive our conceptions of nature” and of human nature. Therefore, “the intellectual origins of the theory of evolution by natural selection are inseparable from social, economic and ideological issues in nineteenth century Britain.” This is the conclusion Young arrives at after substantiating the presence of Malthusian accents¹⁸ in Darwin and fact that his friend and disciple T.H. Huxley on a series of topics (concerning women, blacks, uncivilized peoples, the lower classes, etc.) took “positions that were relatively progressive for their time, but relatively shocking to our eyes” (Young, 1985, p. 958).

And yet, within the “family” of social darwinists, which owes its recurring popularity to its wide range,¹⁹ the differences are very significant. One way of highlighting these differences may be to adopt the definition of “social Darwinism” given by Hofstadter, which is narrower than the previous one. It consists essentially of three propositions:²⁰

1. Identification of the “struggle for existence” among living species competing in a market for products and factors under conditions of scarcity;
2. Identification of “survival of the fittest” in the process of selection occurring in a competitive market, and of its assumed “optimal” outcome: that in nature, as in the social field, the winners will be the most successful competitors in the division of resources, and that the selection process will lead to a continual improvement both of species and of society; and
3. The (metaphorical) assimilation between times and ways of change in the biological and social areas: thus, just as the generation of new species through natural selection requires long times, so social change must proceed by slow, gradual accumulation, without being accelerated or forced.

¹⁸ See for example Darwin (1871), quoted in Young (1985), p. 619. But see also, for examples of distancing from the ideological conceptions of Malthus, La Vergata (1990a), p. 361.

¹⁹ On this point, the reader is referred to La Vergata (1985), p. 958.

²⁰ Cf. Hofstadter (1955). But, according to Hawkins (1997), who tends to include also Darwin in his definition, the “social Darwinists” share the following five propositions: (1) the laws of biology regulate the whole of organic nature, including human beings; (2) the pressure of population growth on resources generates a struggle for existence among organisms; (3) physical and mental traits conferring an advantage on their possessors in this struggle (or in sexual competition) could, through inheritance, spread through the population; (4) the emergence of new species and the disappearance of others is explained by the cumulative effects in time of selection and inheritance; (5) the process of natural selection described in the preceding points can also explain “many, if not all, aspects of culture – religion, ethics, political institutions, the rise and fall of empires and civilisations, in addition to many psychological and behavioural aspects. Social Darwinists, then, endorse two fundamental facts about human nature: that is continuous with animal psychology, and it has evolved through natural selection”.

I think that it would be difficult to find in Darwin's writings any acceptance of these three points; they entail not only considerable simplifications but also a direct correspondence between biological and economic categories. The lack of any such acceptance might be held to be due to Darwin's tendency to caution – the same caution, for instance, that led him to delay publishing his discovery. But this caution reflects, in turn, a complexity stemming from the interaction between the general orientation of Darwin's researches and the analytical categories employed by him. In the first place – and this has to do with the “psychology of discovery,”²¹ – we need to draw a distinction between that general orientation, which led him to discover the mechanism of change through natural selection, and the hypotheses underpinning the specific mechanism that he identified and its significance. With regard to the general orientation, one may accept the idea, put forward by Ospovat (1979), that Darwin's deism initially influenced the view that the general trend of evolution is progressive, and that evolution itself is the means “by which the harmony between the organism and the environment is kept” (La Vergata, 1990a, p. 528). There can be no doubt, however, that the content of his theories in no way depends on these convictions (*ibid.*, p. 528 et seq). For that matter, as regards his general orientation, Darwin himself admits to “oscillations from one and the other extreme,”²² sometimes emphasizing, along with the benefits of nature's harmony, the great suffering and pain that accompany the processes of adaptation and/or destruction of species. Kohn²³ speaks of “ambiguity” (though “fertile”) in Darwin's rhetoric, due to the conflicting intermesh between a natural-theology component, a materialist component in the conception of nature and a liberal-radical component in the sphere of religion.

At the general conceptual level, Darwin's thought is certainly complicated and difficult to condense into effective formulas. But no less rich in articulations, intersections and nuances is the idea of “struggle for existence”, that appears very different from Malthus's concept, as it also diverges from the concept implied in the definition of “social Darwinism” propounded by Hofstadter, where it is assimilated to the competition obtaining in conditions of scarcity. If for no other reason than that, as Darwin himself underlines, he used “this term in a large and metaphorical sense” (Darwin, 1859). The metaphor (La Vergata, 1990a, p. 300), as we saw, is built around multiple significances that, says Darwin (1859), “pass into each other.”

²¹ Here I take up, in a somewhat different context (but for reasons already mentioned I think the reference to Smith is not relevant), a distinction made by Gould (1980): “I believe that the theory of natural selection should be viewed as an extended analogy – whether conscious or unconscious on Darwin's part I do not know – to the *laissez faire* economics of Adam Smith. . . . But the source of one idea is one thing; its truth or fruitfulness is another. The psychology and utility of discovery are very different subjects indeed. Darwin may have cribbed the idea of natural selection from economics, but it may still be right”. On this point see also La Vergata (1985), p. 956.

²² Cf. Darwin (1887), p. 304; quoted in La Vergata, *ibid.*, p. 535, and, more in general, see Chapter 8 of the latter.

²³ Cf. Kohn (1989), p. 214, quoted in La Vergata, *ibid.*, pp. 535–536.

By “struggle” he understands “an effort to overcome a difficulty” (Manier, 1978, pp. 82–83) through relations of dependance, through variations arising at random, and through direct competition for the use of space or goods. “Struggle of organisms for existence,” however, also implies that they exhibit “varying degree of success in the effort to survive and leave fertile progeny, in some single, but complex, environment” (ibid., p. 82). The “struggle,” then, is not merely for food, because resistance to difficulties and “success in leaving progeny” will be equally important, nor is it only a direct contest among species: it is “the mesh of relations who link together all the forms of life nowadays called “ecotype,” which is responsible for possible unexpected outcomes of natural events” (La Vergata, 1990a, p. 287). And this mesh of relations also comprises forms of cooperation and association (ibid., pp. 298–299), which must therefore be included in the idea of struggle, rather than set against it. Manier concludes that the presence of so many meanings attributed to the term “struggle” is in sharp contrast to the employment of the term in the *Essay on Population*, and poses afresh the question of Darwin’s indebtedness to Malthus.

In order to account for phenomena in the natural and human worlds whose wide variability and complexity found no satisfactory explanation, Darwin devised a new language; he employed metaphors which, as Manier has shown, allowed him great flexibility in knowledge, while on the other hand they underlined the inseparable link between the cognitive and the emotional dimensions,²⁴ and thus conveyed a profound harmony with his vision of nature, morality and science.

²⁴ Cf. Manier (1978), pp. 37–40 and 89–96. The significance assumed by the metaphors in Darwin’s thought at the time when the theory of natural selection was formulated can be retraced, in Manier’s view, above all to the influence exerted by the Scottish Realism of Dugald Stewart and by the idea of nature expounded by Wordsworth in *The Excursion*. Stewart had conceived the historical development of language as a metaphorical (‘artificial’) elaboration of a system of communication developing initially from “natural signs” belonging to the affective sphere and linked with primordial biological functions. Discussing the origin and significance of terms like “beauty,” Stewart criticized the “essentialist” thesis that the different meanings must be connected with one another like the species of a common genus. In this, he anticipated Wittgenstein’s idea of “family resemblance,” which, in biology and anthropology, has found its counterparts in polythetic classifications, cf. for example Beckner (1968) and Needham (1975). In his notes on Stewart’s essays *On the Sublime* and *On Taste*, Darwin shows he understands Stewart’s criticism of the “essentialist” idea in his comment: “D. Stewart does not attempt by one common principle to explain the various causes of those sensations, which we call metaphorically sublime, but... it is through a complicated series of associations that we apply to such emotions the same term” (cited in Manier, 1978, p. 39). Wordsworth’s poem proposed “a natural aesthetic-sentimental religion, not a philosophical-theological one, which seeks the sense of life in nature and human experience” (cf. La Vergata, 1990a, p. 313). The young Darwin was thus enabled to obtain, Manier argues, an “antidote” to the positivist methodology, eluding in this way reductionism and mechanism. Darwin’s metaphors also took on an “anthropomorphic” note in his description of animal behaviour – not to say a “moralistic” note, in which the cognitive and the affective dimensions become inextricably intertwined (cf. Manier, ibid., pp. 173, 194–195).

4.3 Darwin, Spencer and the Application of Natural Selection to Social Systems

Among the “social Darwinists” (in the first definition given here), Darwin is the only “real” Darwinian. This underlines that what passes for “social Darwinism” very often derives from Spencer’s ideas on evolution rather than Darwin’s. Spencer’s definition of “evolution” is consistent with the etymology of the word, which implies “unravelling”, i.e. a predetermined “unfolding” by a given entity (an essence). It involves a one-directional movement towards progress through gradual specialization and differentiation. Diversity is the teleological outcome of the process of evolution, not its starting point. The process of evolution leads to an equilibrium (Spencer calls the process “equilibration”) (Hodgson, 1993, p. 83), in which the functional integration of the components of the organism ensures harmony and coherence. This kind of evolutionism, then, displays a strong functionalist imprint. In addition, its conception of the “survival of the fittest,” consistent with the essentialist grounding of its thought – only slightly weakened by his Lamarckism (which admits the inheritance of acquired characteristics) – differs profoundly from the conception of Darwin. The idea that there exists in nature a force that eliminates all variations or deviations is called “eliminationism” by biologists (Mayr, 1982): a static type of conception (essence is immutable, envisaging only “degenerate” variants) that must be clearly distinguished from the dynamic process of natural selection, based on continual variation.

The association of individualism and evolutionist optimism is employed by Spencer – and thereafter by many other authors ranking as “social Darwinists”²⁵ – to justify a policy of *laissez-faire*: while State intervention would reduce pressure on individuals to compete by inhibiting the assertiveness of the fittest, free competition between individuals and enterprises would provide the best environment for social progress. However, the idea of associating “survival” with “optimization from fitness” has come in for much criticism, which may also be directed towards certain interpretations of Darwinian selection.²⁶ In the first place, as Sober (1981) has remarked, it is by no means certain that those held *ex ante*²⁷ to be the “fittest” are actually those who survive. Moreover, the slogan of “survival of the fittest” regards no more than the differential in the death rate, neglecting the birth rate²⁸ (whose importance Darwin stressed). When all is said and done, Hodgson (1993) observes, with the expression “survival of the fittest” Spencer has made no useful contribution to evolutionism: he has overestimated the potential of isolated atomistic

²⁵ While some follow Spencer in exalting the progressive function of “competitive pressure,” others, like Graham Sumner, stress the need not to obstruct the process that leads to the assertion/survival of the fittest by welfare policies aiming to reduce social inequality.

²⁶ For an overview see Hodgson (1993), pp. 94–97.

²⁷ It is important to appreciate the *ex ante* character of the definition in order not to fall into the (very frequent) error of tautologically (and hence *ex post*) deriving the capacity to adapt and survive from the fact of survival.

²⁸ Cf. Gould (1982), p. 101, cit. in Hodgson (1993), p. 95.

organisms, ignoring their interactions with other organisms and with the environment; he has heeded only success and not the creative function of error (or variation, to which Darwin called attention); he has emphasized success and optimized adaptation rather than the capacity for reproduction and transmission, and has observed only conflict and ignored cooperation.

In Darwin's conception, on the other hand, evolution entails three stages: (1) a process that creates variability and diversity among the beings belonging to a population; (2) a process that selects among the entities subjected to the processes of variation; and (3) a differential transmission of the varieties selected to the subsequent generation. Single individuals are not transformed, change occurs at the level of population.

It is important to emphasize that the first two processes are independent of one another: at the individual level, the variability is random (it is largely independent of the context), whereas at the aggregate level of population the change has a direction. Selection tends to reduce the variations, while the processes that create variability and diversity tend continually to reproduce them.

In applying Darwin's conceptions to socio-economic systems, it is taken for granted that evolution will assume a Lamarckian connotation, i.e. the hereditary transmission of acquired characteristics will be allowed. The crucial point, however, concerns the consequences of intentionality, an aspect that in the evolution of living non-human species can (in a first approximation) be bypassed. The intentionality of human actions influences, first, the process of production of the variations, which will not be random and will therefore not be independent either of the process of selection or of the context. This is tantamount to saying that the process of production of the variations will be influenced in particular by the knowledge (and understanding) of the way in which the process of selection works.

But the terms in which the process of selection unfolds, too, are influenced by the intentionality (and the power) of the agents: these terms are negotiated through forms, modalities, channels that are themselves the object of negotiation. Moreover, the actions of the agents that have undergone selection are mediated by the attributions of identity²⁹ of the agents with which they interact. Since these attributions, too, are the result of a negotiation, the actions that derive from them will likewise be the result of a negotiation. And where there is negotiation (and hence influence from the context), the automatic nature of the consequences and the general nature of the outcomes cannot be assumed.

Human intentionality also influences the hierarchy of importance and the significance of the processes of adaptation to environmental variations (in a broad sense). While in the non-human world the trait that succeeds in adapting best to change, considered as exogenous and unmodifiable, tends to survive, human intentionality may propose to alter the environment itself or to wait for processes of co-evolution between environment and society.

²⁹ On these points, see Maxfield and Lane (1997).

The presence of intentionality poses serious problems in transferring the Darwinian processes of natural selection to social systems. More useful indications for reflecting on the processes of change and innovation in human societies have been developed by biologists and palaeontologists whom one might call “heterodox” Darwinists. They have underlined how Darwin himself saw natural selection as “the main but not the exclusive means of modification” (Darwin, 1859, p. 438). Gould and Lewontin comment ironically on the “panglossian” obstinacy of the “adaptationists”, who maintain that adaptation to the environment is the sole evolutionary process and that selection intervenes in the adaptation processes by acting as an optimizing factor. In order to explain forms, functions and behaviours, Gould and Lewontin (1979) have listed at least five ways in which evolution manifests alternatives to immediate adaptation. Among these alternatives we should note *exaptation*, in which the current function of an organ (but, we may add, also that of an artifact or an institution) has appeared at a moment subsequent to its origin. Gould (2002, p. 1214 et seq.) cites Nietzsche’s remark:

... there is for historiography of any kind no more important proposition than the one it took such effort to establish but which really *ought to be* established now: the cause of the origin of a thing and its eventual utility, its actual employment and place in a system of purposes, lie worlds apart.³⁰

Nietzsche went on to say:

[...] purposes and utilities are only *signs* that a will to power has become master of something less powerful and imposed upon it a character of a function: and the entire history of a “thing”, an organ, a custom can in this way be a continuous sign-chain of ever new interpretations and adaptations whose causes do not then have to be related to one another but, on the contrary, in some cases succeed and alternate with one another in a purely chance fashion. The “evolution” of a thing, a custom, an organ is thus by no means its *progressus* toward a goal, even less a logical *progressus* by the shortest route and with the smallest expenditure of force- but a succession of more or less profound, more or less mutually independent processes of subduing, plus the resistances they encounter, the attempts at transformation for the purpose of defense and reaction, and the results of successful counteractions. The form is fluid, but the “meaning” is even more so.

The concept of *exaptation* appears highly important for study of the creation processes of new market systems, connected with the discovery of new functionalities of artifacts. The concept is also of importance in studying the changes in “meaning,” and hence in function, of the institutions that underpin the reproduction of the social system.³¹ The innovations that accompany new functionalities act as a mechanism of change in the economy, creating, and simultaneously destroying, products and jobs.

Contrary to what has often been reiterated, the processes of change are not predominantly underpinned by increases in productivity – which, if at all, follow at a distance from those processes – nor by increased competition, which is actually

³⁰ Cf. Nietzsche (1967), pp. 77–78. Considerations of this kind may have led to Wittgenstein’s scepticism concerning the explanatory power of (evolutionary) historical explanations (Wittgenstein, 1975, pp. 28, 30, 50 and the observations by J. Bouveresse, p. 75).

³¹ On these structures, called “scaffolding structures,” see Lane (2002), p. 71.

attenuated, more or less stably, by the new functionalities. They rest on processes of *exaptation*, i.e. product innovations that lead to new sectors.

As Giuseppe Maione (2003, pp. 10 and 13) has effectively demonstrated, the theoretical perspective that credits productivity³² as the driving mechanism of growth rests on a mistaken transposition to the economic system of a thesis that may (sometimes) be valid for an individual firm. If all firms were to reduce their unit costs and, let us suppose, their prices, there is no certainty that this would lead to a growth in the economy. Firstly, despite the fall in relative prices (with respect to competing countries), collective tastes might move in an opposite direction from that of the products supplied. Furthermore, while a simultaneous shrinkage of all costs translates into savings for firms, it also causes a contraction for the supply firms. Barring exceptional circumstances, this would lead to recession.

The principal defect of such a perspective is that it conceives genuine innovation, i.e. “the discovery of new products and processes [. . .] only as an auxiliary element of productivity” (Maione, 2003, p. 10). Historical experience shows instead that growth, when it occurs, is the result of product innovations (always going hand in hand with process innovations) and of the emergence and gradual ramification of activities indirectly linked with the “new” final production (“new” with respect to the system with which one is dealing). Underpinned by infrastructures and institutions, new sectors of production and marketing come into being, directly and indirectly connected with the innovating sector. In this perspective, the increase in productivity (which, by the way, in the case of product innovation escapes rigorous definition) is, if anything, an element that is auxiliary and subordinate to the innovation, not the other way round.

This, however, is not the – exclusively adaptational – conception put forward by evolutionism that emerges from the exponents of the neoclassical, Walrasian or Hayekian current. Both actually credit the exogenous variations in labor productivity with a driving role in economic growth. This reflects the exclusive importance ascribed to process innovations, and hence of *price* competitiveness as against the existing product – competitiveness viewed as the force that diffuses the sole form of innovation on which they effectively focus.³³ Yet Schumpeter (1992, p. 84, *my italics*) had remarked:

Economists are at long last emerging from the stage in which price competition was all they saw. As soon as quality competition and sales effort are admitted into the sacred precincts of theory, the price variable is ousted from its dominant position. However, it is still competition within a rigid pattern of invariant conditions, methods of production and

³² The concept is, anyway, hard to measure in a service society.

³³ In the neoclassical conception, based on the existence of an aggregate production function, the presence of full employment of labor rules out that increases in output can be obtained through increasing the number of those employed. Since innovations of product are not admitted, production may increase in the long term only through upward translations of the production function, which correspond to innovations of process. It is not of interest here to evidenciate the inconsistency between favoring an idea of competition as process – which will be dealt with in the next section – and, at the same time, arguing for an evolutionism in which adaptation takes on exclusive importance.

forms of industrial organization in particular, which practically monopolizes attention. But in capitalist reality – as distinguished from its textbook picture, it is not that kind of competition which counts but the competition from the new commodity, the new technology, the new source of supply, the new type of organization (the largest-scale unit of control, for instance) – competition which commands a decisive cost or quality advantage and which strikes not at the margins of the profits and the outputs of the existing firms but at their foundations and their very lives. This kind of competition is as much more effective than the other, as a bombardment is in comparison with forcing a door, and so much more important that it becomes a matter of comparative indifference whether competition in the ordinary sense functions more or less promptly; *the powerful lever that in the long run expands output and brings down prices is in any case made of other stuff.*

And here we run up against a paradox: economists of the neoclassical school (to whatever current they belong) tend, on the one hand, to attribute a crucial role to competition in stimulating the growth of the economy while, on the other, they confine their attention to a form of competition so unimportant that, according to Schumpeter, it “becomes a matter of comparative indifference” whether it operates more or less promptly.

Underlying this paradox is a fundamental aspect of neoclassical theory, namely the simultaneous determination of prices and amounts produced. At first glance, this appears to be a powerful tool of knowledge, but further reflection reveals it as a source of unrealistic hypotheses and insurmountable limitations in drawing conclusions. In order to compare the different roles ascribed by the theories to competition and growth, we must briefly review the main differences in the idea of competition as envisioned in the three main currents mentioned above. The origin of these differences can be retraced to the underlying theory of relative prices and income distribution.

4.4 Three Conceptions of Competition

In the notion of competition in classical political economy, prices are “prices of production.” In fact, these are prices of reproduction, i.e. they are fixed in order to ensure the reproduction, year by year, of the process of production and sale of the goods. They are long-term prices, determined in a different way from short-term prices, called market prices, subject to different, less persistent influences. The adjustment of market prices to production prices may occur in a variety of ways: since the long-term positions are not here considered as *functions* of demand and supply, they do not arrive, as in neoclassical theory, along a single route, namely the annulment of excess demand and supply. Given the wage and the expected amounts of output, the prices must guarantee coverage of the normal costs expected and the achievement of the rate of profit; the latter, in a competitive regime (i.e., in this conception, a situation where firms are free to enter and exit from the market), may tend to equal the rate of interest on bonds without risk, plus a reward “for the risk and trouble” of the entrepreneur. That the levels of output are taken as given is important for two reasons. The first is that the formation of prices is held to be an aspect of income distribution, and not, at the same time, an aspect of the

determination of the quantities produced and hence of the growth of nations. The second is that it is acknowledged that the process of determination of the amounts produced is subject to circumstances so complex and changeable (belonging, we should now say, to such different levels of interaction and time and spatial scales) as to elude easy generalizations. In a letter to Malthus of 9th October 1820, Ricardo³⁴ writes:

Political Economy you think is an enquiry into the nature and causes of wealth. I think it should be called an enquiry into the laws which determine the division of the produce of industry amongst the classes who concur in its formation. No law can be laid down respecting quantity, but a tolerably correct one can be laid down respecting proportions. Every day I am more satisfied that the former enquiry is vain and delusive, and the latter only the true object of the science.

This theoretical and methodological formulation has three consequences. First, prices are not indicators of scarcity (scarcity is not the central point of the analysis, but in any case it is not identified by the relation between demand and supply functions). Secondly, equilibrium prices are linked to a particular rule of distribution of the surplus (for example, the rule, peculiar to a competitive regime, by which the rates of profit are equal in every productive sector), but not linked with any connotation of optimality, efficiency or full employment of resources. Thirdly, increasing returns and processes of circular causation are perfectly compatible with the conditions that ensure competition. The separation of the circumstances that determine the amounts produced (among which, for example, the level of aggregate demand) and those that determine prices (for given amounts produced) does away with the need to impose restrictive conditions (on technology, tastes, etc.) that find their justification only because they enable an equilibrium situation to be reached.

A quite different description of competition is provided by the Walrasian school. Here the concept of “perfect competition” is introduced, with atomistic firms, sufficiently small to be unable, by hypothesis, to affect the price (assumed as given for the individual firm, called *price taker*). By maximizing the profit subordinately to technology constraints, these firms demand labor services and sell products. The purchasers are consumers, likewise “atomistic,” who obtain income by selling labor services and maximize their utility subordinately to the budget constraint. The prices of the products and the so-called “production factors” and the quantities produced are determined simultaneously. Based on certain hypotheses (which are actually somewhat restrictive, with regard to the form of technology and the consumer preferences), it is demonstrated that there exists a set of relative prices and quantities exchanged that causes demands and supplies to be equal. These prices and quantities are said to be “equilibrium.” The prices are indicators of scarcity (they reflect the relative abundance, with respect to the demand, of the “factors” and products). At equilibrium, all the resources are employed fully, so that unemployment cannot exist. Then there are two theorems, the first and second theorem

³⁴ Cf. Ricardo (1973), p. 278 and Keynes (1971), Chapter 2, note 2.

of welfare economics, stating that these competitive equilibria have the property of being “optimal” according to Pareto’s definition³⁵ and, conversely, that any optimal situation according to Pareto may accompany a competitive equilibrium, achieved by starting from some initial endowment of resources. In principle, this theorem could be used to reconcile competitive capitalism and egalitarianism: any desirable equilibrium situation could be reached by a suitable redistribution of the initial resources.

The possibility to mix competition (as described in the marginalist theory) and public intervention belongs within a current of “normative marginalism” that, right from the beginnings of this theory, goes hand in hand with the “marginalism of *laissez faire*.” For, starting from the research Walras developed in his mature works, in the neoclassical perspective of equilibrium, two “views” exist side by side: one descriptive, directed towards analyzing the actual working of markets, though abstracting from the frictions and the contingencies as suggested by the methodology of classical physics,³⁶ the other normative, based, in Walras’ case, on transferring principles rationally deduced from natural law to economics.

“From the outset,” according to B. Ingrao and G. Israel (1990, p. 98), “the interest shown by Walras in the theory of exchange value and even the formation of the idea of a general market equilibrium were motivated by the search for a rigorous demonstration of the superiority of free competition as a form of the organization of production and exchange. He identified free competition as an ideal of commutative justice or justice in exchange, in which he recognized a necessary but not sufficient condition for the fulfillment of distributive justice according to the principles deduced from natural law.”

The realization of justice in exchange was attributed not only to the condition of maximum utility of each party but also to the formation of a single price for the same good and to the elimination of losses and non-motivated profits, at the equilibrium price. From this point of view, free competition was for Walras “a normative ideal toward which the actual functioning of markets should be directed” (save for unavoidable instances of monopoly).

At the end of the 1930s, the current of “normative marginalism” re-emerged, as a result of two circumstances: the difficulty of the “descriptive” theory and the proposal, put forward by O. Lange, for a decentralized marginalist “market socialism”, as invoked, by both supporters and critics, by the experience of the USSR’s centralized planning after the 1929 crash.

³⁵ The equilibrium configuration is “optimal” from Pareto’s point of view if, given the technology and tastes, there is no alternative equilibrium configuration that could improve anyone’s situation without worsening anyone else’s situation (where “improve” and “worsen” are defined in terms of the respective preferences).

³⁶ The close analogy between the equations of the exchange proposed by Walras to determine the equilibrium prices and the methods and results of classical mechanics has been most effectively evidenced by Ingrao and Israel (1990), who have also made a reconstruction of the two “views” mentioned in the text.

With regard to the first point, it should be recalled that throughout its development in history the theory of general economic equilibrium has constantly aimed to demonstrate three fundamental results: the existence of equilibrium (i.e. the compatibility among the actions of the maximizing economic actors in a competitive market), the uniqueness of the equilibrium (absence of other equilibrium states – an indispensable condition *inter alia* for performing analysis of comparative statics), and the global stability of the equilibrium (the market forces “inexorably” lead to the state of compatibility of the actions of the actors). Already in the 1930s it was perceived that the conditions necessary for a system of decentralized decisions (as described by the theory) to work are so stringent as to make it impossible in practice for them to exist.³⁷ Attention was focused on two problems in particular.

The first of these concerned the process by which equilibrium was approached, the *tâtonnement*: Walras adduced the figure of the auctioneer announcing the prices to refer to the really existing auction markets; at the same time, by excluding exchanges outside the equilibrium, he removed the obstacles by which such exchanges, with the associated variation in incomes and prices, would hinder the attainment of equilibrium. The process becomes analogous to an abstract mathematical algorithm which can be used, if need be, for normative purposes, but remains light-years away from a description of the effective movements of prices.

The second point concerned an implication arising out of the deterministic character of the theory (namely, the hypothesis that the current state of the economic system univocally determines its evolution in time). Poincaré, the great mathematician to whom Walras had submitted his writings, at once detected this with singular acuteness. Replying to Walras, he noted that, whereas in mechanics friction is disregarded and bodies are treated as infinitely smooth, “you consider people as infinitely selfish and infinitely far-sighted. The first hypothesis may perhaps be admitted in a first approximation, the second may call for some reservations” (Ingrao & Israel, 1990, p. 159). With polite irony, Poincaré put his finger on one of the shakiest points in Walras’ formulation, the hypothesis of perfect far-sightedness, or clairvoyance.

The binding character of this hypothesis³⁸ would clearly emerge in the theoretical research along Walrasian lines in the 1930s, to which Hayek, among others,

³⁷ Subsequent researches (associated mainly with Arrow and Debreu) have found that a formal proof can be given only of the first result. The second has been obtained with very restrictive hypotheses, while the proof of the third, which is of decisive importance, has had a negative outcome. Failure on the last two points casts doubt on what seemed to be one of the strong points of the theory – namely, deriving a series of important results starting from very general hypotheses on the behaviour of the agents. The general logic of this choice is too “flimsy” in terms of structure to be able to impose economically significant restrictions and to avoid, in this way, “perverse” behaviour of the demand and supply functions (i.e. instability and multiplicity of the equilibria). The upshot is that “the endeavor to keep the theory at the highest level of generality, proved to be one of its weakest point” (Ibid., p. 316).

³⁸ The difficulties of the Walrasian theory stem from unsolved problems of capital theory. Since Walras chose to take as given the individual amounts of capital goods (and not their totality in

contributed. In 1935, Morgenstern, weighing up this research, observed “the prevalence in the literature of the belief that the theoretical perfection of equilibrium could not be achieved without the hypothesis of complete foresight on the part of agents” and concluded that the hypothesis “should be drastically eliminated from economic theory” (Ingrao & Israel, 1990, p. 195). He added: “a theory of equilibrium which ‘explains’ only a *static situation, which is given and unalterable* and which, because of this basic assumption, is completely unable to say anything about the economy when a variation occurs, is utterly unimportant from a scientific point of view. It would hardly deserve the names of theory and science.”

The difficulty in addressing an economic crisis with the measures suggested by the orthodox theory (cutting wages to enable competition to operate freely), on the one hand demanded a formulation other than that of general equilibrium (Keynes), while it led even those who remained faithful to the Walrasian formulation to take an interest in planning. For there was indeed another possibility, to interpret competitive equilibrium in normative rather than descriptive terms.

This path was taken in two famous articles of 1936 and 1937 by Oskar Lange. As Ingrao and Israel (1990, p. 253) remark, he must be seen as a “direct heir of the normative tradition” of Walras, though in a perspective oriented towards planning rather than natural law doctrine. Von Mises had argued that only in a regime of private ownership can prices figure as an efficient gauge of scarcity of resources. Lange retorted that a system with public ownership, too, could achieve the equilibrium situation by miming a decentralized market, i.e. by assigning both a price list and maximizing norms of behavior to the units of consumption and production.³⁹ Lange

terms of value, which would have posed insoluble problems of circular argument in determining prices), this led to a short-term general equilibrium in which the rates of profit on the various capital goods were not uniform. In subsequent periods, the effect of competition, tending to uniformity of the rate of profit, would inevitably have induced changes in the composition of the capital stock, and hence in the prices of the goods and productive services. On this point see Garegnani (1976, pp. 25 et seq). This could happen either through an equilibrium on intertemporal base, founded on forward markets (for all future dates) of all the goods and services (as suggested by Hayek in the 1930s), or through the hypothesis of a succession of temporary equilibria with expectations of price realized, in steady state conditions, or of uniform development. The development of the Walrasian formulation led, therefore, to a new notion of equilibrium, temporary and not long-term, subjected to unrealistic hypotheses (the diffusion of forward markets, equivalent to perfect foresight) or so restrictive and contingent (temporary equilibrium with given expectations in steady state) as to call in doubt whether it could act as a guide to identify the basic forces that determine the real working of the economy.

³⁹ This gave birth in the USA to a current of “normative marginalism” which was to involve, with lavish funding (the Cowles Commission, working groups of the Rand Corporation), economists sensitive to the correction of the disequilibria present in the market economies (Arrow, Debreu and Koopmans among the economists engaged in researching the properties of equilibrium, Dorfman, Samuelson and Solow on the theme of including the input output analysis of Leontief within the theory of general economic equilibrium). Within this current of “normative marginalism” the neoclassical synthesis takes its place (synthesis of neoclassical theory and Keynesian theory, with possibility of short-term equilibria without full employment, but support in the long term for the neoclassical theory). Samuelson and Solow figure among the leading exponents of this synthesis.

substantially entrusted a centrally directed market with the task of implementing Walras' process of *tâtonnement*.

While Hayek in 1940 acknowledged that Lange had arrived at a socialist system with a competitive market, he asserted that "it is difficult to suppress the suspicion" that the formulation of a process of trial and error to gradually approach the appropriate solution "has been born out of an excessive preoccupation with the problems of the pure theory of stationary equilibrium. If in the real world we had to deal with approximately constant data, that is, if the problem were to find a price system which then could be left more or less unchanged for long periods, then the proposal under consideration would not be so entirely unreasonable. But this is far from being the situation in the real world, where constant change is the rule" (von Hayek, 1940, pp. 187–188). In attempting to show the superiority, in processes of adjustment, of the competition of the capitalist market as against the fictitious one of planning, Hayek adduced the demands of realism that appear important also outside the (highly abstract) debate on planning – demands that the theory of general equilibrium (along with its application to planning) had entirely disregarded in hypothesizing virtual and/or instantaneous adjustments.

The conception of competition adopted by the Austrian school (especially by Hayek) is opposed (within limits I shall indicate) to that of the Walrasian school, which it sees as a possible antechamber for public intervention in economics and "hence," it argues, for planning. It criticizes that theory for having concentrated on the conditions of static, at the expense of dynamic, equilibrium, i.e. the process by which equilibrium is reached. It underlines the limits of the definition of "perfect competition" – in particular, the hypothesis of perfect knowledge (of markets, technology, the future), the tendency of products and techniques to be homogeneous, the absence of change (in endowments, tastes and technology). In these respects, it recovers a perspective – that of competition "as a procedure of discovery of the new" – that may be likened, *at least at first sight*, to that of classical political economy (or at least to non-marginalist conceptions).⁴⁰ In many other respects, however – and these are the crucial ones – it retains unaltered the categories of marginalist economics, as though they could survive unharmed the criticisms previously made.

⁴⁰ The idea that "competition is a discovery procedure" (Hayek, 1978) has precursors that Hayek himself would not have welcomed (See, for that matter, the reference to the essay by the sociologist L. von Wiese, published in 1929, and cited in von Hayek (1940, p. 179, note 1). Suffice it to recall, for example, the idea that competition generalizes the methods of production (Marx, 1962, p. 259), and the remarks on its "socializing force," to which Simmel (1908) devoted some enlightening pages. For an analysis of these passages see Bonifati (2003), to whom I am grateful for useful discussions on this point. Simmel (1908) writes, for example, that "it is the poisonous, dispersive and destructive effects of competition that are usually emphasized . . . But alongside these there exists an enormous associative effect: competition compels the aspirant who has beside him another aspirant, and who only in this way really becomes an aspirant, to approach and come to grips with the object of the competition, to connect with him, to learn about his weaknesses and strengths and to adjust to them, to look for or to try to construct bridges that might link his own being and his own performance to him." (my translation)

Up to 1937, the year when Hayek published an essay entitled *Economics and Knowledge*, he had “practically identified economic theory with equilibrium theory.”⁴¹ This article marks a turning point, insofar as he distinguishes between an equilibrium at the individual level and an equilibrium for society as a whole. The first of these follows the postulates of the Pure Logic of Choice – therefore, at microeconomic level, of the maximization of the utility subjectively perceived by each agent. The equilibrium at society level, on the other hand, takes for granted the fact that knowledge is subjective (it cannot be known by any central planning mind), dispersed, imperfect, dominated to different extents by ignorance, error and inconsistency. Where the theorist of equilibrium assumes the existence of objective, certain, stable data that enable a configuration of equilibrium to be attained, here the assumptions needed for this to come about are lacking. And yet an equilibrium configuration is reached: hence, Hayek concludes, the market *in practice*, through the system of relative prices, brings about the transmission of information among individuals that, at the level of society, enables coordination of their individual plans. “Competition,” von Hayek (1948) concludes in another essay devoted to the *Meaning of Competition*, “is essentially a process of the formation of opinion: by spreading information, it creates that unity and coherence of the economic system which we presuppose when we think of it as a market.”

How can we trace the outcomes, at the societal level, of an exogenous change that arouses adaptation behavior in individuals? Here the theory of equilibrium is once again of service, not to trace the path, but to describe the goal. For von Hayek (1945) writes: “It is in this connection [of adjustment to changes] that what I have called “the economic calculus” (or the Pure Logic of Choice [that is the theory of the single decision maker] help us, *at least by analogy* to see how this problem can be solved, and in fact is being solved, by the price system”.⁴² The proof of the informative value of the system of prices – understood as a device for transmitting information – therefore relies “by analogy” on the microeconomic theory of equilibrium, while the scheme that should demonstrate how equilibrium is reached at the level of the economic system is criticized for its want of realism.

At the outset of the 1940s, then, Hayek finds himself facing this dilemma: as long as one insists on the concept of general economic equilibrium, one cannot avoid dealing with a normative interpretation of the scheme of general equilibrium that envisages the possibility of using it for planning purposes. As against that, the criticisms of the Walrasian process of *tâtonnement* can be used to criticize the lack

⁴¹ Cf. Caldwell (1988), p. 529 and von Hayek (1940), pp. 33–56.

⁴² Italics added. He goes on to say: “Even the single controlling mind, in possession of all the data for some small, self-contained economic system, would not- every time some small adjustment in the allocation of resources had to be made – go explicitly through all the relations between ends and means which might possibly be affected. It is indeed the great contribution of the Pure Logic of Choice that it has demonstrated conclusively that even such a single mind could solve this kind of problem only by constructing and constantly using rates of equivalence (or ‘values’ or ‘marginal rates of substitution’), that is by attaching to each kind of scarce resource a numerical index which cannot be derived from any property possessed by that particular thing, but which reflects, or in which it is condensed, its significance in view of the whole means-ends structure.”

of realism in the real adjustment processes, but they by no means ensure (unless by “analogy” with the theory of the individual agent) that competition will attain a position in some sense of equilibrium and “optimal” (at least in comparison with alternative systems).

To solve this dilemma, and to contrast the threat of “normative marginalism,” Hayek’s thought (apparently) changes tack: it upholds evolutionism and replaces the concept (or term?) of “general equilibrium” with that of “spontaneous social order”. The competitive process is likened to that of natural selection. The latter works as a process of adjustment to equilibrium, fulfilling the role previously, and unsuccessfully, ascribed to the Walrasian *tâtonnement*.

Several objections can of course be made to Hayek’s way out of the impasse. In the first place, it assumes what should be proved, namely the effectiveness of the markets in coordinating decisions, an effectiveness open to grave doubts owing to widespread situations of crisis and unemployment. Moreover, one may doubt that prices really are indicators of scarcity and the sole, most satisfactory vehicles for transmission of information:⁴³ suffice it to think of all the relations within and among firms channelled through formal and informal agreements, written and unwritten customs, contracts, licences, exchanges based on trust, etc. As Richardson (1972) has argued, the formation of networks among firms where coordination takes place through agreements and not through impersonal relations of price, is particularly widespread when activities are complementary and dissimilar, and when the problem of quality control of the products exchanged is especially important – aspects, these, by no means uncommon in the economic system. Furthermore, the diffusion of these agreements is also fuelled by a circumstance – the uncertainty of the level of aggregate demand, which in Hayek, given the presumed equilibrating efficacy of prices, is entirely lacking. Lastly, and most importantly, we may wonder to which evolutionary mechanism Hayek refers.

In the next section we shall see that in order to understand his position it is necessary to adduce three components of his thought, that can be traced respectively to Menger (organicist individualism), Mach (physicalism) and Spencer (non-Darwinian organicist evolutionism). Though these cultural references may help to explain Hayek, they can hardly provide useful indications for dealing with the theme posited by him. An important problem, for that matter, is how the coordination (or lack of coordination) of decisions works in a world in which the knowledge necessary for individuals to be able to act – but also to construct a balanced set-up at the level of the economic system – “never exists in concentrated or integrated form but solely as the dispersed bits of incomplete and frequently contradictory knowledge which all the separate individuals possess” (von Hayek, 1945, p. 77, italics added).

⁴³ Rosenberg (1976, p. 110) has argued that economic incentives (e.g. the prospect of profit) are a variable “so universal” as to have a very limited power to explain “the particular sequence and timing of innovative activity”. He has therefore evidenced a series of “induction mechanisms” and “focusing devices,” other than price, that underpin possible sequences of technological change.

4.5 The Liberalism of Hayek

Some much-debated questions arise from the curious mixture of contributions that converge in Hayek's conceptions: the falsifiability of theoretical propositions (including those of Hayek himself), the nature of what Forsyth has called Hayek's "bizarre" liberalism, and the possible inconsistency with his professed support for methodological individualism. Here I shall confine myself to some brief schematic hints on the cultural ancestries of these.

Take, for example, the distinction (Donzelli, 1988, pp. 39–40). Hayek proposed between "explanation of detail" and "explanation of principle;" this harks back to Carl Menger and his own particular Aristotelianism,⁴⁴ which was to influence the entire Austrian school of economists. The first of these terms refers to explanations of events about which it is possible to formulate general laws and particular propositions that describe the particular conditions upon which the occurrence of the event studied depends. The second refers to explanations of events "so complex," in the sense of involving a configuration of elementary events connected by such a "complex and inextricable texture" and such a large number of relations and variables that an "explanation of detail" is impossible. In this case, it is the theoretician's job to enunciate the general laws that govern the phenomenon of study. These laws enable one to predict certain qualitative characteristics of the "order of the events" (as in the case of the "spontaneous order" mentioned above), but do not allow the occurrence of particular events to be predicted.

In 1944, Hayek interpreted the equilibrium of an economic system, which Walras had dealt with, as an instance of argument susceptible exclusively of "explanation of principle," in which the complexity of the events studied would hinder (and make nonsense of) attempts at "explanation of detail." It becomes plain that the distinction conduces to discourage any normative proposal (even the collection of data for this purpose⁴⁵) for focusing attention on the qualitative outcomes of general laws.

The attention paid to these general laws – by which Hayek seems to dodge empirical verification of his propositions – has its origin in Menger's Aristotelianism; Menger holds that in reality there are entities that are "strictly universal" ("*essences*") (Smith, 1990, p. 266). These are discovered through the theoretical effort of the economist, whose task it is also to identify the general connection (the "exact laws") among the entities that make up the economic phenomena. On the one hand, the intelligibility of the basic structures of which the reality consists

⁴⁴ According to Smith (1986, p. 9 and 1990, p. 263), the source of Menger's Aristotelianism can be traced back to the school programs of Habsburg Austria, which imposed rigid, uniform textbooks based on simplified versions of the philosophy of Aristotle and the Scholastics, with the aim of insulating Austria from harmful liberal and cosmopolitan influences coming from abroad. Smith adds that this Aristotelianism has been shorn of any reference to the distinction between "act" and "power" (upon which, as we know, Aristotle based his analysis of change). For a discussion of the connections between these components of Aristotle's thought and Sen's theory of capabilities, to which we shall allude later, see Nussbaum (1988).

⁴⁵ Cf. for example the polemic against Leontief (von Hayek, 1978, pp. 242 et seq.).

(e.g. exchange, barter, rent, profit, ownership⁴⁶) depends on the fact that they are “universal,” i.e. they manifest in every society. On the other hand, the possibility to understand them stems from the fact that the observer is himself an individual. This enables him, through introspection, to “put himself in the shoes” of the individuals whose processes of thought and action he studies (Smith, 1990, p. 278). In trying to establish the way the different blocks of the economic reality are connected with one another to form structured social organisms (thus rejecting an atomistic view based on individual independent entities), the economist also tries, by a method called “genetico-compositive”) to identify their origin and modalities of growth and transformation.

In Menger’s view, there is an analogy between nature, the function and origin of natural organisms and that of social organisms. “Natural organisms almost without exception exhibit, when closely observed, a really admirable functionality of all parts with respect to the whole, a functionality which is not however the result of human *calculation*, but of a *natural* process. Similarly, we can observe in numerous social institutions a strikingly apparent functionality with respect to the whole. But with closer consideration, they still do not prove to be the result of an intention aimed at this purpose, i.e. the result of an agreement of members of society or of positive legislation. They, too, present themselves to us rather as ‘natural’ products (in a certain sense), as *unintended results of historical development*” (Menger, 1883, p. 130). Ontological and methodological individualism here meld, since social organisms are explained starting from the preferences, needs and aims of individuals.

Hayek would lean heavily on these themes after his thought took an evolutionist direction in the early 1940s,⁴⁷ but he associated them with Ferguson, rather than with Menger and his methodology.⁴⁸ His reconstruction of the social thought of the

⁴⁶ “Human economy and ownership,” wrote Menger (1871, my translation) “have a common economic origin, because both find their ultimate *raison d’être* in the fact that there are goods whose available quantity is lesser than what is needed, and therefore ownership, like economy, is not an arbitrary invention but, rather, the sole possible practical solution of that problem posed by the nature of things for all goods – namely, the disproportion between the amount of goods needed and the amount available.”

⁴⁷ However, the writings in which the evolutionist perspective is most clearly in evidence belong to the subsequent years. See the essays published in the 1960s (F.A. von Hayek, *Studies in Philosophy, Politics and Economics*, London, Routledge, 1976) and, above all, *Law, Legislation and Liberty*, cit., Vol. I (pp. 23–24 and 152–153) and Vol. III (pp. 154–159 and 199–202). See also *The Fatal Conceit. The errors of Socialism*, London, Routledge, 1988, Chapter 1.

⁴⁸ von Hayek (1978, p. 101) credits Menger of having “resuscitated” this idea “in a form which now [...] seems to have become widely accepted.” Hayek adds: “The point [...] which was not fully understood until at last Carl Menger explained it clearly, was that the problem of the origin or formation and that of the manner of functioning of social institutions was essentially the same: the institutions did develop in a particular way because the co-ordination of the actions of the parts which they secured proved more effective than the alternative institutions which they had competed and which they have displaced. The theory of evolution of traditions and habits which made the formation of spontaneous orders possible stands therefore in a loose relation to the theory of evolution of the particular kinds of spontaneous orders which we call organisms, and has in fact provided the essential concepts on which the latter was built.”

three previous centuries along an individualism/constructivism⁴⁹ dichotomy aims in fact to confront the functionality deriving from the absence of intentionality of social organisms with the abuses arising from the propensity of reason to design these same organisms. The central focus is on the equivalence of a competitive process with outcomes implying order and functionality and a process of natural selection. Selection (competition) occurs at all levels, involving natural, cultural and social phenomena with a unique mechanism (that does not exclude forms of group selection). Cultural evolution is seen as analogous, in some respects, to that of biology, though different in other respects. The analogy has to do, firstly, with the principle of selection (i.e. that of survival or reproductive advantage) (von Hayek, 1991, pp. 25–26). In both contexts (biological and cultural), “we have essentially”, he asserts, “the same kind of process, based on variation, adaptation and competition, however different are their particular mechanisms, especially those concerned with propagation. All evolution is based on competition: not only that, a continuous competition is necessary also to preserve the attained results.”

Cultural and biological evolution also share the fact that they are processes of “continuous adaptation to unforeseeable events” over which no control can be exerted. This is why, says Hayek, in the evolution of human societies cooperation can find no significant place: “cooperation, like solidarity, presupposes a large measure of agreement on ends as well as on means employed in their pursuit. It makes sense in a small group whose members share particular habits, knowledge and beliefs about possibilities. It makes hardly any sense when the problem is to adapt to unknown circumstances.” Hayek sees the solution to this problem in competition, seen as “a discovery procedure.” And “through further competition, not through agreement, we gradually increase our efficiency” (von Hayek, 1991, p. 19). Society is therefore “held together” not by assumed common goals, but only by the rules, the customs stemming from the evolutionary process.

The differences between the two forms of evolution have to do with the presence, in cultural evolution, of the hereditariness – through *inter alia* imitative learning – of acquired characteristics, the reception of traits not only from the family circle but from an indefinite number of “ancestors”, and the much faster transmission of cultural traits. These considerations (especially the last) lead Hayek to distance himself from sociobiology, an extreme form of “social Darwinism” based on genetic determinism. This distancing, however, concerns the refusal to view the gene as the unit on which selection is exerted, not the elements contained in the “economic” definition of social Darwinism given by Hofstadter.

⁴⁹ In the early 1940s, Hayek had sought to identify a dichotomy in the history of thought, between “individualism” of English type (in reality largely Scottish) and “scientism” or “rationalist constructivism” of French type, which he retraces to Descartes. This latter position he viewed as asserting the superiority of the “deliberate design and planning on spontaneous forces of society.” Quoting Ferguson, Hayek argues that institutions (e.g. the market, language, money) are “the results of human action but not of human design” (von Hayek, 1978, p. 96). Hayek’s view, constrained by his dichotomic classification, excludes in advance that the results of human action may also bear an imprint of human designs.

It is indeed the stress on cultural evolution that provides Hayek with the opportunity silently to replace a Darwinian process of evolution based on the *variation* of the composition of a population's characteristics, starting from the sequence variation-selection-transmission – and hence a phylogenetic process, with a quite different evolutionary process, of ontogenetic and essentialist type, that entails development of an organism starting from a set of given and immutable characteristics.

I think we shall not be far off the mark if we set Hayek within a tradition of anti-enlightenment thought that leads from Burke to Carl Menger, passing through the historical school of law of Savigny:⁵⁰ slow accumulation of traditions, wisdom embodied in the institutions, basic optimism about the way the “spontaneous forces” work, and a need to avoid all interference with the free play of these forces.

Here we have a kind of evolutionism that, as has been remarked,⁵¹ seems to hark back rather to pre-Darwinian or Spencerian conceptions than to those of Darwin himself. This is shown by retracing an intellectual ancestry which, along with Mandeville, features authors like Herder, Humboldt and Savigny, whose ideas are a long way from Darwin's, yet who, according to Hayek, “made the idea of evolution a commonplace in the social sciences of the nineteenth century long before Darwin” (von Hayek, 1978, p. 265, cited in Hodgson, 1993, p. 159). Darwin's discovery is not only minimized, it is quite misunderstood when Hayek writes that “Darwin's painstaking efforts to illustrate how the process of evolution operated in living organisms convinced the scientific community of what had long been a commonplace in the humanities” (Hayek, 1991, p. 23, cited in Hodgson, 1993, p. 160). The misunderstanding becomes glaringly obvious with the absence of a crucial aspect of the Darwinian process, the analysis of the circumstances that give rise to the variations upon which selection exerts itself.

One might say that this lacuna reveals (to use Isaiah Berlin's (1969) well-known formula) the form of liberalism to which Hayek belongs: it invokes freedom “from” (state authority, restrictions on free enterprise, etc.) instead of freedom “to” (act creatively, innovate, etc.). But we can pin down Hayek's “bizarre”⁵² liberalism more precisely if we explore the influence on him of Mach (an influence responsible for perhaps the most original elements of Hayek's thought, according to Forsyth, 1988).

Forsyth has highlighted the importance, in this connection, of *The Sensory Order*, published by von Hayek (1952), but based on his studies and writings of 1919–1920. In agreement with Mach (ignoring a few differences), Hayek favors monism, the

⁵⁰ Hayek, for example, quotes with approval a passage by Sir Frederick Pollock, which reads: “The doctrine of evolution is nothing else than the historical method applied to the facts of nature, the historical method is nothing else than the doctrine of evolution applied to human societies and institutions. . . . Savigny, whom we do not yet know or honour enough, or our own Burke, whom we know and honour, but cannot honour enough, were Darwinians before Darwin” (von Hayek, 1982, p. 153). See also Menger (1883), pp. 172 et seq. Note also Menger's stated approval of Spencer, *ibid.*, p. 150, note 55 and p. 198, note 132.

⁵¹ Cf. Hodgson (1993), p. 160 and, for certain aspects of the relations between Hayek, Spencer and Sumner, see Paul (1988).

⁵² See especially the essay *Rules, Perceptions and Intelligibility* (of 1963), now in von Hayek (1978).

idea that “the universe consists of a continuum of physical events, which is *in principle* explicable by one and the same scientific method” (Forsyth, 1988, p. 240). If *in practice* we distinguish between “physical” and “mental” phenomena, that stems only from the limits of our ability to understand their fundamental unity. Hayek claims that all “higher” mental activities, writes Forsyth, “are merely a repetition, at succeeding levels, of processes that already take part at the lowest levels. There is again a continuum, a unity. Consciousness does not mark a significant break, and human thinking of an advanced conceptual nature is only a more complex variety for animal behaviour.” Mental activity, animal and human, conscious and unconscious, can be reduced to a physical mechanism “wholly absorbed in adapting to its environment in order to survive” (Forsyth, 1988, p. 243).

Here we face a paradox that very frequently crops up in the ideas of the exponents of methodological individualism: it is claimed that all social phenomena can be accounted for by taking the individual as the starting point, but the individual is taken as given – an abstract subject, with given interests, desires, abilities – or as having extremely limited possibilities of development. To escape from the risk of abuse of reason, individuals are depicted as almost devoid of reason.⁵³

Despite some marginal differences, there is a close parallel between the evolution of animal organisms and that of the Great Open Society (or free society) as outlined by Hayek in *The Constitution of Liberty* and in subsequent works: animals and members of the Great Society alike are involved in a process of evolutionary adaptation to survive: “in the last resort, it is the relevance of [the] . . . individual wishes to the perpetuation of the group or the species that will determine whether they will persist or change” (von Hayek, 1960, quoted in Forsyth, 1988, p. 248). On the one hand, the Great Society represents the natural, spontaneous process by which human beings adapt to their environment. On the other hand, its working *requires* obedience (not necessarily spontaneous, it appears) to abstract, universal, impersonal rules,⁵⁴ devoid of substantive content (as might be the rules that safeguard competition in the markets). This is a “purely formal liberalism”, Forsyth concludes, “a combination of normless factualism” (the evolutionism of the entire living world) “and factless normativism” (the abstract rules).

The restricted substantial significance that can be associated with the attainment of the simple survival criterion has recently led Sen to criticize the idea of “progress” linked with Darwinian evolutionism, based on the “quality of the species” (i.e. on characteristics that have survived because selected and handed down). By way of alternative, he proposes to resume a “somewhat Aristotelian” perspective (Sen, 2002),⁵⁵ that defines “progress” in terms of “quality of life” – in terms, that is,

⁵³ Cf. von Hayek (1991), pp. 21–22. See more in general, on the limits of individuals, as these emerge from the representations of the theorists of methodological individualism Lukes (1973), p. 153.

⁵⁴ Hayek equates the adherence to the abstract rules of the Great Society with the response of animals to stimuli according to conditioned reflexes. See von Hayek (1973), pp. 74 et seq.

⁵⁵ Sen here alludes to his theory of capabilities. He defines the quality of life in terms of capability to choose and to achieve modalities of doing and being to which he attributes a value (these

of what we *can* do or be. Adopting this point of view, it is by no means obvious that the conditions that facilitate survival and reproduction, in a given context, contribute to making lives more pleasant and more satisfying: “we recognize many virtues and achievements that do not help survival but that we have reason to value; and on the other side, there are many correlates of successful survival that we find deeply objectionable” (Sen, 2002, p. 494).

Not dissimilarly, during another time of globalization, William James wrote: “The entire modern deification of survival *per se*, survival returning to itself, survival naked and abstract, with the denial of any substantive excellence in what survives, except the capacity for more survival still, is surely the strangest intellectual stopping-place ever proposed by one man to another” (James, 1987, p. 359).

4.6 Alchian: Evolutionism as Connecting Link Between Two Laissez Faire Strategies

Both Friedman and Hayek are important exponents of monetarism. This doctrine argues from the markets’ capacity for self-regulation to conclude that the sole source of instability is the wayward variation of the level of prices. This variation is ascribed to a single cause, the “abnormal” change in the money supply. It belongs, however, within two different “rhetorical strategies” used to argue for *laissez faire* (Denis, 2004). At the level of method, Friedman takes a position summed up as “positive” or “instrumentalist” (or “as if”) economics, which at first sight looks a long way from Hayek’s “apriorism.”

According to Friedman, theories are not attempts to describe reality, but tools that enable forecasts to be made. The validity of a theory is not measured, therefore, by the realism of its assumptions but by its predictive capacity: if a theory works – that is, it generates “sufficiently accurate predictions” – the “constructed hypothesis is presumably valid” (Friedman, 1953, p. 20). Friedman, as was said, stands in the Walrasian neoclassical stream. In this “antinormative” perspective, the State is not expected to supply positive contributions to production and society, but effects only distortions or interferences in a mechanism capable of perfect self-regulation; it is therefore possible immediately to draw *general* conclusions (conclusions independent, i.e., of a specific analysis of the situation under consideration) about the advantage of a reduced presence of the State.

modalities are called functionings). This theory differs from utilitarian analyses, which attribute a value to goods or actions on the basis of the utility they procure; and from conceptions that do not attribute a value to liberty of choice. For how this theory relates to Aristotle’s thought (cf. in particular *Nicomachean Ethics*, Book I, §7, and *Politics*, Book III) see Nussbaum (1988), p. 153 and Sen (1993), pp. 46–47. The same authors have evidenced links between the perspective of capabilities and conceptions of Adam Smith (*The Wealth of Nations*) and Marx (*Economic-Philosophical Manuscripts of 1844*). Concerning Smith, see Sen (1993), p. 46; and on Marx, see Nussbaum (1988), pp. 103–104 and Sen (1993), p. 46 and the bibliography therein.

It may be wondered, however, in what way Friedman manages to avoid dangerous “normative” uses of the Walrasian scheme. He does so in two ways. The first concerns the refusal to specify in detail the mechanism of transmission of the money supply to the economic system and, in particular, how, in the short term a determinate variation in the nominal income – induced by a variation in the money supply – is subdivided between variations in the quantities of output and the price level.

This indeterminacy in the monetarist theory – that has led some critics⁵⁶ to speak of “black box theory” – is justified by adducing the methodological stance adopted. Given the complexity of the economic system and the ignorance regarding the multiplicity of individual behaviors, it entails that only simple and general propositions undergo the prediction test, abstracting from the details. There is an analogy with the “explanation of principle” of Menger and Hayek, but in this case the falsifiability of the propositions is jeopardized by the high level of aggregation, which does not enable one to substantiate possible inversions of the causal nexus between the variables (for example, from nominal income to money supply, and not *vice versa*), or between reciprocal interactions.

The second way entails introducing a very specific hypothesis on a non-observable variable, the expectations of inflation of the economic actors. Consistently with the Walrasian scheme, it is assumed that unemployment is exclusively voluntary, depending on the workers’ preferences for free time. Should the government authorities intend to apply an expansionary monetary policy, they could obtain an increase in the amount of labor available on the market by deceit. For, suppose that the increase in the money supply is followed by an increase in the amount of labor demanded by firms, at a higher nominal wage. Suppose, further, that the workers project past expectations of (low or nil) inflation into the future: the workers (unaware of what awaits them) offer more labor since they think that the increased money wage will also translate into an increase in the real wage. But this is an illusion: prices will soon begin to soar and the expansionary monetary policy will have to be promptly abandoned.⁵⁷ In conclusion, Friedman finds the Keynesian policy of reduction in (involuntary) unemployment as, in the words of Parboni (1984, p. 202), “a watered-down version of authoritarian socialism” which forcibly (and not by deceit, as here) compels people to work in poorly productive jobs.

The same conclusions are reached by Hayek, though by a different and more tortuous route. Hayek’s “individualistic organicism” leads him to argue that the social economic system has behaviors and dynamics that do not reflect the intentions and actions of the single individuals. *However*, as we saw in the previous section, we can be confident that the overall outcome of the actions of the individuals will at least be satisfactory (perhaps even optimal). It may be of interest here to clarify the role played by evolutionism in the conceptions of the two exponents of monetarism.

⁵⁶ P. Samuelson (1970), p. 43 has written: “My most serious objection to . . . monetarism is that it is a black box theory”. On this topic see Parboni (1984), p. 153 and note 4.

⁵⁷ In order to avoid waves of instability of price levels caused by the monetary policy, Friedman has proposed a policy of constant expansion of the money stock. Hayek, on the contrary, suggested that the same objective might be achieved by privatizing the production of money.

With reference to an influential article by Alchian (1950) based on the analogy between the natural selection process and the selection process actuated by competition (which would be taken up by Friedman in 1953), Foster (1997, p. 9) writes that “it would be no exaggeration to affirm that the importance of neoclassical theory for real economic policies was justified by the implicit conviction that neo-Darwinian competitive forces are continually in action from the supply side of the economy.”

Alchian’s article is interesting also because it implies a connecting link between the two rhetorical strategies we are dealing with. It intervenes in a discussion in which certain authors (Hall & Hitch, 1939; Lester, 1946), based on questionnaires given to entrepreneurs, called into question the hypothesis that firms maximize profit, instead of following empirical criteria such as the “full cost principle” (a fixed mark-up on the unit costs). Alchian allows that in presence of uncertainty one cannot take the maximization of profit as a guide to the action of a particular firm. If uncertainty prevails, the expected outcome of any action by a firm can only be viewed as a distribution of possible outcomes, and it makes no sense to say that the firm maximizes anything, since it is impossible to maximize a distribution. He holds, however, that the introduction of the criterion of survival enables one to conclude that the maximization of profit is a valid generalization of the behaviors of firms (even though this conclusion concerns the overall firms in an industry,⁵⁸ and not the individual firm).

Faced with a change in the environment, conditions being equal, the firms that are lucky enough to use the optimal *routine* (that is, those adopting the criterion of maximization of profit) will survive, whereas the firms that adopt sub-optimal criteria will be eliminated. What counts for survival is the result (the profit realized *ex post*), not the motivation that the firms had *ex ante*). Hence, industry as a whole moves towards an optimal configuration, not because it *adapts* by changing its behavior versus external shocks, but because the market (the competition) *adopts* (selects) the firms that have by chance followed an optimal *routine* and leaves those that have not done so to perish. The achievement of profit is the result, not of human design (the conscious decision by the individual firm), but of human action (the evolutionary process that performs, in an impersonal way, a selection at the level of the system). It must be emphasized that the evolutionary process selects not only the agents who by chance follow optimal strategies but, in the long term, also the behavioral norms that enable these strategies to be achieved. The pictured traced by Alchian closely recalls Hayek’s interpretation of Ferguson.

It is of interest to note that, in an essay in which he enunciates the “as if” methodology, Friedman on the one hand declares himself in agreement with Alchian’s analysis, while on the other he argues that the hypothesis of the maximization of profit is justified by the fact that, if entrepreneurs did not adopt it, they would soon lose resources or would fail. This argument concerns the expected profits, not the *ex post* ones, and, as against Friedman’s professed support for Alchian’s analysis, seems

⁵⁸ Here, too, the “explanation of the principle” replaces that of the detail, with all its “dangerous” analytic and normative implications.

to accept the very thesis that Alchian had rejected.⁵⁹ But Friedman (1953, *italico added*) goes on to say (somewhat ambiguously): “The process of ‘natural selection’ thus helps to validate the hypothesis – or, rather, given natural selection, acceptance of the hypothesis [of maximum profits] can be based largely on the judgment that it summarizes appropriately the conditions for survival.” Friedman’s contradiction is, however, only apparent: in the scheme adopted, which does not envisage conditions of uncertainty, the criterion of maximization of expected profit and the conditions of *ex post* survival following the “optimizing” process of selection exerted by perfect competition coincide with one another, so that the choice between the two is irrelevant.

As Edith Penrose remarked in an acute comment on Alchian’s article – a comment of more general relevance than the occasion that engendered it – those who employ the biological metaphor in dealing with economic topics have a common characteristic: that of suggesting “explanations of events that do not depend upon the conscious decisions of human beings” (Penrose, 1952). This appears all the more singular as certain branches of biology do in fact deal with processes of learning and decision making, and with “purposive motivation and conscious choice”. The analogies that invoke biology in treating of aspects of economics that do not bear on human motivations and decisions tend to relegate into the background the fact that firms are created by people and serve the purposes of human beings. Penrose (1952, p. 809) writes:

... to abandon [the] development [of the firms] to the laws of nature diverts attention from the importance of human decisions and motives, and from problems of ethics and public policy, and surrounds the whole question of the growth of firms of an aura of ‘naturalness’ and even inevitability.

When these decisions and motivations have been included, any analysis of the effects of a change in environment must take account of the fact that firms will seek “as much consciously to adapt the environment to their own purposes as to adapt themselves to the environment. After all, one of the chief characteristics of man that distinguishes him from other creatures is the remarkable range of his ability to alter his environment, or to become independent of it.” While certain aspects of the environment (an imprecise notion in Alchian’s analysis) appear hard to modify, there are some that are “important for survival which we cannot assume are beyond the

⁵⁹ In this connection, see the critical remarks by Kay (1995). But in a note, Friedman (1953), in the essay quoted, argued that, given uncertainty, the choice among “alternative anticipated probability of receipts or income” depend “on the criteria by which they are supposed to be ranked.” Friedman suggest to “rank probability distributions by the mathematical expectation of the money receipts corresponding to them,” adding, though, that the “methodological issues” involved here were “largely by-passed” in his discussion. In this way, uncertainty actually disappears, since past and future are perfectly fungible and to every possible outcome is associated to a well-defined utility: Alchian’s scheme is returned to the haven of orthodoxy, and the maximization of profit is tacitly reintroduced as the *ex ante* condition for survival. For Alchian, instead, in conditions of uncertainty, the sufficient condition for survival was the *ex post* presence of positive (not necessarily maximum) profits. Eventually, among those attaining *ex post* positive profits, competition would select the “luckiest” (i.e. maximizers) firms.

influence of firms and which can be unpredictably altered by them” – for example, the state of technology and consumer preferences. The requirements for survival, i.e. the profits “attained,” were stated by Alchian on the basis of the condition of *coeteris paribus* (Penrose, 1953, p. 608), but they would be very different if one granted the possibility of action by the firms consciously aimed at altering the conditions of the environment in which they operate.

4.7 The Black Box

In both the exponents of monetarism, evolutionism entails ascribing characteristics of optimality and/or efficiency to the outcome of the competitive processes. In Hayek’s scheme as in Friedman’s one there is a mechanism that is impenetrable to analysis, a “black box.” In the past, Providence ensured the optimality of the social order that conformed to natural law.⁶⁰ Today, the “black box,” by analogy with natural selection (not Darwin’s, in actual fact, but Spencer’s), is competition.

The rhetoric of biological metaphor in economics employed in the first great expansion of the global market in the second half of the 19th century and the one employed in the second phase subsequent to 1980 display many similarities. The uncertainty of boundaries seems to lead not only to a loss of identity but also to the assertion of ideologies at once deterministic and consolatory, not to say – at the level of system – optimistic. In their overbearing simplicity, they invite one to *adapt* to an “environmental” situation viewed as exogenous, in which any intervention, any interference with the process of selection, would be *a priori* counterproductive. As devices to contrast any attempts consciously to alter the *status quo*, they may generally be associated with conservative thought. Quite absent, however (though not entirely so in the case of Hayek, who strives to reconnect with Burke), are other characteristic components of conservative thought, and this fact introduces an element of potential dissonance and instability. As Hofstadter has remarked with reference to the social Darwinism of the second half of the 19th century, we have here “a body of belief whose chief conclusion was that the positive functions of the state should be kept to the barest minimum, it was almost anarchical, and it was devoid of that center of reverence and authority which the state provides in many conservative systems” (Hofstadter, 1955, p. 7). In substance, a “conservatism without religion.” It appears to be terribly difficult to accept that we are part of a social process without direction, without meaning. But once launched on the road of teleology, in its multiple guises, Providence comes to fill up the empty space – and is just round the corner.

⁶⁰ As against the “secularized” interpretation now prevailing, components of natural theology, in Smith’s *Wealth of Nations*, are complementary, not in contrast with the conception of the social order based on the “invisible hand.” See Viner (1972), pp. 81 et seq. and Hill (2001) and the ensuing discussion in the same journal.

Acknowledgments This article is greatly indebted to David Lane. I have learnt much from his essays and from the frequent discussions chaired by him as part of the European ISCOM project, and Section 4.3 bears abundant witness to this. However, he should not be held accountable for what I argue. I am also deeply grateful for helpful comments, criticism and discussion from Lorenzo Bianchi, Giovanni Bonifati, Antonio Ribba and Fernando Vianello. The article was supported by funding from MUIR.

References

- Alchian, A. (1950). Uncertainty, evolution and economic theory. *Journal of Political Economy*, 58, 211–221.
- Bannister, R. C. (1979). *Social Darwinism. Science and myth in Anglo-American social thought*. Philadelphia, PA: Temple University Press.
- Beckner, M. (1968). *The biological way of thought*. Berkeley, CA: University of California Press.
- Berlin, I. (1969) *Two concepts of liberty*. In *Four essays on liberty*. Oxford, UK: Oxford University Press.
- Bonifati, G. (2003). *Concorrenza come processo di interazione, agenti e sistemi del mercato. A note from the Information Society as a Complex System (ISCOM) discussion group*. Department of Social, Cognitive and Quantitative Sciences, University of Modena and Reggio Emilia.
- Caldwell, B. J. (1988). Hayek's transformation. *History of Political Economy*, 20(4), 513–541.
- Darwin, C. (1859). *On the origin of species by means of natural selection, or the preservation of favoured races in the struggle for life*. In *The collected works of Charles Darwin* (Vol. 16). London, UK: William Pickering, 1988.
- Darwin, C. (1887). *The life and letters of Charles Darwin, including an autobiographical chapter* (Vol. I, F. Darwin, Ed.). London, UK: John Murray Publishers.
- Darwin, C. (1871). *The descent of man, and selection in relation to sex*. London, UK: John Murray Publishers.
- Darwin, C. (1958). *The autobiography of Charles Darwin 1809–1882 with original omissions restored* (N. Barlow, Ed.). London, UK: Collins.
- Denis, A. (2004). Two rhetorical strategies of laissez-faire. *Journal of Economic Methodology*, 3, 341–357.
- Donzelli, F. (1988). Introduction to F.A. von Hayek, *Conoscenza, mercato, pianificazione*. Bologna: Il Mulino.
- Forsyth, M. (1988). Hayek's bizarre liberalism, a critique. *Political Studies*, 36(2), 235–250.
- Foster, J. (1997). The analytical foundations of evolutionary economics: From biological analogy to economic self-organization. *Structural Change and Economic Dynamics*, 8(4), 427–451.
- Friedman, M. (1953). The methodology of positive economics. In M. Friedman (Ed.), *Essays in positive economics* (pp. 3–43). Chicago, IL: University of Chicago Press.
- Garegnani, P. (1976). On a change in the notion of equilibrium in recent work on value: A comment on Samuelson. In M. Brown, K. Sato, & P. Zarembka (Eds.), *Essays in modern capital theory*. Amsterdam, New York: North Holland.
- Gould, S. J. (1980). *The Panda's Thumb: More reflections on natural history*. New York: Norton.
- Gould, S. J. (1982). The meaning of punctuated equilibrium and its role in validating a hierarchical approach to macroevolution. In R. Milkman (Ed.), *Perspectives on evolution* (pp. 83–104). Sunderland, MA: Sinauer.
- Gould, S. J. (2002). *The structure of evolutionary theory*. Cambridge, MA: Belknap Press of Harvard University Press.
- Gould, S. J., & Lewontin, R. C. (1979). The spandrels of San Marco and the panglossian paradigm: A critique of the adaptationist program. *Proceedings of the Royal Society of London B. Biological Science*, 205, 581–598.
- Hall, R. E. & Hitch, C. (1939). Price theory and business behaviour. *Oxford Economic Papers*, 2(1), 12–45.

- Hawkins, M. (1997). *Social Darwinism in European and American thought 1860–1945: Nature as model and nature as threat*. Cambridge, UK: Cambridge University Press.
- von Hayek, F. A. (1940). Socialist calculation: The competitive solution. *Economica*, 7(26), 125–149, reprinted in von Hayek, F. A. (1948). *Individualism and economic order*. Chicago, IL: The University of Chicago Press.
- von Hayek, F. A. (1945). The use of knowledge in society. *American Economic Review*, 35(4), 519–530, reprinted in von Hayek, F.A. (1948). *Individualism and Economic Order*. Chicago, IL: The University of Chicago Press.
- von Hayek, F. A. (1948). The meaning of competition. In *Individualism and economic order* (pp. 92–106). Chicago, IL: The University of Chicago Press.
- von Hayek, F. A. (1952). *The sensory order: An inquiry into the foundations of theoretical psychology*. London, UK: Routledge.
- von Hayek, F. A. (1960). *The constitution of liberty*. Chicago, IL: University of Chicago Press.
- von Hayek, F. A. (1978). *New studies in philosophy, politics, economics and the history of ideas*. London, UK: Routledge.
- von Hayek, F. A. (1973). *Law, legislation and liberty, Vol. I: Rules and order*. Chicago, IL: University of Chicago Press.
- von Hayek, F. A. (1976). *Law, legislation and liberty, Vol. III: The political order of a free people*. Chicago, IL: University of Chicago Press.
- von Hayek, F. A. (1991). *The fatal conceit: The errors of socialism* (W. W. Bartley III, Ed.). Chicago, IL: University Of Chicago Press.
- Hill, L. (2001). The hidden theology of Adam Smith. *European Journal of the History of Economic Thought*, (8)1, 1–29.
- Hodge, M. J. S., & Kohn, D. (1985). The immediate origins of natural selection. In D. Kohn (Ed.), *The Darwinian heritage* (pp. 185–206), Princeton, NJ: Princeton University Press.
- Hodgson, G. M. (1993). *Economics and evolution: Bringing life back into economics*. Ann Arbor, MI: University of Michigan Press.
- Hofstadter, R. (1955). *Social Darwinism in American thought, 1860–1915*. Boston, MA: Beacon Press.
- Ingrao, B., & Israel, G. (1990). *The invisible hand. Economic equilibrium in the history of science* (I. McGilvray, Trans.). Cambridge, MA: The MIT Press.
- James, W. (1987). *Essays, comments and reviews* (F. H. Burkhardt, F. Bowers, & I. K. Skrupskelis, Eds.). Cambridge, MA: Harvard University Press.
- Kay, N. M. (1995). Alchian and the ‘Alchian Thesis.’ *Journal of Economic Methodology*, 2(2), 281–286.
- Keynes, J. M. (1971). The general theory of employment, interest and money. In *The collected writings of John Maynard Keynes* (Vol. VII). London, UK: St Martin’s Press [first edition, 1936].
- Kohn, D. (1989). Darwin’s ambiguity: The secularization of biological meaning. *British Journal for the History of Science*, 22, 215–239.
- La Vergata, A. (1985). Images of Darwin: A historiographic review. In D. Kohn (Ed.), *The Darwinian heritage* (pp. 901–972). Princeton, NJ: Princeton University Press.
- La Vergata A. (1990a). *L’equilibrio e la guerra della natura. Dalla Teologia naturale al darwinismo*. Naples, Italy: Morano.
- La Vergata, A. (1990b). *Nonostante Malthus. Fecondità, popolazioni e armonia della natura, 1700–1900*. Torino, Italy: Bollati Boringhieri.
- Lakoff, G., & Nuñez, G. E. (2000). *Where mathematics comes from: How the embodied mind brings mathematics into being*. New York: Basic Books.
- Lane, D. A. (2002). Complexity and local interaction: Towards a theory of industrial districts, in complexity and industrial clusters (A. Quadrio Curzio & M. Fortis, Eds.). Heidelberg-New York: Physica-Verlag.
- Lester, R. A. (1946) Shortcomings of marginal analysis for wage-employment problems. *American Economic Review*, 36, 62–82.
- Lukes, S. (1973). *Individualism*. Oxford, UK: Blackwell.

- Maione, G. (2003). *Le Merci Intelligenti: Miti e Realtà del Capitalismo Contemporaneo*. Milan, Italy: Bruno Mondadori.
- Malthus, T. R. (1798). *An essay on the principle of population, as it affects the future improvement of society, with remarks on the speculations of Mr. Godwin, M. Condorcet, and other writers*. London, UK: J. Johnson.
- Malthus, T. R. (1979). *An essay on the principle of population; and, a summary view of the principle of population* (A. Flew, Ed.). Harmondsworth, UK: UK Penguin Books.
- Manier, E. (1978). *The young Darwin and his cultural circle*. Dordrecht, The Netherlands: D. Reidel Publishing House.
- Marx, K. (1962). *Capital: A critique of political economy*, vol. III. Moscow: Foreign Languages Publishing House.
- Marx, K. (1969). *Theories of surplus value*, vol. 2. London: Lawrence & Wishart.
- Maxfield, R., & Lane, D. A. (1997). Foresight, complexity and strategy. In W. F. Arthur, S. Durlauf, & D. A. Lane (Eds.), *Economy as a complex, evolving system II*. Redwood City, CA: Addison Wesley.
- Mayr, E. (1982). *The growth of biological thought. diversity, evolution and inheritance*. Cambridge, MA: The Belknap Press of Harvard University Press.
- Menger, C. (1871). *Grundsätze der Volkswirtschaftslehre*, Vienna, Austria: Wilhelm Braumüller, Authors Translation.
- Menger, C. (1883). *Untersuchungen über die Methode der Sozialwissenschaften und der politischen Ökonomie insbesondere*. Leipzig Drucker & Humblot; L. Schneider (Ed.), *Problems of economics and sociology* (in English). Urbana, IL: University of Illinois Press, Urbana, 1963.
- Nardozzi, G. (2004). *Miracolo e declino. L'Italia fra concorrenza e protezione*. Bari, Italy: Laterza.
- Needham, R. (1975). Polythetic classification: convergence and consequences. *Man*, 10, 349–369.
- Nietzsche, F. (1967). *On the genealogy of morals: A polemic*. New York: Vintage Books.
- Nussbaum, M. C. (1988). Nature, function and capability: Aristotle on political distribution. *Oxford Studies in Ancient Philosophy* (Suppl. Vol., pp. 145–188). Oxford, UK: Clarendon Press.
- Ospovat, D. (1979). Darwin after Malthus. *Journal of the History of Biology*, 12, 211–230.
- Parboni, R. (1984). *Moneta e Monetarismo*. Bologna, Italy: Il Mulino.
- Paul, E. F. (1988). Liberalism, unintended orders and evolutionism. *Political Studies*, 36(2), 251–272.
- Penrose, E. (1952). Biological analogies in the theory of the firm. *American Economic Review*, 42(5), 804–819.
- Penrose, E. (1953). Biological analogies in the theory of the firm: Rejoinder. *American Economic Review*, 43(4), 603–609.
- Picchio, A. (1992). *Social Reproduction. The Political Economy of the Labour Market*. Cambridge, UK: Cambridge University Press.
- Ricardo, D. (1973). Letters 1819 – June 1821. Vol. VIII of P. Sraffa and M. Dobb (Eds.), *The works and correspondence of David Ricardo*. Cambridge, UK: Cambridge University Press.
- Richardson, G. B. (1972). The organisation of industry. *Economic Journal*, 82(327), 883–896.
- Rorty, R. (1979). *Philosophy and the mirror of nature*. Princeton, NJ: Princeton University Press.
- Rosenberg, N. (1976). *Perspectives on technology*. Cambridge, UK: Cambridge University Press.
- Samuelson, P. (1970). Reflections on the recent Federal Reserve policy. *Journal of Money, Credit and Banking*, 2(1), 33–44.
- Schweber, S.S. (1977). The origin of the *Origin* revisited. *Journal of the History of Biology*, 10(2), 229–316.
- Schumpeter, J. A. (1992). *Capitalism, socialism and democracy*. London, UK: Routledge.
- Simmel, G. (1908). *Soziologie*. Leipzig, Germany: Dunker & Humblot.
- Sen, A. K. (1993). Capability and well-being. In A. K. Sen & M. C. Nussbaum (Eds.), *The quality of life*. Oxford, UK: Clarendon Press.
- Sen, A. K. (2002). *On the Darwinian view of progress* (1991). In A. K. Sen, *Rationality and freedom*. Cambridge, MA: Belknap Press of Harvard University Press.

- Smith, B. A. (1986). Austrian economics and Austrian philosophy. In W. Grassl & B. Smith (Eds.), *Austrian economics: historical and philosophical background*. London: Croom Helm.
- Smith, B. A. (1990). Aristotle, Menger, Mises: an essay in the metaphysics of economics. In B.J. Caldwell (Ed.), *Carl Menger and his legacy in economics*. Annual Supplement to *History of political economy* (Vol. 22). Durham, NC: Duke University Press.
- Sober, E. (1981). *Holism, individualism and the units of selection*. In E. Sober (Ed.), *Conceptual issues in evolutionary biology: An anthology*. Cambridge, MA: MIT Press.
- Viner, J. (1972). *The role of providence in the social order*. Princeton, NY: Princeton University Press.
- Wittgenstein, L. (1975). Note sul Ramo d'oro di Fraser. Milano: Adelphi. Italian translation of Bemerkungen über Frazers The Golden Bough, *Synthese*, XVII, 1967.
- Young, R. M. (1985). Darwinism is social. In D. Kohn (Ed.), *The Darwinian heritage* (pp. 609–638). Princeton, NJ: Princeton University Press.

Chapter 5

Innovation in the Context of Networks, Hierarchies, and Cohesion

Douglas R. White

5.1 Introduction

This chapter attempts to bridge two different worlds, that of substantive social science theory and that of formal mathematical theory. It therefore has to be read and understood from both these perspectives simultaneously. Substantively speaking, this chapter explores how diverse, multi-level, sparsely (or densely) interconnected, complicated, and loosely (or tightly) integrated the structures and network processes involved in innovation may be. To study such intricate systems and processes, with very different forms and types and degrees of complexity, we need to reconsider and reconfigure the very basic concepts we use. That is where formal theories come in. The other goal of this chapter is to further the integration of formal theories of network dynamics with substantive theories of socio-historical dynamics and to outline how certain types of innovation play out within these dynamics.

Some readers will inevitably ask, “Why make things so difficult?” The reason is that by using substantive and formal theory in tandem, we can generate bodies of data that have substantive relevance but can be coded in formal and relational terms, enabling, on the one hand, potentially explanatory formal theories to be substantively tested, and, on the other, allowing the incorporation into substantive theory of proofs of outcomes that follow logically from definitions.¹ Such proofs, in turn, may help us to grasp complex relational phenomena in ways that are otherwise difficult to construct with any great precision within substantive theory.

The two worlds we are talking about are also distinguished by the fact that the social sciences have focused on static, structural descriptions of social organization that were primarily concerned with the *position* of individuals in the organization,

D.R. White (✉)

School of Social Sciences, University of California, Irvine, UK

¹ Michael Kearns et al. (2007) has shown in behavioral experiments that networks that are very hard to solve optimally, from a computer science/algorithmic perspective, were relatively easy for subjects paid according to each person’s performance, while those networks that were algorithmically easy were appreciably harder for the subjects to do well. The significance is that concepts like structural cohesion relevant to substantive problems should not be avoided just because they are computationally “hard” or mathematically “hard to understand.”

whereas the kind of formal approach we propose is, in essence, dynamical and focuses on *relationships* between individuals to understand the *organizations* that emerge and their evolution. In terms of Chapter 1 of this book, the social sciences traditionally *looked at the distribution of populations, power, riches, etc. in order to describe and categorize organizations* ('population thinking'), whereas the formal approach to which we refer here focuses on *looking at the underlying dynamics of organizations to understand relationships between individuals* ('organization thinking').

If we consider that events occur in a world in motion, inherently relational concepts benefit enormously from being well specified and measurable. This is particularly crucial for the study of the processes that affect innovation and for the study of change in general.

Our perspective aims to lay out conceptual and analytic foundations iteratively for a specific kind of relational social theory ("network theory"), while at the same time providing commensurate tools for the analysis of empirical phenomena. We do so by identifying entities and their various attributes and relations through the study of their appearance, transformation, and disappearance at different scales of resolution, rather than predefining entities and then studying how they relate. This allows us to identify the dynamics of interactions and to clarify how formal aspects of structures and movements operate in relation to one another. However, the intertwining of these two aims does not make the chapter easy to read. We have, therefore, decided to use this introduction to present you with a 'roadmap' of the paper.

The Section 5.2 of this chapter begins by clarifying specific formal definitions, used in network theory, that enable the study of forms of *hierarchy* and *cohesion*, and their various empirical manifestations to clarify and measure different aspects of the relational nature of agents and agency and their connections.²

In the Section 5.2.1, we present a general formal model of hierarchy that might be instantiated in very different substantive contexts. When we instantiate it, there immediately arises a multiplicity of substantive phenomena that could give that model very diverse contents. Conversely, multiple substantive phenomena could be mapped into the same formal model of hierarchical relations. The beauty of relational models is that they allow us to analyze such pluralities. They examine the richness of the many dimensions and structures involved in different contexts of interaction at multiple spatiotemporal scales. This enables us to improve our understanding of the dynamics involved and the consequences of these dynamics

² In mathematics and in physics, the concepts associated with the term "hierarchy" are well specified, conveying the precise relational structure of an ordered or pre-ordered set. They can also distinguish between types of cases and levels. Similarly, the concepts associated with the notion of "cohesion" specify and measure the specific strength and nature of bonds and linkages. But they also provide information on the consequences of the disconnection or removal of an element from other elements in the structure – and this is key for the study of networks and the measurement of the relationships within them.

for the outcomes of complex interactive processes – including the understanding of innovation and change at scales both large and small.³

The following sections illustrate different kinds of hierarchical structures from the (anthropological, archaeological, sociological, and business) literature in (formal) terms of varieties of *hierarchical structure*, *routes of traversal* (through which matter, energy and information flow throughout the network), and the *specific attributes of individual nodes and directions, qualities, and weightings of the linkages between them*.

Several of these examples link various manifestations of structural cohesion as a formal concept with substantive measures. They provide a variety of examples and arguments for why we should expect to find cohesive structures among the substantive underpinnings of hierarchy, and how structural cohesion is necessarily involved in the very construction of certain types of hierarchies. The aim along the way is to resolve some of the ambiguities of hierarchy by examining the distribution of attributes viewed as the outcome of network processes.

These sections flesh out some of the disparate evolutionary trajectories and discontinuities that are inherent in network phenomena. The empirical and historical examples for the existence of social hierarchy and its role in innovation that we present, vary from simple and small-scale to large-scale, complex urban systems. As the chapter proceeds, the reader will see that this approach also allows us to explore networks conceived as a multiple layering of social phenomena in overlapping but ambiguous hierarchies and the structurally cohesive groups that might support them. The possible ambiguity of what might be understood in any particular social context by an overlay of different hierarchies (along with other forms of relations), emphasizes the importance of understanding such systems in terms of their genesis and the outcomes of the complex interactions that are responsible for them.

Therefore, Section 5.2 ends with a discussion of validation criteria relevant to interpretation. It considers how the consequences of multiple processes – multiple outcome variables resulting from network interaction – can help us make more precise inferences about these processes.

Section 5.3 looks at the relationship between structure and dynamics. It begins by discussing Menger's Multi-Connectivity Theorem, in which the hierarchical and traversal aspects of cohesion unite to form a single concept ('*Structural cohesion*') with two independently measurable aspects: (a) degrees of invulnerability to node removal or outside attack on members that will disconnect the group (structurally) and (b) ease of (traversal) communication and transports within the group. It turns out that not only are (a) and (b) isomorphic, but that structural cohesion is scalable without increase in overhead, which enables the genesis of very large networks that are relatively invulnerable to disconnection until they encounter escalating

³ At a detailed level of analysis, for example, aspects of a narrative approach can be intrinsically re-conceptualized and then recoded into selective modes that yield to network approaches that restrain the lens of theory, even if partially, to focus on general aspects of explanatory frameworks and to clarify how empirical specificities are defined and measured in ways that enable feedback from relevant and possible tests of theory.

opposition. Such opposition is based on the same very rapid scalability of structural cohesion among initially under-connected sets. Examples range from high-school groups and identities in the US, to social class formation, to the expansion and demise of Empires.

Next, structural cohesion is linked to innovation, and, in particular, to ‘innovation waves’ by pointing out that high levels of structural social cohesion and interaction among people with similar interests in innovation occur rarely, but, when they do occur, the diffusion of ideas within the group and the success of its members are significantly enhanced compared to isolated individuals or groups with lesser cohesion.

An example in the biotech industry is then used to illustrate the theorem that structurally cohesive groups in a hierarchy are always contained in the same (or a larger) set of nodes that are structurally cohesive at the next lower level, thus creating cohesive network hierarchies. Cities effectively consist of multiple (partly) overlapping structurally cohesive hierarchies. That characteristic allows them to scale up almost indefinitely, up to limits of collapse.

In the section “*Networks affecting hierarchy: the scale-up of city networks*,” a novel perspective is presented on the relationship between the exponential world population growth prior to, and power-law growth dynamics of cities in the last two millennia. With the up-scaling of “multi-connectivity” through the spread of structural cohesion in large groups of the population, the flows of matter, energy, and information can now reach further and further, through trade, through the diffusion of technological innovation, and through attraction to cities, which allows higher “storage” of dense populations while the (higher) reproductive rates of rural areas replace and compensate for the out-migrants to cities.

It is then argued that the inherent instabilities in the network dynamics linking cities over long distances cause regional cycles of “power-law growth spurts” alternating with periods of relative stability, and that these cycles are linked to each other across continents. Hence, we need to apply a co-evolutionary perspective to cities and urban networks. When such a perspective is applied to size distributions of cities across Eurasia between 500 and 2000 CE, it appears that there are non-random oscillations between values that alternate significantly below or above the Zipf curve.

Section 5.4 tentatively searches for the causalities behind the dynamics thus outlined. It approaches this search first by focusing on the time lags involved between different, far away, parts of the system, but also between the evolution of larger and smaller cities. These indicate that technologically advanced areas are, in effect, diffusion sources for less advanced ones, and that the carriers of the correlations (and of innovative technologies) are the trade networks (and technology of war).

Borrowing from Turchin (2005a, 2005b), the Section 5.4.1 looks at density-dependent changes in population/resource pressure and socio-political instabilities in China and in Europe, to distinguish endogenous from exogenous dynamics in the context of a two-equation time series. The dynamics turn out to be dominated by a combination of endogenous cycles and – in work that goes beyond Turchin – exogenous shocks from major warfare events. A further prediction from

the endogenous part of the model is that high rates of innovation occur in a relatively closed region, in the periods *following a downturn* after which total population is rising faster than resources. The actual innovations implemented in these periods often use inventions that trace back to periods in earlier cycles as well as trade networks expanded earlier.

This last section ends with a case study of innovation cycles in the contemporary U.S. that points out how, late in a cycle, the cross-links between innovations in different sectors are subject to hierarchical dynamics that keep these sectors separate, rather than fostering interactions that would lead to a new wave of innovation.

5.2 Understanding Network Structure

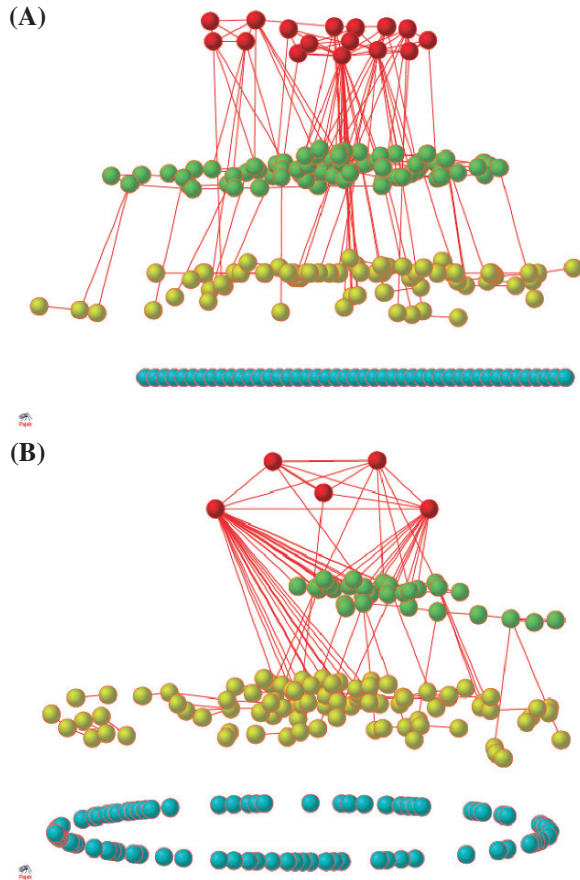
5.2.1 *Hierarchy and Reverse Hierarchy*

When a network consisting of a single kind or set of interpersonal linkages is portrayed as a graph, consisting of nodes and directed or bi-directed edges between them, the generic form of hierarchy is described as a *directed asymmetric graph*, or *dag*, where every edge is directed (asymmetric, nonreciprocal – or, as a refinement, at least between different levels) and there are no directed cycles (see Fig. 5.1A, omitting cycles within a level). The extent to which particular observed networks approximate the properties of a *dag* can be examined for different kinds of relations and contexts: inter-group, interpersonal, inter-positional, and inter-site relations, for example. A theorem for *dag* structures is that the nodes can be ordered in a minimal number of levels so that every directed edge orients in a consistent direction (Fig. 5.1B).

If a social network has the *dag* and a pyramidal structure with fewer nodes at the top and more toward the bottom (beware: not all *dag* structures have this property), then it defines a *social hierarchy*. The usual conception of a social hierarchy is that those (fewer) nodes at the “top” have greater power, prestige, income, information, expertise, etc. Such *dominance hierarchies* appear in many species as one of several basic and often co-occurring forms of social hierarchy. In network terms, individuals in such a hierarchy dominate others in acyclic directed chains, that is, *dags*, often with additional structural properties, such as density of links between different pairs of levels.

Another organizational form involving tendencies toward acyclic directed chains is the *reverse hierarchy*, identified by anthropologists, drawing on both field experiments and a huge array of ethnographic examples. Such hierarchies, based on the recognition of achievements or prestige, seem to be unique to humans and a crucial part of human culture (Henrich, 2001; Henrich et al., 2005). The reverse hierarchies reviewed by Boehm (1993, 2001) include altruistic relations such as sharing, aid, and redistribution. They might be thought of as hierarchies in which there are more obligations than privileges, such as a hierarchy of philanthropists in which those at the top are the ones who give most.

Fig. 5.1 Chinese women migrants' network discussing children, contraceptives. **(A)** Levels of advice are directed downward in the hierarchy; **(B)** levels of advice have laterally symmetricities



Flows of benefits and liabilities in terms of different overlapping network hierarchies form rather complex systems that are often difficult to unravel: are they hierarchies, reverse hierarchies, a mixture of both, or something else entirely? Do they have conditions for stability or instability that depend on tradeoff of benefits between different ranks or levels? Is inter-level mobility based on achievement and tradeoff of benefits or unfair competition, monopolies of power, and structural inequality?

To deal with these ambiguities, we can conceptualize networks and the dynamics in them simultaneously in terms of three fundamental characteristics: their *structure* (e.g., *dag*, pyramidal), their *routes of traversal* through which (directed or two-way) flows of information, materials, or energies are channeled, and the *specific attributes of individual nodes and of the linkages between them*. “Benefits,” for example, may flow *both up and down* in multiple dimensions. For a network in which nodes *give* something to others, *giving* may be any of a whole array of possibilities, from taxes

or tithes to punishments or gifts. In addition, if one relation forms a pyramidal *dag* hierarchy, it may be complemented by others, including *dag* hierarchies of the same or different types. The alignment of *dag* hierarchies and flows of various kinds then become research questions.

Ambiguities in the ways that hierarchies are constructed are a common source of instability in economic, social, and political systems. Even if we have data on multiple intersecting networks and processes that include sufficient indicators of hierarchy to understand the workings of these processes and their outcome ('population') distributions, the complexities of the interactions inflected by competing hierarchical structures, and the many direct and indirect traversal patterns that they allow, make prediction of trajectories difficult. These are complex nonlinear systems, and the best we can do probably is to model how their structural and flow properties provide an understanding of their dynamical instabilities.⁴

These instabilities at the core of complex systems may show up in market fluctuations between major periods of demand exceeding supply (scarcity, downturns) and those of supply exceeding demand (booms, bubbles), i.e., in business cycles or in the Kondratieff (1984) cycles of investment and degradation of infrastructure. They also occur in the Kuznets cycles of migration (in response to wages and benefits for labor) and, of course, in the cycles of leading sectors of innovation (Modelski & Thompson, 1996). As we will see, instabilities associated with inequalities in hierarchies and the ambiguous and often unpredictable quality of flows provide crucial points of entry into discussions of city and trading system dynamics.

One reason that markets require regulation by norms and laws is the need to maintain level playing fields in the face of such instabilities. Often, this is achieved by reverse hierarchies regulated by moral obligation and reciprocity.⁵

Figure 5.1 shows two examples of reverse social hierarchy, each with three levels of advice given by women migrants to others (directed either downward in the hierarchy or laterally symmetric, i.e., a semi-order) in a housing complex for urban workers (Du et al., 2007). Moral orders of this sort, in which leaders act on behalf of others, are common among the urban poor (Lomnitz, 1977) but are also very common in pre-state societies.

⁴ Elements of my discussion of ambiguity and complex dynamics in multirelational interactions involving hierarchical elements are often referred to as "heterarchy." Definitions of heterarchy are somewhat elusive and may refer to distributed authority, hierarchy that is flattened by design, exploitation of overlap of different evaluative principles (Boltanski & Thévenot 1999), self-organization (Morel & Ramanujam 1999), organization of diversity, organizational self-similarity or fractality (Abbott 2001: 165), and the like.

⁵ There are, of course, variants on the *dag* model, whether for hierarchy or reverse hierarchy. A *dag* structure may have transitivity or partial ordering as an additional property, and may be strict (irreflexive) or non-strict (reflexive). The latter terms may also describe a containment hierarchy or taxonomy. Variants of these sorts of hierarchies, such as semi-orders, may also have symmetric ties or directed cycles within levels, and partial orders between levels.

Network Example 1: A Hierarchical Network of Generalized Reciprocity

Not all human social hierarchies are based on power, domination, or authority. Survival on the margins of urban society might require, and surely mobilizes, networks of social support that are outside formal economic institutions (Lomnitz, 1977). Du, White, Li, Jin, and Feldman (n.d.) were surprised to find evidence of network hierarchy in their analysis of networks of assistance in everyday activities and of discussions concerning problems of social support (marriage and family, childbearing and planning, contraceptives, old age) among 200 women migrants in a Chinese industrial city. Some of the women's reports of their outgoing ties were mirrored symmetrically by others who reported outgoing ties back to them. When symmetric ties were excluded, however, the asymmetric ties in five of the seven relations reported for the network (all but emotional support and social activities) showed a nearly perfect hierarchical structure of the *dag* variety (above) in which the tendency is for help and advice to flow down from woman to woman through ordered levels, as in Fig. 5.1. Even when all five of these relations were combined into one, the resulting symmetries were only between pairs of women in adjacent or nearby levels of hierarchy. Women at the upper levels of these hierarchies of giving were those with higher status, higher income, and more family connections, including marriage. This is generalized exchange, and it exemplifies a social support system that is self-organizing. 'Generalized' reciprocity may take effect in the course of the life cycle but it is neither balanced in the short run nor necessarily even in the long term. The ranking in these hierarchies is one of responsiveness to a morality of obligation. Such hierarchies are not atypical of the moral orders of pre-state, pre-urban, and many village societies, in which leadership is more of a *primus-inter-pares* obligation to help others rather than one of domination. What is remarkable in this case is that when each woman's ranking in each of the five networks is normalized to a scalar measure between 0 and 1, the common variance among the five rankings in a principal components analysis is 78%, and the variance structure is single-factor. The morality of obligation transfers very strongly across these hierarchies of giving.

5.2.2 Assessing Hierarchy in Reverse, Zipfian, and Intermediate Cases

One way to assess inequality in social hierarchies, along with *dag* levels if they can be computed, is to measure the *outcomes* of resource access. For example, among the Cheyenne, leaders were expected to give certain gifts in response to requests. The distribution of horses among men was relatively flat in relation to the network or ranking of prestige differences, so leaders and officials were *primus inter pares* (Llewellyn & Hoebel, 1961). Inequality indices are useful in the assessment of how resources distribute on hierarchical networks. If, for a pyramidal *dag* hierarchy, the

ranks r in the hierarchy differ in their occupancy probability for resource measure x , so that a resource distribution is unequal but evenly divided by equal intervals in the log of rank, we have the Zipfian rank-size “law” $p(x) \sim r^{-1}$ (or cumulative occupancy probability $P(x) \sim r^{-2}$). If rank 1 has x resources, rank 2 will have $x/2$, the third $x/3$, and so forth. For k occupants with the basic resource, $k/2$ occupants are expected with twice the resources, $k/3$ with three times the resources, and so forth. The occupancy distributions of numbers of women at different levels in Fig. 5.1 follow a power law $p(x) \sim r^{-1.5}$ in part A, with more pronounced inequality than the Zipfian, but follow a linear progression in part B that is non-Zipfian.

Structural properties of a network often provide the basis for hierarchy as if, for example, there is an advantage to having more ties. Figure 5.2 shows a pyramidal *dag* for buyer-seller links in a Tokyo industrial district (Nakano &

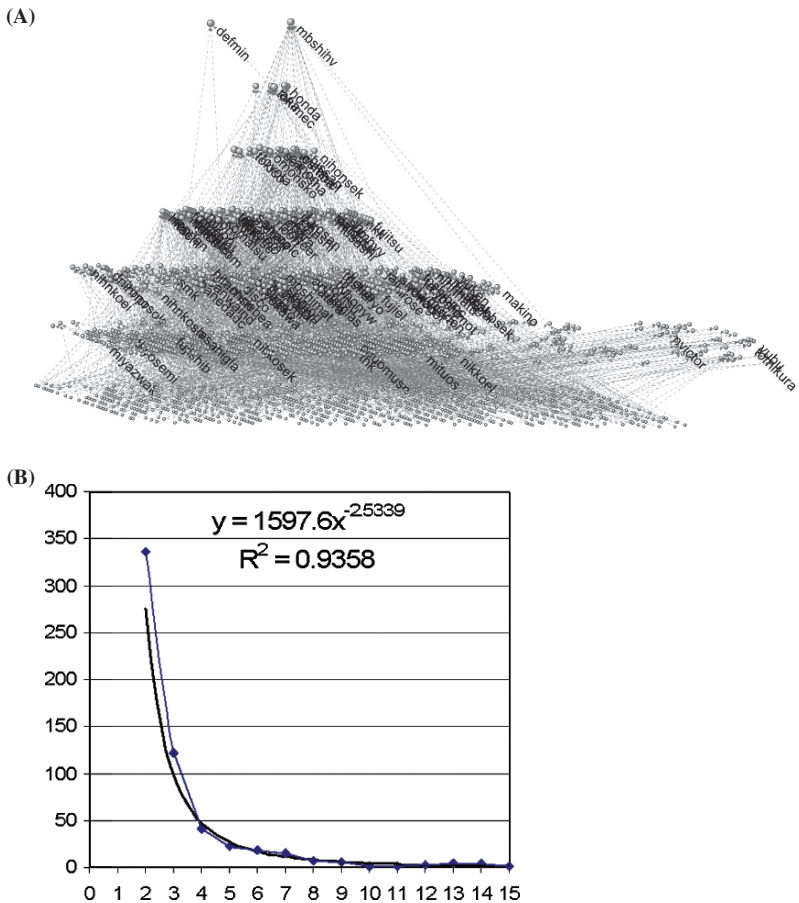


Fig. 5.2 Tokyo Industrial District network. (A) bi-component and (B) numbers of suppliers for the largest set of firms (1609) wherein every pair is connected by two or more independent paths. Over 300 firms (such as the Matsucom watch company) have two ties, 120 three ties, and so forth

White, 2006a, 2006b, 2007) in part A and a thinner-than-Zipfian power-law tail, $p(x) \sim x^{-2.5}$ ($P(x) \sim x^{-3.5}$) in part B. The upper firms have more ties than expected in the rank-size law, an amplified inequality that matches that of their capital and ability to organize suppliers in ways that can generate pricing advantages from them.

Figure 5.2B shows the type of distribution expected in a “preferential attachment” bias where a firm entering the hierarchy will connect to an existing node having in-degree k with probability $p(k) \approx 1/k$. Barabási and Albert (1999) call networks generated in this way “scale-free” because they have power-law degree distributions. In the industrial hierarchy, the contracting to suppliers is organized top-down, however, and network generation is not a “preferential attachment” of seller to buyer but an “organizing activity” bias of buyers recruiting suppliers. Still, if $p(k) \approx k^{-1}$ is an accurate generating probability for such a hierarchy (a finding that would require monitoring growth through time), the actual degree distribution will be a power-law, where the frequency of nodes with degree k will be $p(k) \approx k^{-\alpha}$ where $\alpha \rightarrow 3$ as the number of nodes grows large. The Tokyo distribution, for example, is within the range of the “scale-free” preferential or organizational attachment model.

A more general index of inequality is the Pareto II distribution, where x is the resource variable, y is the fraction with that resource level, and the probability of a cumulative distribution of having resources x or greater is $y = P_{\Theta, \sigma}(x) \approx (1+x/\sigma)^{-\Theta}$. As $\sigma \rightarrow 1$, the probability $P_{\Theta}(x) = (1+x/\sigma)^{-\Theta}$ converges to the Pareto power-law $P_{\Theta}(x) \approx x^{-\Theta}$. In this case $\Theta = 2$ is the Zipfian distribution in which every sum of resources is equal within equal log intervals. The P_{Θ} distribution becomes more exponential (converging to e^{-x}) as $\Theta \rightarrow \infty$, resembling a distribution of possessions that is random or independent of rank or resource level. The phenomenon of inequality generated by this random process is of interest because of the ambiguity of producing an apparent hierarchy in which “virtue” is attributed to winners of what might be a random lottery, as described for example by Farmer, Patelli, and Zovko (2005) for the rankings of hedge funds on the basis of their earnings.

5.2.3 Transitions from Reverse to Level Hierarchies

Another way to assess inequality in social hierarchies is to look at the structure of flows and transactions. Big-man systems (such as on New Guinea), chiefdoms, kingdoms, and early states provide useful examples for evaluating the initially ambiguous concepts of hierarchy and reverse hierarchy by the larger contexts in which many different kinds of entities, relations, and attribute distributions intersect. These examples move from kinship to nonreciprocal exchange structures.

Wright (2004) examines ethnographically well-described chiefdoms and archaeologically and historically known instances of the emergence of early states. The concept of reverse hierarchical ranking as developed here applies to his description of chiefdoms, marked by a single leader who does not delegate authority, in contrast to states with a delegated division of labor for authority. In his view, a stage-wise developmental process from chiefdom to early state is inappropriate. This is corroborated by the fact that pristine states emerged at Lake Titicaca (e.g., Griffin and

Stanish (2007, p. 24) after a long period of cycles in which multiple chiefdoms climb the population size gradient only to be fragmented by fission. In primary centers, the increase in productivity and population occur sporadically during that period, but without synchronization. Next, in one rapid burst, these unsynchronized features – including bipolarity of centers – emerge synchronously, passing a probability threshold for fission. As a result, a pair of networked city-states emerges, with extensive trade networks that include lesser centers.

A simulation, repeated hundreds of times, reconstructs the fission probability relative to settlement size as a humped or threshold distribution. It may be represented as the product of increasing resistance to political consolidation (which is roughly linear with size) and the ratio of resistance to chiefly strength (which decays exponentially with size).

It is important to note that this *synchrony*, as is often the case, is not an instance of homogeneous change (e.g., within the region), but a *network phenomenon* that occurs *between* the chiefdoms in the region, involving competitive innovation. As the new state polities consolidate, they assimilate the former chiefdoms by introducing the new, distinctive feature of hierarchical representatives of these old chiefdoms into the power structure.

Wright's (2006) *Atlas* clarifies the contrast between chiefdom and state in a way that is consistent with this dynamic of threshold transformation in level and scale. His sample of primary and secondary states all show *three or more* levels of mobilization of resources upwards through officials that are designated in a hierarchy of divided offices, but, furthermore, *state offices outlive and recruit their occupants*. Chiefs and paramount chiefs, in contrast, may govern subdivided territories with local leaders and ritual specialists but there are at most *two* levels of chiefly resource mobilization conducting directly to the chief and *all political decisions are integrated into the chiefly persona*. Cultural concepts govern how others connect to that person through relational ties and levels that embody differences of rank and reciprocal expectations. Succession is a matter of competition by personal exertion of power and negotiated rank rather than recruitment to office in a division of official political labor. While there is pronounced hierarchy, the notion of reverse hierarchy applies in several ways. The hierarchy is constructed interpersonally rather than through offices that concentrate the ability to command resources, and interpersonal orders of rank embody obligations to redistribute resources to counterbalance the processes whereby resources were concentrated through interpersonal network ties. The system of rank, resource mobilization, and redistribution operates as a segmentary system,⁶ focusing on a single apex, in the person of the chief, as a unique element holding together the hierarchical flows and counter-flows of goods. Intermarriage and exchange within the chiefdom, however, integrate the chiefdom through cohesive hierarchies and cross-cutting ties and structural cohesion. Chiefdoms have a dynamical tendency for their cohesive integration

⁶ In a segmentary lineage organization, minimal lineages (such as nuclear families) are encompassed as segments of minor lineages, minor lineages as segments of major lineages, and so on.

to fall apart as growth occurs and a tendency to fission at times of crisis, especially if these coincide with issues of political succession. There is no tendency in these dynamics for permanent offices to form within the political hierarchy or for any gradual evolution toward state organization. The mosaic of sub-chief territories mapped into the chiefly ranking are segments that recurrently fall apart and re-form in successive periods of political change (see Griffin and Stanish, 2007).

We might go so far as to think of paramount chiefs (two levels of offices), chiefdoms (one level of office), big-men systems (emergent leadership without office *per se*), and all combinations thereof and gradations within them as variants of personalized leadership systems that are distinct from the three-plus levels of less personalized offices of state systems. The emergent-leader systems that White and Johansen (2005) have studied in the Middle East constitute a more segmentary version, one with less recognition of leadership in the form of feasting and gift-giving (as in New Guinea) and a greater emphasis on giving advice for community decision-making. However, these are graded differences and not a solid basis for constructing a typology.

Wiessner & Tumu (1998, pp. 299–301) summarizes the results of a huge oral ethno-historical reconstruction project for the region of Highland Papua New Guinea, where the extensive Enga *Tee* ceremonial exchange systems developed following the introduction of the sweet potato. She details the network dynamics and constraints of systems in which exchanges, operating through kinship obligations and reciprocities and their extensions, not only expanded with new surpluses, but also supported innovations in the means of funding more distant exchanges by including in the network chains of clans linked through marriage along exchange routes.

The structure of these exchanges, following the trade routes along valleys in the first period of productivity and trade expansion (called the period of the Great Wars), was markedly segmentary. Big men competed for exchange advantages and accumulation of wealth to fund ceremonials, draw together support groups, wage wars at the cut-points in the segmentary structure, and, then, engage in regional ceremonials of reparations and peacemaking. But this did not lead to the emergence of hierarchical leadership. “The largest Great War . . . gave way to the *Tee* cycle in the late 1930s/early 1940s” (Wiessner & Tumu, 1998, p. 298), consisting of three phases of exchange ceremonies, moving from east to west and then back again, in which “thousands of pigs, goods, and valuables were distributed through a series of linked public festivals. . . . Prestige, influence, and the necessary resources for polygamous marriages accrued to those who were most successful in channeling wealth to themselves and their clans.” Together with the shift from linear chains in the trade-route links of the Great Wars period to locally radial chains operating through polygyny in the *Tee* ceremonies, the overall social integration shifted from segmentary to cross-cutting, but without a change in political leadership structure. Sons of Big Men still had to compete for prestige, as did everyone else, and permanent political office never emerged.

Hence, both Big Man and Chiefdom systems, consistent with Wright’s view, do not develop into pristine states with political offices. The development of states, arguably, is connected to the rise of divisions of specialized labor, production, and

exchange within and between cities in the evolution of a very different kind of hierarchy. Ties of exchange in early states are somewhere *between* connecting and separating ties; they are competitive, but also cooperative.

White and Johansen's (2005) account of the ethnogenesis and sociopolitical dynamics of a pastoral nomadic clan shows a social system that is scalable through structural cohesion, and that expands through reproduction, kinship alliances, and fissions, and overcomes internal conflicts and those with neighbors along routes of migration through pastoral ecozones that promote high levels of human health and demographic reproduction. The resulting networks constitute a generative demographic engine for large sibling groups and extensive cooperation within and between these groups, which is constructed through reciprocal ties of marriage. The organization, in this case, blends scalable segmentary organization (through the fission/fusion dynamics of male groups) with the scalable structural-cohesion dynamics of marriage alliances. The hierarchical aspects of male rank are flexible, variable in terms of age-succession and migration options, highly egalitarian and achievement-oriented, with historical leadership in the local group not decided on an elective or ascriptive basis, but attained through emergence processes involving the leader's position in the structural cohesion hierarchy. In other words, it is not based on lineage but on a combination of segmentary and cross-cutting lineage and alliance dynamics.

This lineage and alliance study looks under the hood, so to speak, of one case of what is a widespread cohesive dynamic of certain segmentary systems with a complex, multi-hierarchical multinet. Some would call this a form of heterarchy, and it stands between hierarchy and reverse hierarchy in ways capable of rapid adaptations that easily assimilate and generate innovation. Of special interest is how these segmentary nomads adapt to the spread of the market economy and the shrinking of free space, while a form of egalitarian leadership endures, outside the state system, as a form of reverse hierarchy and acts in the broader interests of the larger community.

Economics might be seen here as an attempt to tame the opposition in terms of opposite tendencies between connecting and separating, competitive but also cooperative ties that might settle into equilibrium because of this opposition, which allows equilibrium states to be defined, whether stable or unstable. But the resolution of this opposition through pricing, with its alternations between over-supply and excess demand for investments, at many different spatiotemporal scales, leads, almost inevitably (although not uniquely as many other resolution mechanisms, like vengeance killing, are unstable) to recurrent near-equilibrium fluctuations. These become part of the story of innovation with the startup of city networks and states.

5.2.4 Segmentation and Cross-Cutting Integration

Once we consider complex multi-nets (multiple overlapping networks and attributes) such as the above, it may be better to use higher order terms for the *effects* of

multi-net overlap rather than rely on the terms hierarchy and cohesion, which have well-defined formal properties but whose properties are poorly understood once they are overlaid in different ways. This is consistent with the argument I have made, so far, that “from their effects shall we know them,” i.e., from their dynamics, including both in structure and process.

The contrasting terms for segmentary and cross-cutting organization have long been used in the social sciences and are widely understood. Comparative studies by Thoden van Velsen and van Wetering (1960) found correlations for contrast between segmentary and cross-cutting forms of organization in pre-state societies: while segmented male descent groups, each resident in their own localized communities, have high levels of feuding and internal warfare, while groups that are cross-cut by different organizational and residential affiliations have internal peace.

In a multi-net-multi-attribute perspective, *segmentation* highlights separability, efficiency, specialization, restricted access, and uniqueness of traversal, among other features. For example, Ojibwa society was historically segmented during six to eight months of the year because men would set off individually on long trapping expeditions. With so little male cooperation in the means for economic success, there were almost no overarching social, political, or religious institutions for larger groups of people, other than temporary and fluid aggregations when men came back in the summer. Two people in a canoe could harvest wild rice in the fall. Women cooperated in their activities in small groups.

In contrast, in the great Pueblo villages, the underlying productive technologies were highly collaborative, and there was every kind of social, political, or religious institution that cross-cut and integrated the many cohesive subgroups within the society that supported cooperation (White, 1969). *Cross-cutting integration* involves multi-nets that intersect and overlap in cohesive patterns opposite to those of segmentation. Intersecting social circles and overlaps in group affiliations were theorized and studied by Georg Simmel (and later Peter Blau) as the basis for urban social integration.

Finally, urban societies can be characterized in terms of segmentary versus cross-cutting transport patterns. The usual downtown city grid is cross-cutting relative to urban spaces, exposes diverse groups to one another, and typically is regarded in the history of cities as the peaceful zone of integration of diversity in interactions. Gated communities and branching street or road structures (Low, 2003), however, create community segregation at the cut-points where one branch separates from another. Segregation is augmented when there are structurally cohesive routes of traffic within a cut-point-segregated community. Further, while structurally cohesive residential streets connect to form traversal patterns that support social cohesion, secondary, primary, and freeway routes disconnect the zones on either side of them, and structural cohesion on these larger routes *does not* create traversal patterns for inter-community cohesion (Grannis, 1998). Grannis (1998) and others have quantified some of the relevant network measures for effects of the built environment on daily life and sociological outcomes. The correlates of these patterns turned out to be the same as those found by van Velsen and van Wetering (1960): more segmentation leads to more interethnic violence, more crime at both ends of

the economic spectrum, and higher segregation on every major index of inequality (Grannis, 1998).

Two of the key developments still outstanding, however, are how to measure structural cohesion in a multi-net and how to identify more precise measures of segmentation versus cross-cutting integration in overlapping relational structures. The first extension of structural cohesion measurement to multiple relations is to study which relations are correlated across pairs of people. Clustered relations with similar kinds of content transmission can be identified as those that form interactive communities. Once those community-forming sets of relations are identified, they can be analyzed together by simple aggregation, and these aggregate networks can be analyzed for boundaries of structural cohesion. The communities may overlap, forming cross-cutting integration, or may segregate, forming segmentary structures. How societies and their transport and communication systems are constructed may have massive impacts on the types of conflicts and areas of cooperation that occur within contemporary social and information societies.

5.3 Relating Structure and Process

5.3.1 Hierarchies are Supported by Structural Cohesion

Biological organisms are supported by branching networks (blood vessels, bronchial and nervous system pathways) that often have a hierarchical *dag* single-root tree-structure, but they are also embedded inside cohesive structures of cells and organs and have return pathways that form various types of cycles, depending on whether we are considering pathways in the brain with lots of local cycles, circulatory vesicular pathways with separate branchings for vein and artery systems, or bronchial inhalation/exhalation two-way respiratory pathways. Social hierarchies are extra-somatic, but they also require such cohesive supports for multi-connected cooperation, independent of hierarchy.

Menger's multi-connectivity theorem⁷ provides precise ways to define how cohesive network structures relevant to social organization and hierarchy are constructed. The multi-connectivity theorem is especially important because it connects the cohesive *structures* that occur in networks with the capacities for multiple independent *routes of traversal* within these very same structures. Recall that structure and traversal are the two fundamental kinds of properties that combine in networks. The multi-connectivity theorem makes it possible to understand scalability – i.e., the ability to scale-up the size of an organizational structure with roughly constant or diminishing cost per increment – of a radically different sort than “scale-free” networks. This is possible through *structural cohesion*, which I will

⁷ In the mathematical discipline of graph theory and related areas, Menger's theorem is a basic result about connectivity in finite undirected graphs. It was proved for edge-connectivity and vertex-connectivity by Karl Menger in 1927.

explicate here. Scalability of cohesion seems like an impossible benefit of networks, but is achieved quite simply once the concepts that define structural cohesion are understood.

A *component* of a network is a largest possible connected structure; but as it may be minimally chain-link connected, it is potentially vulnerable to disconnection by removal of a single node. A *bi-component* is a largest possible connected structure that cannot be internally disconnected by single-node removal. Components, bi-components, tri-components, and k -components (vulnerable to disconnection only by removal of k or more nodes) are instances of the general concept of *structural cohesion* but differ in level of cohesion. Menger's theorem is that k -components (k measuring the extent of *structural cohesion*) are equivalent to largest sub-networks with at least k redundant pathways between every pair of nodes (i.e., all pairs have a minimum level k of cohesive *traversal* between them). The structural and traversal aspects of cohesion unite to form a single concept with two aspects: (a) structural cohesion as invulnerability to node removal of outside attack on members that will disconnect the group (structurally) and (b) ease of (traversal) communication and transport within.⁸ If we identify k -components in social networks as sub-graphs of nodes and the edges between them, i.e., structural cohesion *sets*, we can also call them by the term structural cohesion *groups*. Their cohesive properties of multi-connectivity involve hierarchically stacked subgroups. When the cycles that generate the cohesion are sufficiently overlapping and short the redundant paths of information passing between them (and the diffusion of shared attributes) can be assumed to entail knowledge about core members' identities.⁹

Structural cohesion groups can grow very large while minimizing the possibility that parts of the network will disconnect with the exit of a limited number of members. In this kind of structure, a sub-network can add nodes and links to strengthen traversability within and to reinforce immunity to disconnection from without. Scalability is achieved when most of the nodes in an expanding group stay close to a minimum threshold of k ties per node without diminishing the number of independent traversal pathways that this enables between every pair.

Unintended cohesive optimization occurs under these modest constraints even when the ties of new members are random. Scale-up may even diminish the increments of cost involved in connecting other nodes to every organizational node *inside the organization*, especially when the average network distance within a structurally cohesive group does not automatically entail longer spatial distances. Alternatively, the attenuation of spatial distance effects can be achieved by innovations in long-distance transport or communication.

⁸ Every science has a concept of entities as cohesive units having sufficient internal bonding among all of their parts to provide resistance to external perturbation. It is important to note that *edges, per se*, are conceived of as entities as are nodes and possibly emergent structurally cohesive units.

⁹ Such groups do not have to be named or formally organized, and mutual identification need not be assumed beyond core members. It is sufficient to assume that members of the less cohesive sectors of the cohesive hierarchy (not in a higher component but in a lower component that contains it) will have weaker identifications for members of the group cores.

5.3.2 *Structural Cohesion as a Scalable Concept*

Structural cohesion groups have obvious benefits for cooperation and collaboration, and, if they can scale up with little or no institutional cost or additional cost per participant, we should find that they are implicated in a wide variety of contexts, including social stability and the formation of social units on the one hand, and in innovation, social movements, competitive rivalry, and conflict on the other. How can k -components scale-up to have *more nodes* – conserving overall strengths of structure and traversal – *without* increasing the number of edges per node? How are organizational costs reduced when adding nodes and edges? The simple answer is that a constant *cost per node* can be distributed over benefits to *pairs*, i.e., by dividing the costs of k edges per node by the *square* of the number of nodes in the network. Because every node in a k -component has k *independent paths to each of the others*, a growing network that conserves this property, merely by the placement (random or not) of edges, can reduce costs accordingly. Successive growth can occur in which the edge/node ratio is constant while the n^2 node-node pairs with k independent pathways grow exponentially. Structurally cohesive groups have a network externality (Arthur, 1994), defined as marginal returns for items interconnected in the network, that increases with the number of nodes in the network. An example is a set of word processors that can exchange documents: the greater their fraction in a network, the more valuable it is to possess one of this set rather than one of an incompatible set.

Many networks contain cohesive sets with relatively low k -cohesive index numbers. There are several reasons why it is so poorly understood that scalability of a structurally cohesive organization permits it to grow indefinitely. First, intuitions about structural cohesion are often incorrect if they are based on the analogy of growth of clique sizes within a network that entails analogy of increases in network density and in the number of average edges per node. The growth of groups of constant k -connectivity can occur while there is a reduction in overall network density, and edge/node ratios can remain constant within k -components. Second, the Menger theorem is one of the more intricate parts of graph and network theory, and the necessary ease of computation to deal with it has developed only recently (White & Harary, 2001; White & Newman, 2001; Moody & White, 2003). Third, because this is a new area of investigation, there are still relatively few studies that show how common are the properties of cohesive hierarchies that would allow them to be considered structurally cohesive groups rather than arbitrary sets of nodes and edges.

Structurally cohesive groups with informal organization can occur entirely outside of formal institutions. They can build on friendship networks, networks of trust, social movements, co-alignment based on opposition to other groups, structurally equivalent linkage to other groups, and the like. Hence, up-scaling of structural cohesion can occur without increasing institutional overhead. To see the potential power of structurally cohesive groups, we can look to Turchin's (2005a) documentation of how historical empires expand repeatedly, only to encounter boundaries where identities of groups exposed to conquest change sufficiently to provoke

solidarity, so that further expansion is blocked by escalating opposition. Turchin found that interethnic frontiers tended to be the sites where emergent solidarity allowed groups that were relatively disorganized at the start to oppose and stop assimilation by the empire. The U.S. is currently experiencing such a situation in Iraq, and we may see the same resurgence against occupation by the U.N. in Afghanistan. Often, such resistance movements come to overthrow empires that were initially much more powerful. Alternatively, by repeatedly attacking settled populations at their boundaries, groups with low population density like the historical Mongols develop the cohesion to overrun powerful state systems.

A *small-world* is a large connected network in which the edges around certain sets of nodes tend to form denser clusters (something that is possible but not required in structurally cohesive groups) and average distance between nodes is at or below that of a network formed by random rewiring of edges (Watts & Strogatz, 1998). A large clustered network requires only modest random rewiring to retain the small-world property. Similarly, only a modest amount of random rewiring is needed to create a cohesive k -component among almost all the nodes in sub-networks with k or more edges per node. This explains the mechanism that underlies Turchin's findings – that when ethnically distinctive groups are attacked by intrusive empires, initial disorganization may quickly transform into a scale-up in size and resistance of structurally cohesive groups, whose size is virtually unlimited as more contestants are brought in against a common outside enemy. This can be achieved almost effortlessly even when new bonding occurs randomly, as in the chaos of war, and may produce high levels and wide dispersal of structural cohesion. It gives rise to solidarity and *asabiya*, or emergent cooperation, as defined by historian and Muslim scholar Ibn Khaldun (1332–1406).

White (1996) was the first to theorize, and Brudner and White (1997) to demonstrate that cohesive social class boundaries could be formed by structural cohesion.¹⁰ Following them, Fitzgerald (2004) showed cohesive breaks corresponding to (and partially defining) social class boundaries in one of the historical religious communities in London. These studies challenge views about pluralist political systems within contemporary states. It is shown that structurally cohesive marriages often link the political class into a coherent social unit, or power elite, through structural *endogamy* (White, 1996). Moody & White (2003) found as a general result for American high school students, that attachment to school may be predicted, other things being equal, as a function of the degree of structurally cohesive membership in school friendship circles. This result was replicated precisely in each of a sample of twelve randomly chosen high schools, among those surveyed in the AdHealth

¹⁰ In these studies, structural *endogamy* was used as a special case of structural cohesion that describes the cohesion added to forests of genealogical trees by the facts of common ancestries in earlier generations and intermarriage or offspring of two genealogical lines below. A specialized literature has developed in ethnology around this concept that will not be reviewed here. The line of research leading from structural endogamy to the more general concept of structural cohesion was supported by NSF Award BCS-9978282, "Longitudinal Network Studies and Predictive Social Cohesion Theory," to D.R. White, 1999–2003.

study of 100 representatively selected American high schools, and was validated through a multiple regression using all grades and all available network and attribute measures as competing and control variables.

5.3.3 *Structural Cohesion and Innovation*

As for innovation, the Roman historian Velleius (ca. 19 BC–AD 31) noted “I frequently search for the reasons why men of similar talents occur exclusively in certain epochs and not only flock to one pursuit but also attain like success.” (Velleius, 1924, p. 5). Could this be because high levels of structural social cohesion and interaction among people with similar interests in innovation occurs rarely but, when it does occur, the diffusion of ideas within the group and the success of its members are significantly enhanced compared to isolated individuals or groups with lesser cohesion? Kroeber’s (1944) *magnum opus* on creativity and cultural growth found Velleius’s observation to hold for the vast majority of innovations in philosophy, science, philology, sculpture, painting, drama, literature, music, and national growth over the course of history covered in his review. The one regularity he found was that innovations and creativity – individual “genius” – occur in temporal spurts in certain places involving social interactions in specific groups. In addition, diffusion occurs in and around those groups in spatial and temporal clusters. Florida (2005) has recently introduced this phenomenon in the economic literature by explaining inventiveness in certain cities as due to the aggregation of creative individuals into a ‘creative class’ linked by ‘spillovers’ between different creative activities.

Allowing for network interaction that is not necessarily spatially proximal, Kroeber’s finding is no less true today. It is not unusual to see cohesive hierarchies of collaborations in new industries that grow out of innovations, which are enhanced by the multiplier effects of multi-connectivity, another name for structural cohesion. The biotechnology network below, studied through its collaborative inter-industry contracts during the early florescence of the industry (1988–1999), illustrates structural cohesion in early stages of a developing industry operating in recruitment for technological innovation and knowledge discovery.

Network Example 2: Organizational Cohesion in the Biotech Industry

Here, networks of innovation emerged out of the potentials provided by changes in the laws governing patents and commercialization that had developed out of innovations in Silicon Valley and the IT industry. These provided the scaffolding for organizations and institutions to build networks of scientific collaborations that expand the knowledge base and technology for development of new products. Rather than competing with large pharmaceutical companies for improved mass-market products for illness, healing or health-enhancing treatments, specialized niches were sought through collabo-

rative ties. Rather than protecting new privately held knowledge following the pattern of the pharmaceuticals, the industry built through successes as parent firms spawned specialized offspring. Emergent organizational couplings between complementary firms resulted in rapid emergence of a single broad structurally cohesive core of organizations linked through collaborative ties in a division of labor, including capital and marketing. In so doing, core organizations in the field benefited from broader diffusion of knowledge but, by bringing new recruits into a mixing process for new knowledge potentials and skills, also managed to prevent stultification that could result from sharing and homogenization. This involved a dynamic of periods of heavy recruitment of newcomers alternating with periods that focused on integration of knowledge as between new recruits and established practitioners (Powell et al., 2005; White, 2004). Balancing broad cohesion with innovation, the success of the industry led to scaling-up of overall size and dispersion of the industry sufficient to enable internal niche diversity to withstand competition by the pharmaceuticals, instabilities of the marketplace, and volatility of R & D funding sources. In the biotech industry study, variables measuring structural cohesion were pitted against several hundred variables other entered into a time-series database used to test hypotheses about predictors of new contracts between biotech firms and partners. McFadden's (1973, 1981) discrete choice model for each of 11 time-lagged periods showed two major sets of predictors that beat all other competitors in these predictions: levels of structurally cohesive multi-connectivity and measures of diversity. Powell et al. (2005) statistically demonstrated that new link formation and repeat links in the biotech industry showed proportionality effects for multi-connectivity: the more multi-connectivity, the more the attraction for potential collaborators, leaving only very slight effects for Barabási and Albert's (1999) model of preferential attachment to hubs.

5.3.4 Cohesive Network Hierarchy

Another theorem of structural cohesion is that every k -component is contained in the same (or larger) set of nodes that are structurally cohesive at the next lower level: a tri-component within a larger bi-component, for example, with the bi-component within a component, and the component within the total network, which, if not connected, is a zeroth-order component of itself. Structurally cohesive groups form hierarchies in this sense of ordered containment or *dag* structures of inclusion relations, not excluding the possibility that two k -components of higher order are contained in one of lower order.

The short take on the network of contracts in the biotech industry is that they form a cohesive hierarchy with a single cohesive group at the apex, which, in the

period of study, oscillated between a 6-component and an 8-component. The regression results are summarized by two findings: a firm's formation of new contracts favors those higher in the cohesive hierarchy, but those firms higher in the cohesive hierarchy *also* favor contracts with new entrants that bring innovation to the industry. These two tendencies counterbalance greater diversity within a firm's network with greater integration, while avoiding greater homogeneity. As for traversal of information, connection into the higher k -components also provides firms with a broader window on diversity within the industry (for advantages of diversity see Page, 2007).

5.3.5 Urban Cohesive Hierarchy

In a division of labor, it should be obvious that for large-scale organizations to get things done, they are dependent on those at lower scales, and that this is one basis for such scaling laws that exist. Another is that as organizations scale up in size, they require exchanges at greater distances, including those with the larger organizations at a distance. They do so as mediators of exchange or flows in the multi-relational network that sustains their existence, as producers of goods they consume, and as consumers of goods that they produce. City systems necessarily entail diverse levels and boundaries of structural cohesion and, thus, constitute cohesive network hierarchies. Unlike the single focus of cohesive hierarchies in the world biotech industry, linked by high-speed travel and communication, the geographic sites of urban cohesive hierarchies worldwide will be multiple and overlapping, with multiple apical cohesive groups.

The inequality indices discussed earlier (cf. Network Example 1 and Section 5.2.1) are relevant as outcome distributions for processes in urban cohesive hierarchies. Cumulative city populations at different city size levels tend to be Zipfian (Pareto slope 2) for larger cities but Pareto II overall. A Zipfian distribution implies that as you successively multiply a given level x of resources by a constant (additive log interval), the quantity of resources added at that new level is constant. The network itself, and its boundaries and levels of cohesion, provides a basis for hierarchy.

5.3.6 Network Dynamics Affecting Hierarchy: The Scale-Up of City Networks

City systems and their growth provide a useful example by which to examine network dynamics in a familiar context. World population growth rates from the Paleolithic through the Neolithic agricultural revolution are difficult to estimate and compare to growth rates after the rise of cities. Although mitochondrial DNA (mtDNA) "nuclear sequence data do not support a simple model of recent population

growth, we nonetheless know that a drastic population expansion occurred at least 12 kya with the advent and spread of agriculture. Furthermore, archaeological evidence suggests that human population sizes have expanded over the last 40–50 kya or more” (Walla & Przeworski, 2000). If the trend of that early expansion were characterized as a small positive constant rate of world population growth (such as 1.0001 or .01% per annum), the growth curve would be exponential.¹¹ Once interdependent cities arose through networks of trade, however, world population growth displayed distinctive bursts, in which the rate of growth varied with size of the population. These are bursts of power-law growth,¹² and they are followed by periods of leveling or no-growth (Fig. 5.2). Although the existence of such bursts has not been ruled out prior to the advent of cities, the ability of cities to attract an influx of population and to sustain high population densities provides an explanation for these bursts whenever rural areas have the ability to replenish their population. The advantages and attractions of cities are multiplicative rather than simply additive with size (Algaze, 2005). The earliest urban settlements also show Pareto size distributions, hierarchies of size, that tend to follow power laws even in the earliest periods of urbanization (Adams, 1981).¹³

With the rise of cities, total world population shows concave functions mirroring growth rates with increase proportional to size, unlike constant rate or simple exponential growth. This is expressed as $N_t = C(t_0 - t)^{-\alpha}$, where t is historical time. The t_0 is a time in the future, a “singularity” that acts as an absolute historical limit to the power-law growth trend. Such growth is unsustainable and requires readjustment as population explosion precipitates a crisis of energetic and material resources.

The first to identify power law population growth¹⁴ were Foerster, Mora, and Amiot (1960), who gave a “structural” explanation in terms of increase in network connectivity in the last millennium. When the interconnected component of a population grows to a considerable fraction, they argued, it is possible to conceive of a growth law such as $dN/dt = N^2/K$ where it is some proportion of connected *pairs* in the population of N elements that have an effect on growth. Only when connectivity or traversability through large portions of a total population is possible can there be diffusion effects that result from that connectivity. Adaptations learned

¹¹ The population curve through time for a fixed rate of growth is an exponential, i.e., population $N = N_0 e^{kt}$ at time t starting from an initial population N_0 , where k is close to zero. Percentage growth on a savings account is an example of fixed rate growth.

¹² Power-law growth, $N_t = C(t_0 - t)^{-\alpha}$, given a constant C and coefficient α , cannot continue past t_0 because N would become infinite, so that power-law growth benefits are self-limiting.

¹³ Here, the frequency $f(s)$ of cities of size s varies inversely to size, $f(s) \sim s^{-\alpha}$. In power-law growth, the savings account analogy would entail a rise in interest rise proportional to size of the account (the rich get richer faster).

¹⁴ It is easy to confuse this exponential growth – occurring at a steady growth rate – which will outstrip linear growth, say, for resources, as argued by Malthus. Power-law growth – with a growth rate proportional to some power of the current population – is super-exponential even if the power is 1.

in one part of a population can then diffuse to others, provided that communication is possible. This cannot occur if the smaller groups do not connect to allow diffusion from local regions to larger portions of the total population. (Exponential growth, $dN/dt = aN$, in contrast, can operate without this interactive component.) I would add that the ramp-up effect of multi-connectivity on population growth does not scale with n by n interactions in a population of n elements, but at the much lower cost and density of scalable cohesion, involving the spread of structural cohesion rather than more links per node. The effect is transmitted through trade, through the diffusion of technological innovation, and through attraction to cities that allows higher “storage” of dense populations while the higher reproductive rates of rural areas replace the out-migrants to cities.

Figure 5.3 shows, from Kremer’s (1993) world population figures, population growth data as presented in von Foerster et al.’s (1960) log-log format and a semi log plot of the same data. Because the log-log plot in part A includes logged “time to singularity” in an infinite succession of divisions by 10, it entails a population growth to infinity as the time scale expands to the left beyond what is shown here. It is also somewhat deceptive in suggesting greater linearity in power-law slope after 1500 because the data-points are increasingly compressed later in time. The semilog plot in part B shows the same data in equal-interval historical time from the Neolithic to the present. Here it can be seen that roughly exponential growth gives way to what might be power law growth – i.e., increasing growth *rates* – after the rise of cities in the fourth millennium BCE. Each of the intermittent precipitous growth trends after the rise of cities hit a limit and then flattened to start a new but similar growth trend. Although there is certainly room for improvement on Kremer’s data as a basis of these generalizations, the empirical and theoretical basis for power-law urban growth has now been given by Bettencourt et al. (2007), and I think that its prevalence will be verified further with improvements to the population data available.

Accelerating population growth *rates* associated with urbanism occur first at circa 500 BCE during the crisis of archaic civilizations, several thousand years before the singularity that would have occurred if the world population growth trend of that time had continued without change. Food production is a predominant factor limiting the acceleration of population growth. The findings of West et al. (this book) and Bettencourt et al. (2007) show that cities have to run to stay in place and to import managers and skilled workers in production to maintain their position *vis-à-vis* other competing cities in a regional and/or world-system economy. Extrapolating these findings back in time might explain the peculiar population growth phenomena observed even in the earliest phases of urban civilizations (Algaze, 2005). Page (2007, pp. 329–335, 340) argues for the super-additivity of diverse tools for problem solving that cohabit in cities as the driver of city growth. Structural cohesion provides a scalable model of low-cost network mixture in cities that facilitates such super-additivity and, hence, innovation through problem solving. While Tainter (2000) explains how the demise of large empires might be due to the increasing cost of problem solving, this lack of scalability may be due to the

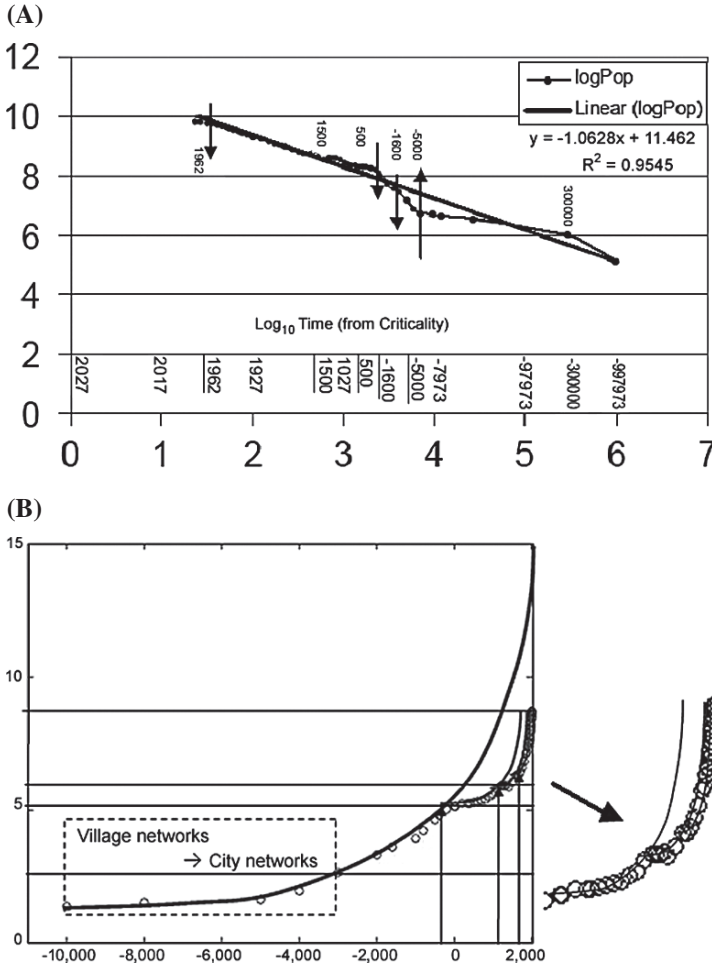


Fig. 5.3 World Population Power-Law Growth Spurts and Flattening as shown in: (A) von Foerster et al.'s (1960, redrawn) log-log plot with updated population estimates back to 1 million years BCE, with downward arrows marking the start of reduced growth rates, and (B) a semilog plot of the same data with successive power-law fits

structure of bureaucracies, replacing decentralized structural cohesion with competitive segmentary groups as the vehicle for network organization and collaborative problem solving.

After 500 BCE, there are three phases shown in Fig. 5.3 where world population growth flattens but then begins a precipitous rise toward singularity. Flattening occurs at about 1250 CE, then again at about 1650; and again at about 1860 (data from Kremer 1993). Each time, after the flattening, a new growth resumes similar to power-law growth, each curve pointing roughly, to the same singularity at t_0 . These growth spurts and halts are due partly to the growth of cities *per se* but may also be due to the new technologies and communications thresholds developed out of

leading urban economies and world powers that provide significant innovations and productive breakthroughs (van Duijn, 1983, pp. 176–179).¹⁵

5.3.7 Co-Evolution of Cities and Urban Networks

The concept of power-law tails of urban population distributions of leading cities has long predominated in urban studies and gives a false impression (noted by Batty, 2006) of cities as self-organizing systems (Bak, 1996) with universal rank-size scaling. City sizes are not scale-free and city-size distributions do not conform to invariant or even fully scale-free power laws. There is no commonly accepted explanation for this discrepancy in the urban studies literature, from Zipf to Sassen, nor for the Zipf “law” itself (Krugman, 1996). For a test of whether city size distributions are stable over time or have characteristic oscillations, data can be found in Chandler’s (1987) inventory of sizes for the world’s largest cities. Figure 5.4 tells the story by showing, for European, Chinese, and Mid-Asian city systems, two parameters for

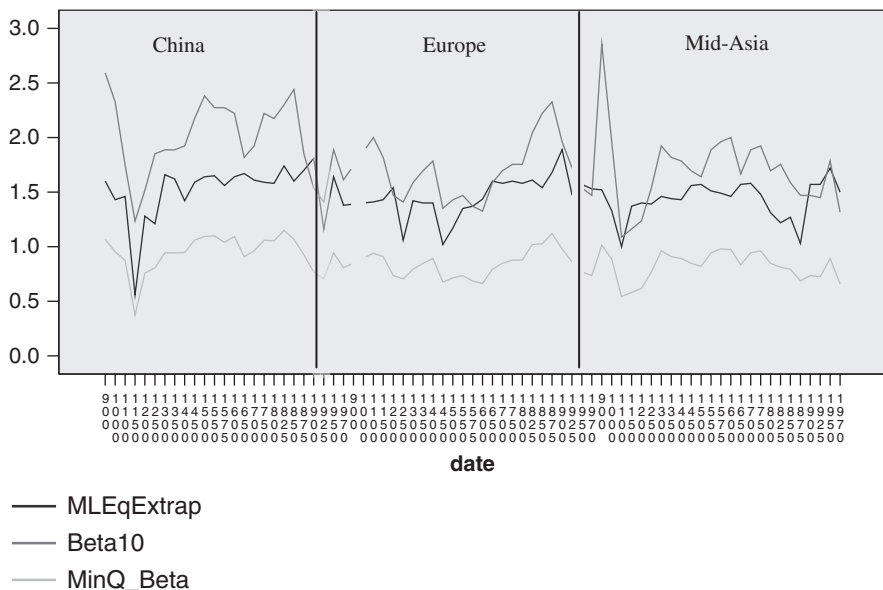


Fig. 5.4 Fitted parameters for city size distributions in Eurasian Regions (β in Pareto I tail, q in Pareto II body, and their normalized minimum)

¹⁵ It would be useful to explore the extent to which these include innovations in agricultural production that can support more population and organizational forms that can help cities overcome centrifugal tendencies. There may be exceptional agricultural innovations in rural England in the 18th century and the US in the 19th century but these are taking place within urban industrializing economies with market supports of agricultural innovation.

different curve fits to the city size distributions for each of 25 historical periods (White and Tambayong & Kejžar, 2007). The upper lines show the linear log-log slope β coefficients for the linear Pareto I log-log plots for the top 10 cities in each region and for each time period. The middle lines show variations in q , a shape coefficient for a Pareto II distribution¹⁶ that has a power-law tail but a crossover toward an exponential distribution for smaller city sizes. The equation fitted for the Pareto distribution is $\Pr(X > x) \approx 1/x^\beta$. The Pareto II equation that is fitted is $\Pr(X > x) \approx 1/(1-(q-1)x/\kappa)^{1/(q-1)}$ (Shalizi, 2007), where κ affects the crossover, and the slope of the distribution asymptotes in the tail to $1/(q-1)$. This asymptotic slope would equal $-\beta$ if the tail and the body of the distribution were consistent, which they are not. This tells us that larger cities and smaller cities are affected historically in differing ways, possibly, because large cities suffer declines with a foreign invasion, while their size distribution elongates with imperial expansions and broader world trade. Size distributions for smaller cities may be affected by regional trade volumes.¹⁷

The q coefficient varies around 1.5, for which the Pareto II curve has a Zipfian tail. The Pareto coefficient varies around 2, also conforming to the Zipfian distribution. The lower line is the Zipf-normalized average of the two, dividing β by 2 and q by 1.5. Each of these measures shows statistically significant runs of values that remain significantly above or below their mean for long historical periods. Their oscillations are nonrandom, and they define historical periods of ups and down in the tails and bodies of city distributions. Variations in β reflect historical fluctuations of larger cities, and q those of smaller cities.

Thus, as shown in Fig. 5.4, between 900 and 2000 CE there are alternating periods of high and low β and q values for each of our three regions. Between different major historical periods, they change between high- q periods (with thicker power-law tails and greater heterogeneity) and low- q periods (more exponential in the shape of the distribution and egalitarian in terms of more thin-tailed city size differences). Low q , low β periods represent collapses in city systems; low β being collapse of large cities, low q the collapse of smaller cities.

To gauge the effects of major wars on the irregular oscillations shown in Fig. 5.5, we defined a dichotomous variable W that measured sufficient SPI magnitude, duration, and extent of conflicts that may be considered as external shocks affecting



Fig. 5.5 Socio-Political Instability (SPI) recoded as dichotomy (W) for warfare shock to cities

¹⁶ Defined in the first section as $y = P_{\Theta, \sigma}(x) = (1 + x/\sigma)^{-\Theta}$ the only difference being a linear transformation of these parameters (σ , Θ) into κ and q (see Shalizi, 2007).

¹⁷ We have also found that our MLE estimates can give unbiased estimates with reasonable confidence limits for total population in China in different historical periods, and for Europe and other regions as well.

urban population movements, based on a 0/1 coding of Lee's data for China, and comparable data from multiple sources on European wars, as diagrammed for the historical intervals in Fig. 5.5.

Using W as a measure of exogenous shock, we then modeled the dynamics of alternation between periods of city system collapse and city-rise or growth as shown in Fig. 5.4, for both Europe and China, by two time-lagged equations that summarize the puzzling relation between β and q and the effects of major disruptive wars (the C parameters are regression constants).¹⁸

$$\beta_{t+1} \approx -q_t + q_t \beta_t + C_\beta \text{ (overall } R^2 \sim .79, \text{ China } \sim 0.75, \text{ Europe } \sim 0.69) \quad (5.1)$$

$$q_{t+1} \approx -\beta_t + q_t \beta_t + C_q - W_t \text{ (overall } R^2 \sim .57, \text{ China } \sim 0.54, \text{ Europe } \sim 0.66) \quad (5.2)$$

These equations express a pair of time-lagged multiple regressions that predict, with a lag of a single generation (25 years), that

1. β and q affect one another inversely, i.e., the health or decline of one end of the size distribution ramifies to the other,
2. Their product, with a magnitude of disruption in consistent growth, affects each, and
3. Major wars affect q adversely, i.e., they are most directly disruptive to the smaller cities.

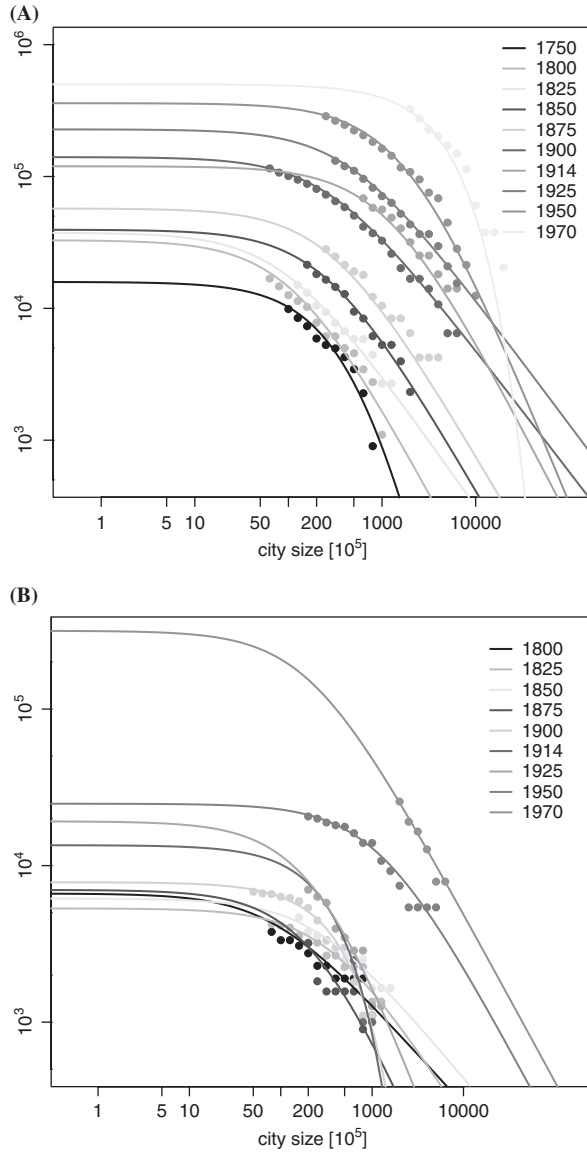
These effects occur with a time lag. Disregarding the effect of war, these interactions are positive feedbacks, and without war, they imply a dynamic that would lead to a city size distribution equilibrium close to the Zipfian.¹⁹ Major wars disequilibrate city size distributions, according to these findings (White et al., 2007).

Figure 5.6 shows some of the distributions of cumulative population by city size for China (part A) and world cities (part B). The curves in the log-log plots for the bodies of these distributions bend toward the horizontal as they approach the y -axis, consistent with the q -exponential. These fitted curves show dramatically how the city population curves differ from power-laws (which are straight lines in log-log; but this departure is more evident when the number of cities is cumulated), and how much the distributions differ from one period to another. On the y -axis, in logged units of 1,000, are the cumulative city populations of the logged city-size bins in the x -axis. Large differences are evident in the shapes of the curves between 1800, 1875, 1914 and 1950 for China, for example, and in the case of the cities of the world, one is struck by how 1700 differs by degree from 1850 and 1900. In their

¹⁸ W for China and SPI for Western Europe were coded by Tambayong. Equations (5.1) and (5.2) were also his work.

¹⁹ A two-equation reciprocal time-lag model such as Equations (5.1) and (5.2) produces fluctuations if the signs of the right hand elements are opposite, but convergence or divergence if they are the same. This can be verified in difference equations using initial values that generate a full time series.

Fig. 5.6 Fitted q -exponential curves for World and Chinese city distributions. **(A)** World Cities 1750–1970 and **(B)** Chinese Cities 1800–1970. Largest sizes are limited by a diagonal for cities at the lower right of these graphs. Projection lines beyond that limit are meaningless.



cumulative distributions all these curves asymptote at the tail, to a power law with slope β measured by $1/(q - 1)$ if $q > 1$, and approach the total urban population as they asymptote toward the y -axis when city sizes on the x -axis approach 4,000 people.

Many of the actual data points in the tails, however, differ from the slope expected from fitted q (here fittings were weighted by the cumulative number of cities in each bin, and, since these weights increase from right to left, they give weighting

priority to the body of the distribution). A key feature of these distributions is that of systematic exaggeration of the sizes of the very largest city as compared to what is expected from the rest of the distribution. This occurs up to 1950 for China, and to 1850 for the largest world cities. These exceptional sizes of historical cities tend to disappear with globalization. Thus, the exceptional primacies of capital cities in the historical periods may reflect their integration into cross-regional trading networks – trading connections to the hubs of other regions – an integration that is not shared by smaller cities until the economy becomes fully globalized.

5.4 Causality and Dynamics at the Macro Level

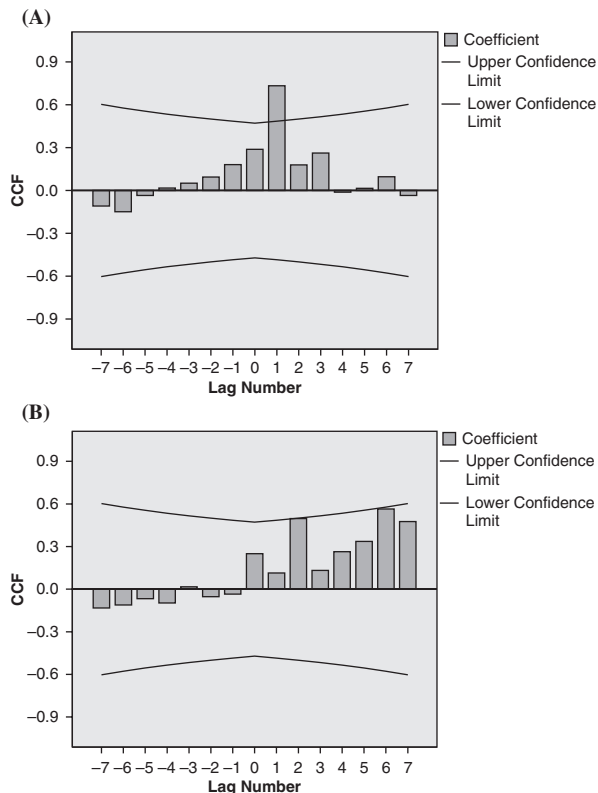
Networks and node attributes have all the elements for causal analysis, as implemented by Pearl (2000, 2007), H. L. White (2007), and White and Chalak (2007), and in simpler forms for historical dynamics by Turchin (2005b). Our use of MLE estimation (unbiased even for small samples) of q and β supports the validity of one of the most striking pieces of evidence of causal effects in our cities dataset: a time-lagged regression model, which replicates well in all three of our regions, shows that q values reflecting rise and fall in the smaller cities sector of the distribution at time t affects commensurate changes in the power-law tail at time $t + 1$ ($R^2 = 0.50$ for all three regions). This is relevant for innovation if we interpret the health of the small cities' economy, trade network, and city size distribution as indicative of innovative and competitive industries and exports, with the larger city distribution benefiting with a generational time lag.

The results illustrated above are only the beginning of our city system oscillations studies. Further early results show strong long-distance correlations between the parameters of city size distributions whose time lags indicate diffusion effects: changes in q in one region that have a temporally lagged correlation with changes in q in another. Temporal cross-correlations are shown in Fig. 5.7A for changes in the Islamic region (Mid Asia) that precede those in China 50 years later, and in Fig. 5.7B for changes in China that precede those in Europe 100 years later. The Mid-Asian dynamics only weakly link to changes in Europe 150 years later, which suggests that these effects are mediated by the Chinese cross-continent trade through the silk routes (Fig. 5.8).

It is clear in these time lags that when regions that are more technologically advanced as a diffusion source, in a given historical period (early on, Mid-Asia over China over Europe, later on, China over Europe), show a rise in q , a less advanced region rises in q with a time lag, but only when there are major trade flows between them. The same is true for the lag when q declines in the source.²⁰ This result supports the hypothesis that the carriers of positive long-distance correlations

²⁰ These graphs reflect improvements in accuracy in our scaling estimates using maximal likelihood estimation (MLE) methods for Pareto I and II developed in collaboration with Shalizi (2007). Pareto II scaling gives equivalent results as q exponential scaling.

Fig. 5.7 Cross-correlations for temporal effects of one region on another. (A) changes in the Islamic region (Mid Asia) and (B) changes in China that precede those in Europe 50 years later in China 100 years later



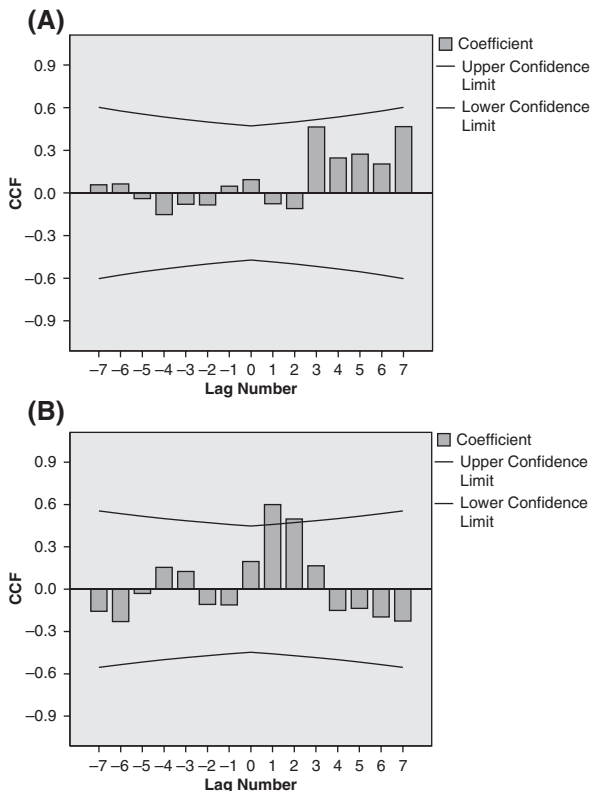
in city sizes are the intercity networks, operating primarily through trade relationships (cf. Turchin & Hall, 2003) and, with secondary effects in the smaller cities, as primary producers and mediators of trade networks. For example, these correlated oscillations, when lagged by 100 years, begin in Eurasia after the Song invention of national markets and credit mechanisms in the 900s and first diffuse to Europe through the silk routes.

Mobile states or groups, such as the Mongols, play significant roles both in establishing and furthering long-distance trade and in its disruption (Chase-Dunn et al., 2006). World-system wars that lead to disruptions of trade networks have a negative effect on city population distributions. Trade, warfare, and diffusion of information and innovations (or materials that can be used for innovations) play significant roles in city-rises and city-quakes at the population distribution level.

5.4.1 Innovation, Socio-Political Upheavals, Secular Cyclicality and Trade Networks in Complex Urban Societies

The processes examined here are among those that require a coarser, aggregated level of analysis to measure statistical and causal relationships. For a statistical

Fig. 5.8 Time-lagged cross-correlation effect. (A) Effect of Mid-Asia q on Europe and (B) effect of the silk road trade on European β (ten top cities)



model, a *sufficient statistic* is one that captures the information relevant to statistical inference within the context of the model, including the size and composition of the units of study. A *sufficient unit* is one for which a random sample of aggregate statistics are sufficient in the statistical sense that validity of inference is demonstrable in comparison with analysis using units and variables at smaller levels of aggregation.

We will here apply statistical models that relate the dynamics of changing population size relative to resources in units of a size and composition sufficient to capture interactive effects with other variables. If the population variable were broken up, in these examples, into its components – births, morbidity, deaths, and migration – we can easily lose sight of primary effects such as Malthusian pressures and get lost in a mass of detail that entails particular areas of demographic study but fails to see the forest for the trees.

In particular, we will begin by focusing on large, “sufficient” regions that are relatively unaffected, for long periods, by exogenous shocks such as external wars, and we will apply models of how resources bear on reproduction, how population pressure on resources bears on social conflict or cooperation, on differential well-being, on growth or decline of social inequality, on social hierarchy, and on social mobility. Then, we continue, at the level of comparably large regional units,

to consider the rate of innovations relative to historical periods of growth, decline, or conflict. Lastly, we turn to how innovations and growth in transport systems within regions in such periods affect the augmentation of structural cohesion in regional trade networks and increase the capacity and likelihood of growing intraregional trade volumes, which connect dynamically to long-distance trade, and in some instances to innovations at the regional level such as monetization.

Using the sufficient statistics approach, scientists like Turchin have been able to take what can be transferred from the study of animal ecologies to specify and amplify problems in historical dynamics that do permit use of available data to examine hypotheses about the dynamics of change and periods of innovation. The key move is to look at the demographic variables per area (such as sheer number of people) as they bear on available resources. Large bounded areas may serve as a control for migration if there is little migration in and out of the region as a whole. Total numbers of people per area – and how these numbers change relative to resources – allow empirical study by means of models of direct and delayed *density-dependent* mechanisms of change (Turchin, 2003, p. 138). In this context, *endogenous* refers to *density-dependent* feedback mechanisms (e.g., population/resources; predator/prey) and *exogenous* refers to such factors as affect population (e.g., climate, carrying capacity) but are not affected by population growth rates in the short to medium run.

The British historical records, comprehensively analyzed by the Cambridge Group for the History of Population and Social Structure, now provide “sufficient statistics” that are aggregated over successive intervals for total populations, total bushels of wheat, etc., to provide Turchin and others with a plenitude of data for successive periods in English history. Historical research on the Roman Empire and various dynasties in China provides a similar breadth of data that enables us to compare regional statistics over successive historical periods.

Access to these kinds of sources provided relatively reliable regional data on *density-dependent* changes in population/resource pressure and sociopolitical instabilities (the so-called socio-political violence index, SPI) for Turchin (2005b, 2006) to identify *endogenous dynamics* in a number of Eurasian populations. His findings on interactive dynamics between population pressure and sociopolitical instability and conflict have been replicated by Kohler, Cole and Ciupe (2009) for one of the best-studied regions of the Southwestern Pueblos. Turchin (2005b, p. 10) notes for his two-equation time-lagged regression that “the statistical approach . . . can yield valuable insights into the feedback structure characterizing the interactions between different aspects of the studied dynamical system – when data are reasonably plentiful (for example, at least two or three complete oscillations), cover different aspects of the system, and the measurement errors are not too large.”

Such a two-equation time series regression is useful when there is a known length of time lag, like generational reproduction, for one variable to affect the other positively (e.g., population pressure → transfer of competition levels into overt sociopolitical conflict), while fit to the second equation examines the reverse (negative) lag (e.g., sociopolitical conflict → lowered reproduction reducing population pressure in the next generation) to check for an endogenous time-lagged feedback loop. The

competing hypothesis is that each variable has an endogenous cycle length and so is predictable from an inertial rather than interactive feedback process (Turchin, 2005b, p. 16).

Turchin's tests of the alternative hypotheses show that interactive rather than inertial dynamics account for the historical data for England and the Han and Tang dynasties in China. Panels (a) and (b) in Fig. 5.8, reprinted from Turchin (2005a), show how fluctuations in population pressure (i.e., population divided by bushels of grain) lead those of SPI by about a generation. Each period repeats a similar endogenous dynamic, defined by the negative time-lagged feedback between P , population density per resource, and SPI (data from Lee, 1931), fitting, with appropriate constants, a 2-equation model:

$$SPI_t \approx +P_{t-1} \text{ (Population change drives change in SPI)} \quad (5.3)$$

$$P_t \approx -SPI_{t-1} \text{ (SPI has a reverse effect on Population)} \quad (5.4)$$

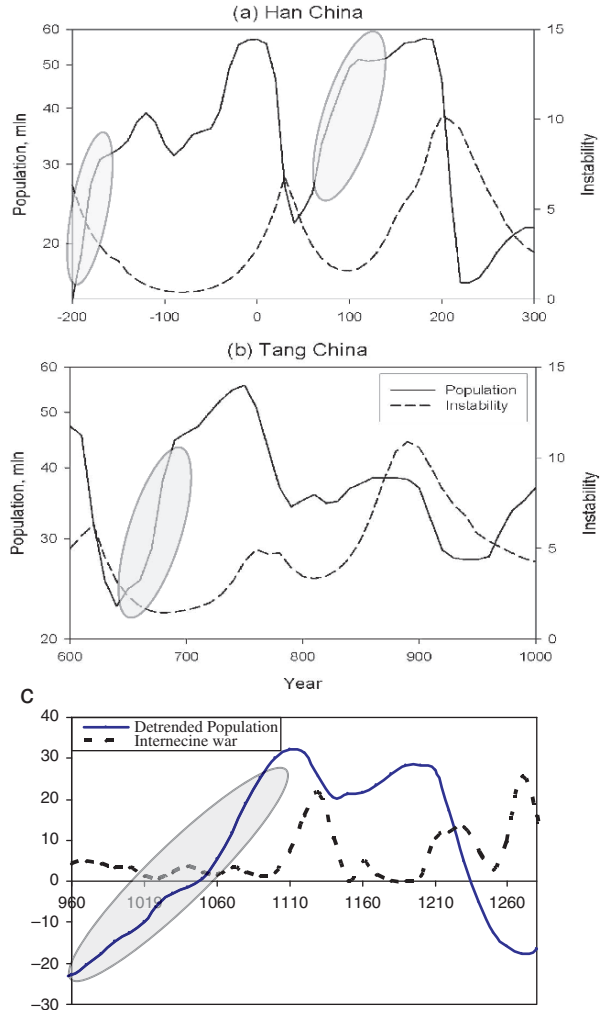
Part C of Fig. 5.9 shows that Turchin's model is compatible with the Song dynasty data (960–1279) on population (Zhao & Xie, 1988), and with Lee's (1931) SPI index of internecine wars. The time-lagged correlation is consistent with Turchin's cycles, and we see in the period 960–1250 CE that the sociopolitical violence index (SPI) lags population growth by a generation.²¹

We saw earlier in Equations (5.1) and (5.2) that the "shocks" of major wars, through sociopolitical instability W , tend to have an immediate effect on the relation between the size scaling of smaller and larger cities for Europe and China, a shock that affects q but not β . Without the effect of SPI, these two equations would predict positive feedback between β and q that would result in either a convergent or a divergent time series. It is only the "shock" index, W_i given the locations of SPI events in the Turchin dynamics (which we also estimated for Europe from historical data) that acts as an external shock and makes the predicted time series oscillatory, often clocking with Turchin's endogenous dynamic. With the warfare "shock" included, however, Fig. 5.10 shows how the Equations (5.1) and (5.2), applied recursively as a difference equation with only initial q and β values, form an oscillatory pattern between q and β that would eventually converge to equilibrium if there were no further wars. Optimized difference equations lead to an oscillating trajectory that converges to $q = 1.43$ and $\beta = 1.75$ for both China and Europe toward the middle *third* millennium, i.e., still a long time into the future, and with an assumption of no further dis-equilibrating wars.

The prediction from the Turchin model is that high rates of innovation occur in a relatively closed region, in the periods following a downturn in which total

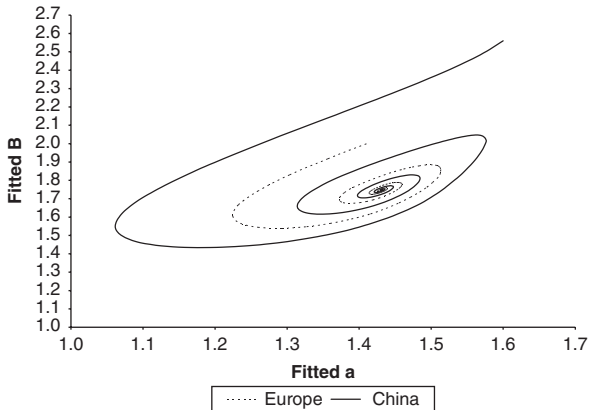
²¹ The research groups doing these population estimates (Zhao and Xie, 1988) and extrapolations (Korotayev et al., 2006), may have used peasant rebellions and wars, including invasions from without, to infer population decline, and the reconstructions of other authors will differ from these estimates, but the general pattern in Turchin's data for Han, Tang, England, Rome and elsewhere is that after dividing population by resources (e.g., bushels of grain), population shifts lead SPI shifts by about a generation.

Fig. 5.9 China's interactive dynamics of sociopolitical instability (broken curve for internecine wars, after Lee (1931) and population (solid curve: Zhao & Xie, 1988): (A) Han (200 BCE to 300 CE), (B) Tang (600 to 1000 CE), from Turchin (2005a), with population detrended by bushels of grain; and (C) Song period population (960 to 1279 CE), divided by successive trend values



population is rising faster than resources. This is reported in research by Nefedov (2003), Turchin and Nefedov (2008), and Korotayev, Malkov, & Khaltourina (2006). These authors find that the major periods of innovation in these historical fluctuations are the periods of *growth following prior collapse*, as shown by the shaded areas in Fig. 5.9. The prediction is also borne out for the dates of major innovations during the Song period (Temple, 1986; see Modelski & Thompson, 1996, pp. 131–133, 142–145, 160–176), as shown in part C of Fig. 5.9. Periods of innovation tend to be growth periods with ample resources and stability for a rebuilding and expansion of infrastructure and productivity, following depopulating and disruptive regional and city-system conflicts. The actual innovations implemented in these periods often utilize inventions that trace back to the periods of earlier cycles, many of which diffuse to other regions before they are adopted, often in different ways

Fig. 5.10 Phase diagram for difference equations (5.1) and (5.2), with exogenous shocks to 1970 and convergent tendencies modeled to 2500 CE, barring disequilibrium by warfare



than originally intended. This is well documented, for example, in the history of Chinese technological inventions that take hundreds of years to become European innovations (Temple, 1986), but the same occurs between other pairs of regions.

In a closer examination of the relation between innovation, Turchin or ‘secular’ (centuries-long) cycles, and trade networks, our historical project carried out a detailed European region study for 1100–1500 CE using the data of Spufford (2002), and found, like Spufford, that in medieval Europe, innovations were most pronounced not only in the periods of population rise relative to resources, following previous crises, and before the period of stagnation following the rise, but relate to *monetization* as one of the primary dynamics of innovation and economic reconfiguration (White & Spufford 2005).

Monetization had come earlier to Sung China, with its period of demographic rise, and later to Medieval Europe, with the spread from China, especially along the silk routes, involving credit and paper instruments and currencies. Demographic/conflict cycles may have been going on every time that new forms of money (M2, M3, M4, M5, etc.) developed. The relation to secular demographic/conflict cycles is that elites gain advantage as the population/resource ratio becomes Malthusian, labor becomes cheap, landlords and owners of productive property benefit, so that the rich grow richer and the poor poorer, if only for a period.

In Medieval Europe, elites exporting goods for cash in periods of scarcity used their monetary income to move to cities, hire wage retainers, kick peasants off their estates, fuel conspicuous urban demand for long-distance luxury trade, etc. The benefits of increased trade and profit margins for merchants were enormous, and, as the volume of trade expanded, merchant organizations, rural estates, and all kinds of other organizational structures expanded. In doing so, they recurrently surpassed thresholds at which innovative reorganizations of transport, production, accounting and all the other related areas occurred in nonlinear explosions.

In this study of Medieval European transformation (White & Spufford, 2005), we also identified how differences in network properties – such as betweenness and flow centrality for individual cities, or structural cohesion for regions – affected the trading benefits and growth of population or the wealth of different cities. The

growth of merchant capital and attendant innovations, for example, were affected in the short run by the relative betweenness centralities (Freeman, 1977) of cities in which the finance capital was affected in the longer run by flow centrality. As a general consequence, the building of roads and their contribution to improvement of structural cohesion in regional trade routes had massive effects on trade, innovation, and development (see Spufford, 2002).

Finally, to illustrate how the historical cycles approach applies to the contemporary problems of political and business leadership, one last case study of innovation cycles in the U.S. may be useful. Policymakers, rather than address global warming, the scarcity of nonrenewable fuels, overpopulation, and multiple insurgencies, have, until very recently, maintained a defensive, nationalistic, competitive, partisan, and expansionist approach characteristic of strictly segmentary organization, fostering the collapse of cohesive or cross-cutting integration at the national level and level of international alliances. Segmentary competition among elites and exploitation of the strategic advantage of the elite position in periods of scarcity is historically characteristic of large-scale polities that are in the resource shortage and crisis phase of Turchin's secular cycles (drafted in 2006, this sentence is consistent with Turchin's model of collapse for the U.S.A. in the present period and the actual financial meltdown of fall 2008).

There is a vast amount of innovation in the American economy, but it is directed at high-end medical technologies that will benefit elites over common citizens, at pharmaceuticals, at military technology, surveillance technology, and genetic and information technology battlegrounds over private ownership and patents and open source technology. We have a huge number of inventions that could be mobilized toward solutions of major world and national problems that go beyond issues of security and competitive rivalry over ownership. These inventions are not being mobilized to create the innovations that would solve these problems. Until recently, our business leaders, from Bill Gates to IBM, are concerned with stamping out the open source and collaborative open access technologies where innovations might occur outside the competitive frameworks of corporations, partisan political blocks, and the vested interests of elite factions.

5.5 Conclusions

This chapter has looked at structure and dynamics in human groups – and problems such as social organization and behavioral patterns – within a network formulation that is illustrative of aspects of social and historical theory that bear on concepts of innovation. Re-conceptualizing cohesion and hierarchy within a dynamical network framework so as to extend the foundations of social theory has been the central task of the chapter.

In spite of the ambiguity of social hierarchy, large-scale urban systems provide empirical and historical examples that validate the existence of social hierarchy but also explore major fluctuations and instabilities in the indicators of hierarchy. I argue that multinet and node/edge attribute analysis can be employed at

multiple spatiotemporal scales to improve our understanding of dynamics of hierarchy and cohesion, of how innovation is fostered, and how conflict – our variables SPI and SPI “shocks” (variable W) in city size dynamics – can work for or against one’s own group and, at times, for the benefit of those one chooses to make war against. I argue that such studies of the dynamics of hierarchy and cohesion are important to understanding innovation because they create the contexts for variable rates of occurrence of inventions and later innovations involving prior inventions.

Which conclusions about innovation can we draw from this chapter? We can see abundant evidence that there is a great deal of variability in pre-state and pre-urban innovation and transformation in the structures and dynamics of social organization and exchange, a diversity that tended to be constrained within the moral regulatory universes associated with kinship, obligation, reciprocity and reverse hierarchy. Although there is enormous human inventiveness and spread of innovations in forms of exchange, growth of “surplus” does not produce automatically the kinds of hierarchy seen in development of pristine states and their elites.

Instead, what is seen is the startup of networks of cities as ensembles with internal division of productive labor oriented toward a division of labor in the exchanges between cities, along with specialized roles that facilitate external predation and trade, on the one hand, and organizational channeling of production into export trade, on the other. Superadditivity in interaction of diverse problem solving approaches might be the key to urban competitiveness, but surfeits of inequality and competitiveness may also destabilize entire city systems, as shown in the historical investigations of Section 5.5. Temporally shifting imbalances between supply and demand, even in the absence of market pricing systems, are already implicit in these intercity trading/warfare systems. By any means, the attraction of population into cities to fuel labor demands as well as to staff specialized oversight, military, and other specialized roles, leads to a related set of instabilities, those between growing population and the eventuality of limited resources. This, in turn, leads to alternating periods of scarcity and plenty caused by population pressure, on the one hand, and the alternation of blowback effects of inequalities and internal sociopolitical instabilities reacting to these imbalances (population pressure producing scarcities, lowering of effective wages for workers while advantaging elites) on the other. Sociopolitical instability reacts to the eventual reduction in population pressure. This is the secular cycles model of agrarian state politics. In this chapter, I have sought to extend some of the findings and insights gained from such models to broader classes of problems where multi-relation and multi-attribute dynamics can be examined within a network approach to social theory.

Some of the work summarized here, on the historical and network dynamics in regional and interregional city systems, shows that urban system fluctuations and crises are at least negatively coupled through the sociopolitical instabilities (SPI) and shocks that follow from the structural-conflict cycles of agrarian states. On the network side, these studies also show strong interregional diffusion effects for inventions (e.g., China to Europe, as just one example) and eventually for innovations made possible by the diffusion of inventions, which are often put to very

different uses in later times and different regions. Further, the growth periods of these great swings in dynamics, in urban and regional economies, and in the consequences of politico-military interventions abroad, provide conditions and incentives for innovation.

If these results (and subsequent ones in this line of work) do show that the kinds of innovations that help solve problems in the crisis periods of civilizations and city systems do not come into play in those actual crisis periods where they are needed, but usually only after the fall following the crisis, then my contention would be that we need to rethink our understanding of the dynamics that affect innovation and the policy implications that we derive from them.

Further, conflicts that serve as shocks to city systems seem, according to quantitative historical dynamic studies, to drive city systems away from settling into a more stable equilibrium between the smaller cities that often have been the most dramatic sites of creativity, and the largest cities that often seem to draw out the economic vitality of larger regions. One is left to consider and to hypothesize the potentially beneficial effects of reducing sources of conflict while increasing diversity, especially in contexts where peaceful diversity can lead to cooperation and greater problem solving rather than innovation mainly for the sake of competitive advantage. If the approaches reviewed in this chapter have been productive beyond expectations, perhaps they will prove efficacious in investigations that could shape avenues toward a more benign set of futures than the ones we face now.

Acknowledgments Special thanks to David Lane, for his suggestions on organizing this chapter, and to Sander van der Leeuw, for rewriting the introduction and doing a superlative job in editing.

References

- Abbott, A. (2001). *Chaos of disciplines*. Chicago: University of Chicago Press.
- Adams, A. E. (1981). *Heartland of cities*. Chicago, IL: University of Chicago Press.
- Algaze, G. A. (2005). The Sumerian takeoff. *Structure and Dynamics*, 1(1), 5–48.
- Arthur, W. B. (1994). *Increasing returns and path dependence in the economy*. Ann Arbor, MI: University of Michigan Press.
- Bak, P. (1996). *How nature works: The science of self-organized criticality*. New York, NY: Springer-Verlag.
- Barabási, A.-L., & Albert, R. (1999). Emergence of scaling in random networks. *Science*, 286(5439), 509–512.
- Batty, M. (2006). Rank clocks. *Nature (Letters)*, 444, 592–596.
- Bettencourt, L. M. A., Lobo, J., Helbing, D., Kühnert, C., & West, G. B. (2007). Growth, innovation, scaling, and the pace of life in cities. *Proceedings of the National Academy of Sciences*, 104(17), 7301–7306.
- Boehm, C. (1993). Egalitarian behavior and reverse dominance hierarchy. *Current Anthropology*, 34(3), 227–254.
- Boehm, C. (2001). *Hierarchy in the forest: The evolution of egalitarian behavior*. Cambridge, MA: Harvard University press.
- Boltanski, L., & Thévenot, L. (1999). The sociology of critical capacity. *European Journal of Social Theory*, 2(3), 359–377.

- Brudner, L. A., & White, D. R. (1997). Class, property and structural endogamy: Visualizing networked histories. *Theory and Society*, 26, 161–208.
- Chandler, T. (1987). *Four thousand years of urban growth: An historical census*. Lewiston, NY: Edwin Mellon Press.
- Chase-Dunn, C., Niemeyer, R., Alvarez, A., Inoue, H., Lawrence, K., & Carlson, A. (2006). *When north-south relations were east-west: urban and empire synchrony (500 BCE–1500 CE)*. Paper presented at the International Studies Association meeting, San Diego, CA.
- Du, H., White, D. R., Li, S., Jin, X., & Feldman, M. W. (n.d.) “Network Models for Chinese Women’s Rural-Urban Migration Experience” Ms.
- Farmer, J. D., Patelli, P., & Zovko, I. (2005). The predictive power of zero intelligence in financial markets. *Proceedings of the National Academy of Sciences*, 102(6), 2254–2259.
- Fitzgerald, W. (2004). *Structural and non-structural approaches to social class: An empirical investigation*. PhD Dissertation, Program in Social Networks. University of California Irvine.
- Florida, R. L. (2005). *The flight of the creative class: The new global competition for talent*. New York: HarperCollins.
- Freeman, L. C. (1977). A set of measures of centrality based on betweenness. *Sociometry*, 40, 35–41.
- Grannis, R. (1998). The importance of trivial streets: residential streets and residential segregation. *American Journal of Sociology*, 103(6), 1530–1564.
- Griffin, A. F., & Stanish, C. (2007). An agent-based model of prehistoric settlement patterns and political consolidation in the Lake Titicaca Basin of Peru and Bolivia. *Structure and Dynamics*, 2(2), Article 2.
- Henrich, J. (2001). Cultural transmission and the diffusion of innovations: adoption dynamics indicate that biased cultural transmission is the predominate force in behavioral change and much of sociocultural evolution. *American Anthropologist*, 103, 992–1013.
- Henrich, J., Boyd, R., Bowles, S., Gintis, H., Fehr, E., Camerer, C., et al. (2005). ‘Economic man’ in cross-cultural perspective: ethnography and experiments from 15 small-scale societies. *Behavioral and Brain Sciences*, 28, 795–855.
- Kearns, M., Suri, S., & Montfort, N. (2006). An experimental study of the coloring problem on human subject networks. *Science*, 313(5788), 824–827.
- Kohler, T. A., Cole, S., & Ciupe, S. (2009). Population and warfare: a test of the Turchin Model in Puebloan societies. In S. Shennan (Ed.), *Pattern and process in cultural evolution*. Berkeley, CA: University of California Press. Also as a Sante Fe Institute Working Paper, 06-06-018, Preprint: <http://www.santafe.edu/research/publications/wpabstract/200606018>
- Kondratieff, N. D. (1984). *The long wave cycle* (G. Daniels, Trans.). New York: Richardson and Snyder.
- Korotayev, A. V., Malkov, A., & Khaltourina, D. A. (2006). *Introduction to social macrodynamics: Secular cycles and millennial trends*. Moscow, Russia: URSS Press.
- Kremer, M. (1993). Population growth and technological change: One million B.C. to 1990. *The Quarterly Journal of Economics*, 108(3), 681–716.
- Krugman, P. (1996). Confronting the mystery of urban hierarchy, *Journal of Political Economies*, 10(4), 399–418.
- Lee, J. S. (1931). The periodic recurrence of internecine wars in China. *The China Journal* (March–April): 111–163.
- Llewellyn, K. N., & Hoebel, E. A. (1961). *The Cheyenne way: Conflict and case law in primitive jurisprudence*. Norman, OK: University of Oklahoma.
- Lomnitz, L. (1977). *Networks of marginality: Life in a Mexican Shantytown*. New York: Academic Press.
- Low, S. M. (2003). *Behind the gates: Life, security and the pursuit of happiness in fortress America*. New York: Routledge.
- McFadden, D. (1973). Conditional logit analysis of qualitative choice behavior. In P. Zarembka (Ed.), *Frontiers in econometrics* (pp. 105–142). New York: Academic Press.

- McFadden, D. (1981). Econometric models of probabilistic choice. In C. Manski & D. McFadden (Eds.), *Structural analysis of discrete data: With econometric applications* (pp. 198–272). Cambridge, MA: MIT Press.
- Modelski, G., & Thompson, W.R. (1996). *Leading sectors and world powers: The co-evolution of global economics and politics*. Columbia, SC: University of South Carolina Press.
- Moody, J., & White, D. R. (2003). Social cohesion and embeddedness: A hierarchical conception of social groups. *American Sociological Review*, 68(1), 1–25.
- Morel, B., & Ramanujam, R. (1999). Through the looking glass of complexity: the dynamics of organizations as adaptive and evolving systems. *Organization Science*, 10(3), 278–293.
- Nakano, T., & White, D.R. (2006a). The large-scale network of a Tokyo industrial district: small-world, scale-free, or depth hierarchy? *COI Working Paper June 2006*. New York: Center on Organizational Innovation, Institute for Social and Economic Research and Policy, Columbia University.
- Nakano, T., & White, D. R. (2006b). *The “visible hand” in a production chain market: A market equilibrium from network analytical perspective*. Santa Fe, NM: The Santa Fe Institute. (Santa Fe Institute Working Paper 06-05-015).
- Nakano, T., & White, D. R. (2007). Network-biased pricing mechanisms of complex production-chain markets: Positions and roles affecting the equilibrium. *Structure and Dynamics*, 2(4), 130–154.
- Nefedov, S. A. (2003). Theory of demographic cycles and social evolution of the ancient and medieval societies of the East. *Vostok Oriens*, 3, 5–22 (the English translation is available at <http://repositories.cdlib.org/imbs/socdyn/wp/wp4/>).
- Page, S. E. (2007). *The difference: How the power of diversity creates better groups, firms, schools, and societies*. Princeton, NJ: Princeton University Press.
- Pearl, J. (2000). *Causality*. Cambridge, UK: Cambridge University Press.
- Pearl, J. (2007). Comments on *A unified framework for defining and identifying causal effects*. H. White & K. Chalac, University of California Complexity Videoconference. [http://eclectic.ss.uci.edu/~drwhite/center/ppt.pdf/Pearl\(2007-Feb09\)_Cause.pdf](http://eclectic.ss.uci.edu/~drwhite/center/ppt.pdf/Pearl(2007-Feb09)_Cause.pdf)
- Powell, W. W., White, D. R., Koput, K., & Owen-Smith, J. (2005). Network dynamics and field evolution: The growth of interorganizational collaboration in the life sciences. *American Journal of Sociology*, 110:1132–1205.
- Shalizi, C. (2007). *Maximal likelihood estimation for q-exponential distributions*. http://arxiv.org/PS_cache/math/pdf/0701/0701854v2.pdf
- Spufford, P. (2002). *Power and profit: The merchant in Medieval Europe*. Cambridge, UK: Cambridge University Press.
- Tainter, J. A. (2000). Problem solving: complexity, history, sustainability. *Population and Environment*, 22, 3–41.
- Temple, R. (1986). *The genius of China: 3,000 years of science, discovery, and invention*. New York: Simon and Schuster.
- Thoden van Velsen, H. U. E., & van Wetering, W. (1960). Residence, power groups and intra-societal aggression. *International Archives of Ethnography*, 49, 169–200.
- Turchin, P. (2003). *Historical dynamics: Why states rise and fall*. Cambridge, UK: Cambridge University Press.
- Turchin, P. (2005a). Dynamical feedbacks between population growth and sociopolitical instability in agrarian states. *Structure and Dynamics*, 1(1), 49–69.
- Turchin, P. (2005b). A primer on statistical analysis of dynamical systems in historical social sciences (with a particular emphasis on secular cycles). *Structure and Dynamics*, 1(1), 70–81.
- Turchin, P. (2006). *War and peace and war: The life cycles of imperial nations*. New York: Pi Press.
- Turchin, P., & Hall, T. D. (2003). Spatial synchrony among and within world-systems: insights from theoretical ecology. *Journal of World-Systems Research*, 9, 37–66.
- Turchin, P., & Nefedov, S. (2008). *Secular Cycles*. In contract for publication. Princeton, NJ: Princeton University Press
- van Duijn, J. J. (1983). *The long wave in economic life*. London, UK: Allen and Unwin.

- von Foerster, H., Mora, P. M., & Amiot, L. W. (1960). Doomsday: Friday, 13 November, A.D., 2026. *Science*, *132*, 1291–1295.
- Velleius, C. (1924). *Velleius Paterculus: The Roman History*. (F. W. Shipley, Trans.), Cambridge, MA: Loeb Classical Library.
- Walla, J. D., & Przeworski, M. (2000). When did the human population size start increasing? *Genetics*, *155*, 1865–1874.
- Watts, D., & Strogatz, S. (1998). Collective dynamics of ‘small-world’ networks. *Nature*, *393*, 440–442.
- Weissner, P., & Tumu, A. (1998). The capacity and constraints of kinship in the development of the Enga Tee Ceremonial exchange network (Papua New Guinea Highlands). In T. Schweizer & D. R. White (Eds.), *Kinship, networks, and exchange*. (pp. 277–302). Structural Analysis in the Social Sciences Series. New York and Cambridge: Cambridge University Press.
- White, D. R. (1969). *Cooperation and decision making among North American Indians*. Ann Arbor, MI: Dissertation Reprints.
- White, D. R. (1996). Enfoque de redes al estudio de comunidades urbanas. *Estudios Demográficos y Urbanos*, *26*, 303–326.
- White, D. R. (2004). Ring cohesion theory in marriage and social networks. *Mathématiques et Sciences Humaines*, *168*, 5–28.
- White, D. R., & Harary, F. (2001). The cohesiveness of blocks in social networks: node connectivity and conditional density. *Sociological Methodology*, *31*, 305–359.
- White, D. R., & Johansen, U. C. (2005). *Network analysis and ethnographic problems: Process models of a Turkish nomad clan*. Lanham, MD: Lexington Books.
- White, D. R., & Newman, M. E. J. (2001). Fast approximation algorithms for finding node-independent paths in networks. *Santa Fe Institute Working Papers 7035*, <http://www.santafe.edu/sfi/publications/wpabstract/200107035>
- White, D. R., & Johansen, P. (2005). *Civilizations as dynamic networks: Medieval to modern*. UCI Irvine Social Dynamics and Complexity Working Paper.
- White, D. R., Tambayong, L., & Kejžar, N. (2007). Oscillatory dynamics of city-size distributions in world historical systems. In G. Modelski, T. Devezas & W. Thompson (Eds.), *Globalization as evolutionary process: Modeling, simulating, and forecasting global change*. pp. 190–225. London: Routledge.
- White, H. L. (2007). A comparison of Pearl’s causal models and settlable systems. *UCSD Department of Economics Working Paper W32*.
- White, H. L., & Chalak, K. (2007). A unified framework for defining and identifying causal effects. Submitted to *Journal of Machine Learning Research*.
- Wright, H. (2004). *Raising civilization: Three Stanislaw Ulam memorial lectures on DVD*. Santa Fe, NM: Santa Fe Institute.
- Wright, H. (2006). Atlas of chiefdoms and early states. *Structure and Dynamics*, *1*(4), 738–758.
- Zhao, W., & Xie, S.-Z. (1988). *History of Chinese population*. Peking, China: People’s Publisher (in Chinese).

Part II
Innovation and Urban Systems

Chapter 6

The Organization of Urban Systems

Anne Bretagnolle, Denise Pumain and Céline Vacchiani-Marcuzzo

6.1 Introduction

Cities¹ are a major form of the material, social, and symbolic organization of societies. They are persistent and adaptive structures that fulfill a variety of social functionalities: habitat, production, services, political control over people and territories, as well as technical and symbolic mediation between nature and culture, groups and individuals. For a long time, and especially since the first industrial revolution, they have been both the spaces where most innovations occur and spaces for which most innovations are designed (Chapter 8). Location in space and time has to be considered as a very important feature for these systems, since societal evolution is much more rapid than its biological counterpart is. Being adaptive, cities alter through a variety of intentional actions that may produce non-intentional persistent features: urban structures are produced partly through design and planning, partly through self-organization. This chapter is an analysis of a conception of urban systems as complex systems. We develop a theory of the organization of cities and systems of cities as multilevel networks, including two main observable levels, where specific emergent properties can be observed. Compared to the earlier seminal conception by Berry (1964) of “cities as systems within systems of cities,” we emphasize a multilevel organization perspective and specific dynamic features that produce this structure. We also add an evolutionary perspective, which is based on observations about the way cities co-evolve within urban systems, through a variety of social interactions. Empirical evidence of the corresponding patterns and processes is provided for three main styles of urbanization, exemplified in four different parts of the world.

A. Bretagnolle (✉)

Université Paris I Panthéon Sorbonne, UFR de Géographie, Paris, France

¹ The word “city” is used here to designate generic urban entity defined as a geographically and functionally consistent urbanized area, whatever the administrative or political boundaries.

6.2 Cities and Systems of Cities as Different Levels of Social Organization

The organization of urban systems can be described on three main levels (Pumain, 2006): the micro level represents elementary units (individual people, firms, institutions) that are living together in a city, the meso level corresponds to the city itself (defined as a consistent geographical entity), and the macro level is the system of cities, made up of a large number of towns and cities which interact under unified control (in a national political territory or a global economic network). This organization is shaped by interactions operating on different spatial and temporal scales of observation.

6.2.1 *The City: A Collective Evolving Entity*

In a very general way, a city can be seen as a collective entity whose specific properties, although mostly produced by intentional agents at the individual level, cannot be simply explained or predicted from these intentions, nor derived by summing the characteristics of its inhabitants. Concepts such as urban function, centrality, or morphology were invented by urban scientists to understand cities' emergent properties. Produced through reiterated interactions between individuals, these are only defined at the aggregate level of the city and cannot be estimated from the mere summing of the attributes of the individuals that compose the urban entity. The collective character of a city also emerges over historical time because the evolution of cities is extremely coherent. Each of them has a specific history that defines its identity. Even if they participate in the same historical events and trends, cities each follow an original trajectory, which is influenced mainly by some of their political options or by their successes and failures in the process of their socio-economic co-evolution. The concept of a city as an urban entity is a relatively autonomous and persistent system of locally dense and frequent daily interactions. It is rooted (as a material and a symbolic entity) in the collective process of specification and identification. This is constrained by two distinct trends, bottom up by way of internal interactions (among urban citizens, firms, and institutions) and top down by way of external interactions (among competing cities), that form interlocking networks.

6.2.1.1 The "One-Hour" Traveling Time Constrains the City's Development

For a relevant analysis of the dynamics of this complex system, we have to translate this concept into proper measurements of city size that are comparable over space and through historical times. Measures of city size are difficult because cities are expanding not only inside the boundaries of a fixed territory, but over these limits. We have then to consider an entity according to a common reference in time-space. In this framework, we can visualize a city as the envelope for the daily activities

of its inhabitants and the buildings hosting them. The concept of a city that we use is not restricted to elementary political entities governed by a municipality or an administrative circumscription with fixed boundaries. Indeed, in the course of time, it is frequently observed that urban growth crosses these boundaries and spills over into neighboring circumscriptions, sometimes including other nearby cities and towns in “conurbations,” as noted by Patrick Geddes as early as 1905. In order to correctly identify a city as a consistent geographical entity, and considering that its spatial expansion as well as its *in situ* development are part of its growth, we define it as a “daily urban system” (this concept was the basis for the definition of Standard Metropolitan Areas in the US in the fifties, and its name has been chosen because it allows for frequent social interactions across the boundary, which usually take place within one day). Recent surveys in European, American, or African cities demonstrate that during a working day, an urban dweller typically needs to connect to three or four different places of activity, and devotes on average about one hour each day to traveling for these purposes (Crozet & Joly, 2004). It seems that this travel time, representing about 15% of a 24 hour working day, has remained rather stable over centuries (in the literature on transport, this is known as “Zahavi’s law”).

While the maximum spatial development of cities seems always to have been constrained by this typical one-hour time-budget, improvements in the speed of the means of transportation available have enabled the average commuting distance between places of work or urban services and places of residence, to be multiplied by ten in the last two centuries. In London, for instance, the city maps of the 17th century (Hollar, 1667) represent an urban agglomeration whose radius is 3 to 4 km. This radius has evolved (6 to 8 km in 1830, 20 km in 1900) according to the innovations in transportation technologies, for which London was a pioneer in Europe (1829, horse omnibus; 1836, steam railway; 1863, steam metropolitan; 1905, electric metropolitan). Today, a remote sensing image of London representing built-up areas (*CORINE Landcover*, EEA, 2000) shows a first perimeter of compact constructions over a circle of 20–30 km radius, with expanding branching developments reaching out to zones located nearly 70–80 km from the centre (Fig. 6.1). The historical adaptation of an urban entity, whose expansion was driven by economic success and social pressure towards the maximum possible spatial extension of the time (and such despite policies aiming at “the containment of urban England” – Hall, Gracey, & Drewry, 1973), reveals the general limitation of a one-hour traveling time, which constrained its development. A behavioral parameter defined for spatial interactions at the individual level is reflected in the organization of the urban entity at a higher level.

6.2.2 The System of Co-Evolving Cities

Cities never actually develop as isolated entities; they are always engaged in many types of interactions with other cities. They are linked through a variety of social networks, among which there are not only visible infrastructures, such as roads,

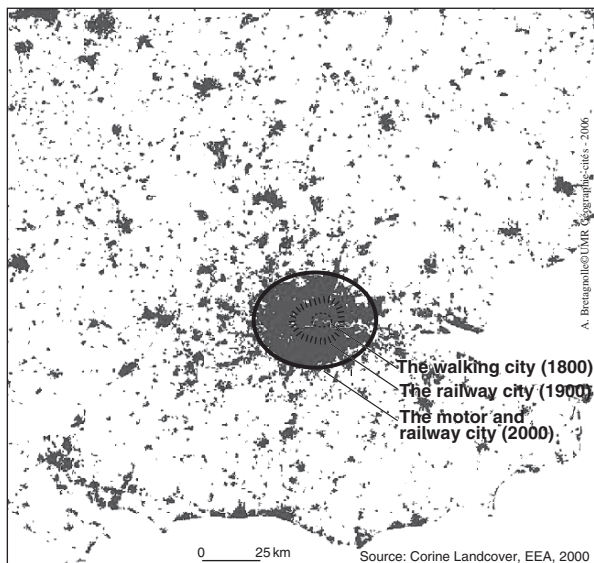


Fig. 6.1 The agglomeration of London, an evolving spatial entity delimited by a “one-hour” constraint of traveling time

railways or airlines, and the kind of exchanges usually accounted for, such as population migrations or trade of goods, but also more “invisible” networks: capital and information flows, for example. (The latter are more involved particularly in the process of innovation diffusion.) The resulting *systems of cities*, although mostly self-organized, ensure the same social functionality: to control territories or networks. The nature of the control has evolved through historical times, mainly from political to economic. The political control involved, which at first may have coincided with the scope of the market area of one city, spread to kingdoms or empires and then to national territories associating several cities. Economic control, expressed first by local entrepreneurs, broadened to national then multinational firms. Individual cities, which in the past were the main actors in their own development – seeking to overcome the limitations of their immediate environment by building networks to exploit distant resources – have become instruments of control of wider territories and networks, a control that is now assumed by national or supranational actors using the networks that cities have built through their interactions. *What we call systems of cities are evolutionary objects that may include subsets of cities connected by long-distance networks or cities belonging to unified political territories.* The general trend is a historical increase in the number of cities that are integrated, through ever more intense and more frequent interactions, although political events or economic crises may reduce their number and impact locally for some time. The precise identification of systems of cities is very difficult, due to the changing nature of the interactions that need to be considered, and the fluctuations in their spatial extension.

6.2.2.1 An Integrated System of Co-Evolving Cities Through Multilevel Spatial Interactions

The direct interactions, which could be named *first order interactions*, produce strong interdependencies in the evolution of cities and give rise to a macro-organizational level: the system of cities. This organization of “city systems” that develop specific emergent properties was noticed a long time ago (the first mention of a “system of cities” can be found in the writings of a French Saint-Simonian engineer, Jean Reynaud, as early as 1841) and now is part of a geographical theory on urban systems (Berry, 1964; Pred, 1977; Pumain, 1997). Systems of cities always show a differentiation of city sizes according to several orders of magnitude (today, from a few thousands to tens of millions of inhabitants), which follow a very regular statistical distribution, which is lognormal or of the Pareto-Zipf type. This hierarchy of sizes also corresponds to a hierarchy of urban functions and to more or less regular settlement patterns. These regularities were summarized, for a while, by central place theory (Christaller, 1933; Berry, 1967), which became embedded into a more general evolutionary theory of urban systems. Systems of cities also are characterized by their functional diversity and the co-evolution of their socio-economic profiles (see Chapter 8), as well as by distributed growth that can be summarized in first approximation by a Gibrat model (1931).

These first order exchanges, although very often reciprocal, are not fully symmetric. Asymmetries in the interaction flows between cities, when reiterated, produce a diversity of quantitative and qualitative differences in terms of city size, economic specialization, social composition, cultural features and urban landscapes. These differences lead to *second order interactions*, which are effects of *selection*, constraining a city’s development according to its rank in the hierarchy of city sizes, its economic specialization, and its “image” in the individual and collective representations of agents. Second order interactions can be observed indirectly by comparing the evolution of qualitative urban features and measuring the evolution of relative sizes in large sets of interacting cities (Pumain, 2006). More specifically, some types of networking are, at least momentarily, restricted to small samples of very large or very specialized cities.

In order to delineate a system of cities geographically, according to a “classic” definition of a “system,” we have to identify a set of interacting urban entities, including cities that have more intense or more frequent interactions with cities inside the system than with cities outside. This is often implicit when systems of cities are considered within the boundaries of national territories. National boundaries delimit a community of political and social rules as well as cultural features that reduce external interactions. However, even in the frame of this static view, cities can exceed these limitations, according to their size and function. In particular, national capitals or cities that are specialized in international activities have a broader range of exchanges than smaller or less specific towns. This wider opening of some cities can be observed in their first order interactions, but, most of the time, it is better illustrated by the interdependencies in their evolution that are produced by second order interactions. Therefore, the envelope of systems of cities should not

be seen as a strictly delimited boundary, but rather as a membrane that is more or less permeable according to the size and function of the cities inside the system. Because the range of interactions usually is strongly influenced by city size, the ability to develop outside interactions and to co-evolve with other urban systems is closely related to the hierarchical differentiation inside the urban system. But there often are, in a system, a few smaller cities that specialize in international activities, for example finance or tourism.

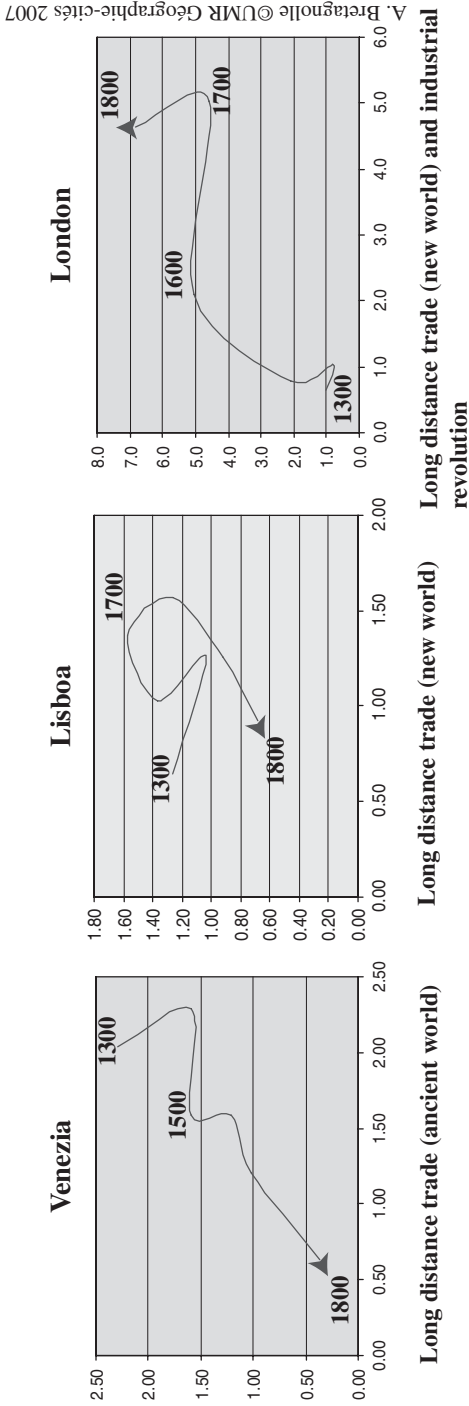
6.2.2.2 Global Cities Since the Middle Ages

At all times, a few outstanding cities have dominated the exchanges in the interconnected “world” formed by systems of cities, crossing the political boundaries. Considering an integrated “system of cities,” defined as a coherent set of cities already engaged in commercial and political competition, it is possible to follow the evolution of the weight of a given city within the system over several centuries and interpret its trajectory in relation to the greater or lesser success of the city in having a share in successive innovations, whether or not these lead to specialization. A few cities, for long periods of time, exhibit coefficients of allometric growth that remain systematically greater than one, when compared to the evolution of the other cities. This dynamic behavior is explained by the fact that these cities have much larger interaction networks than those of the other cities. Figure 6.2 presents three examples of such “global cities” that bypass the European urban field at three historical periods. We have computed the weights of Venice, Lisbon and London within the European total urban population at each date and represented the trajectory of this relative weight in phase space (ordinate: the weight of the city in the system for a given year; abscissa the weight² for the city the preceding year). The trajectories are illustrative of the successive success of Venice, in Mediterranean and Hanseatic trade during Middle Ages, then Lisboa in Atlantic maritime trade with the New World in the 15th century, and then London as the center of a colonial Empire and the industrial revolution at the end of the 18th century.

6.2.3 A Crucial Point in Urban Ontology: The Consistent Delimitation of Urban and Territorial Entities

The dynamics of urban systems have been interpreted in many different or even contradictory ways because the underlying observations of the differential growth of cities are based on different measures from one author to another, varying with the chosen definition of the city, the delimitation of the territory involved, and the samples of cities selected within these territories. Our theoretical conception of systems of cities makes it possible to gain a better understanding of the dynamics

² Europe can be considered approximately as forming a system of cities since the reopening of long distance trade in 12th or 13th century (cf. Pirenne 1927)



X-axis: $x = P_i, t / P_u, t$ P_i, t : population of the city i at time t
 Y-axis: $y = P_i, t+1 / P_u, t+1$ P_u, t : total population of the system of cities at time t

Fig. 6.2 Trajectories of individual cities in the European urban system. The ordinate is the weight of the city in the system for a given year; the abscissa is the weight for the city the preceding year

involved by relating them to their social functionalities and to the artifacts that make these functionalities feasible. As we have seen, at the scale of the city, the main functionality is the coordination of day-to-day activities, while at the scale of a city system the issue is the control of a territory or multiple networks, mainly by means of political power at the outset, and economic power subsequently. We have defined comparable databases within the frame of national or continental boundaries and we shall return to the consequences of this choice (Chapter 12).

Comparing city systems over time and in different areas of the world requires careful preparation of databases. These are never directly exploitable as they are collected and produced by statistical institutes. Official databases have their limitations, first, in terms of the indicators available. Since successive industrial revolutions have widened the gaps in standard of living, the best measure of the success of a city might be its economic power. Most of the time, we must be content with data on numbers of inhabitants of cities, the only indicator that can be mobilized on scales of time and space of this magnitude. However, even regarding population, harmonization is required so that entities termed “cities” are comparable from one country to another and through history (Pumain, Saint-Julien, Cattán, & Rozenblat, 1991; Bretagnolle et al., 2007).

Today census bodies use two main approaches to define the city. The first outlines the *urban agglomeration*, which is formed by the continuity of built-up area and by minimum population or density threshold values; the second outlines the *urban area*, which is much wider than the agglomeration since it also includes peri-urban rings that send part of their working population on a daily basis to the functional pole of the agglomeration. To make urban populations comparable from one country to another and at different times does not necessarily mean adopting identical criteria. For instance, it is considered today that in Europe the minimum functions that can be associated with a city characterize aggregates of more than 10,000 inhabitants, while in South Africa, more recently urbanized and involving a smaller surface area, it is more reasonable to lower this threshold to 5,000 inhabitants. Likewise, a strong feature of the USA is found in the extremely wide spatial encroachment of its cities, arising from the tendency of the population to undertake long-range daily commuting.³

When defined for highly integrated territories with intense, frequent interactions (for instance, with exceptions, most national states today), city systems present a certain number of characteristic emergent properties. Among the main “evolutionary laws,” we insist on four major trends: historical path dependence, competitive expansion and distributed growth, and reinforced urban hierarchy.

³ The different urban definitions that we have adopted in our comparative study are as follows: in Europe and India, agglomerations of more than 10,000 inhabitants; in South Africa, functional agglomerations of more than 5,000 inhabitants (including the white city and the non-white townships that are economically linked to the city by home-to-work commuting); in the USA, the populations of cities and towns up to 1940, then the *Standard Metropolitan Areas* (known today as the *Metropolitan Statistical Areas*) and the *Micropolitan Statistical Areas* up to 2000 (for a more detailed description of these databases, see Bretagnolle Pumain, & Vacchiani-Marcuzzo, 2007; Bretagnolle, Giraud, & Mathian, 2007).

6.3 Historical Path Dependence

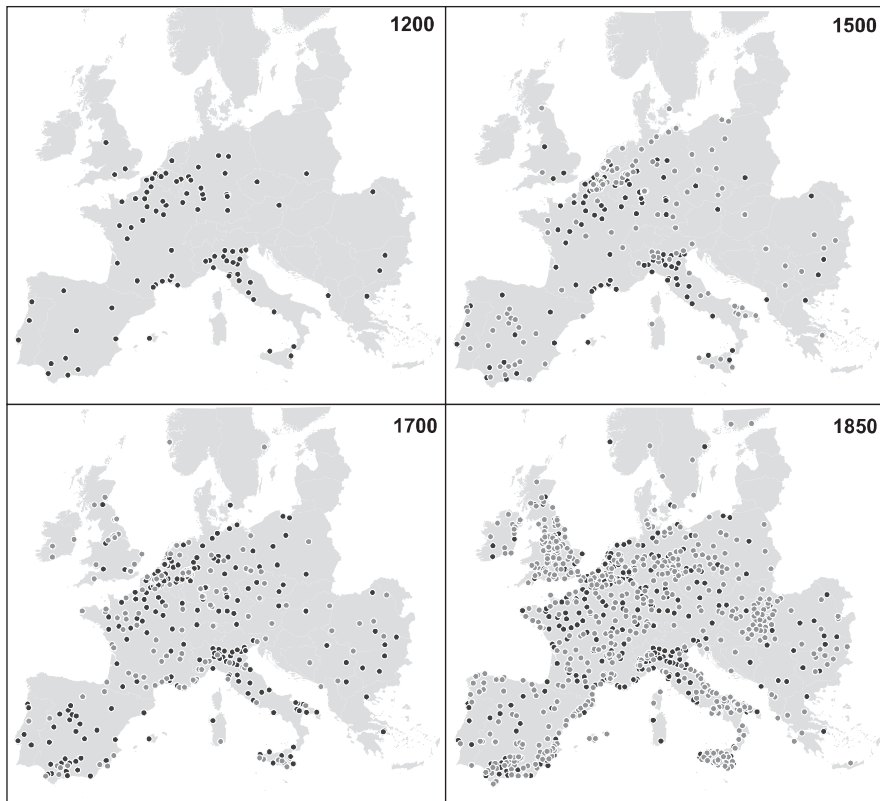
Despite the multiplicity and apparent diversity of national urban systems in the world, three major styles are recognizable in their hierarchical and spatial organization. These styles are differentiated because of different historical trajectories (Table 6.1). Their properties vary according to their period of emergence (technological conditions during the urban transition determine space-filling parameters) and according to any major exogenous impacts (such as colonization).

Table 6.1 Three major styles of urban systems in the world

Urban System	Typical world region	Selected example
Long-standing urbanization, slow and regular evolution	European urban system National urban systems in Europe	European urban system
Long-standing urbanization, major exogenous impact (colonialism)	Asia Black Africa	India
New countries and waves of urban creations	U.S.A. Australia South Africa	U.S.A, South Africa

6.3.1 Diversity in Morphogenesis

In countries with long-standing and continuous settlement processes, cities emerge more or less simultaneously all over the territory, and the city systems are characterized at once by the long-standing nature of their urbanization and by the regularity of their development over time. European countries are a good illustration of this. From Antiquity, long-range inter-city exchange networks became established, and several metropolises had more than 100,000 inhabitants (Ancient Athens) or even one million (Ancient Rome). After periods of relative stagnancy in the Early and Late Middle Ages, European cities were once again undergoing vigorous growth and the reactivation of exchange networks in the Early Modern period, following a “national” or “international” logic (at least for the largest or most specialized). Figure 6.3 shows a marked persistence of the distribution over geographical space of towns and cities since AD 1200 in the western part of Europe and since AD 1500 in eastern Germany and central Europe. The urban growth that is characteristic of the industrial revolution (the number of towns and cities doubled) did not significantly alter spatial distribution, even if there was a strong densification effect along the coasts and in the major mining and metalworking basins. A similar stability also characterizes the top of the urban hierarchy. The majority of the large European metropolises already dominated national networks since the end of the Middle Ages, despite the fact that the major economic centers holding sway over the Mediterranean in the Middle Ages gradually gave way to the trading cities on the



- Cites larger than 10 000 inh.
- Cites larger than 10 000 inh since the date of previous map.

Source: *Bairoch et al. 1988*
A. Bretagnolle ©UMR Géographie-cités - 2006

Fig. 6.3 First urban system style – continuity (Europe)

Atlantic coastline, until the industrial revolution sent Britain to the front (Braudel, 1967; de Vries, 1984).

A second style of urban system is represented by countries where urbanization is long-standing but where disruptions have been experienced, like Asia or Black Africa. The example of India is illustrated in Fig. 6.4. This country underwent a major phase of urbanization in Antiquity and a second phase in the Middle Ages with Muslim expansion. At these times, the country was divided into some twenty kingdoms, and the main cities in existence (Agra, Delhi) were located inland rather than on the coasts, because the main activities were related to domestic trade and territorial control. While the first Portuguese, Dutch, French, and later British trading posts were set up from the 16th century (Bombay in 1532, Madras in 1639, Calcutta in 1690), they did not have any decisive impact on the urban patterns in India before the 19th century. However, when India officially came under the British crown in 1847, the administrative and economic orientations changed radically and altered the distribution of Indian cities in a durable manner. From then on, the main

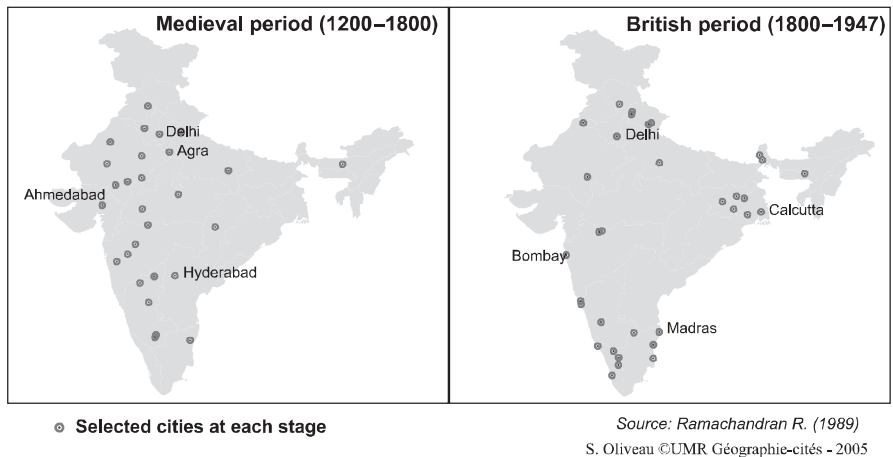


Fig. 6.4 Second urban system style – disruptions (India)

cities were created or developed along the coasts and the main rivers to facilitate exchanges with Britain.

The third type of urban system is that of “New World” countries, where towns and cities were imported by settlers and spread in successive waves of penetration, either driven or accompanied by canals and railways. In the USA, the occupation of space by cities took place in relation to waves of settlement moving in from the coastlines. A first frontier moved west, reaching the Mississippi in the 1850s, the Rockies in the 1870s, and the western coastline in the 1890s (Fig. 6.5). This frontier started from the first communities established in the 17th century along the east coast (Manhattan in 1614, Philadelphia in 1654) which developed slowly up to the start of the 19th century. In the 1790 census, only five cities were larger than 10,000 inhabitants, one of which was New York with only 33,000 inhabitants at this date (while London had reached 948,000 inhabitants). The exponential urban development that followed independence is out of all proportion with the growth that characterized Europe over the same period. In less than one century, New York became the second largest city in the world, with more than three million inhabitants in 1900. A second settlement front opened up in the 18th century on the Pacific coastline when the Spanish moved up from Mexico. The first towns of more than 5,000 inhabitants were San Francisco and Sacramento in 1850.

In South Africa, the process is slightly different: between 1652, when Cape Town was founded, and the start of the 20th century, several waves of Dutch and later British settlements spread along the coast towards the east, within a mainly agro-pastoral economy. The decisive impetus to the creation of urban settlement was the discovery of diamond resources (Kimberly in 1867) and gold (Johannesburg in 1883) in the central province of Gauteng, leading to a major shift of the center of gravity from the coastline towards the interior. From this period onwards, the cities grew more through internal migrations than through international immigration. The formation of the South African state in 1910 and the development of exchanges

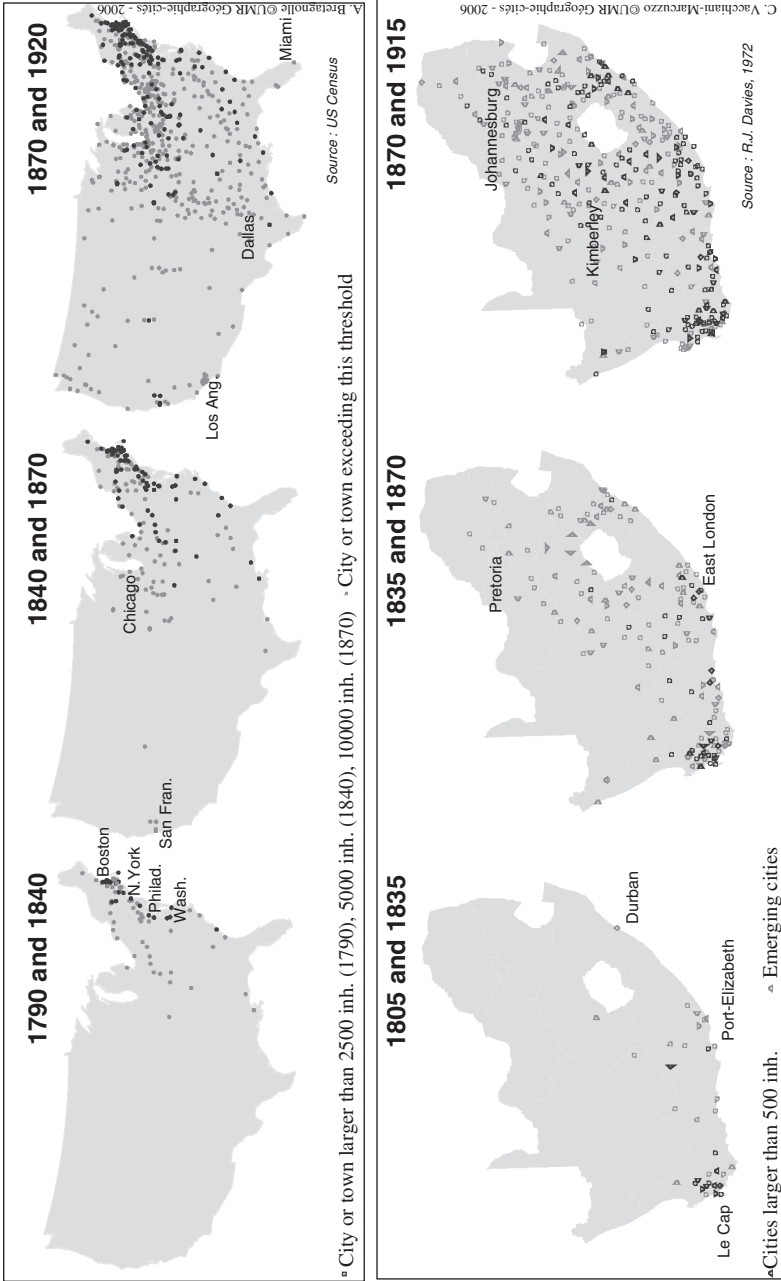


Fig. 6.5 Third urban system style – waves of urban creations (U.S.A., South Africa)

within an economy that had become industrial favored the emergence of a genuine city system in the 1950s (Fig. 6.5).

These three major types of morphogenesis do of course considerably oversimplify the diversity of situations observed worldwide. It would be interesting, in particular, to analyze in a more precise manner the particular place of the city systems in Latin America, North Africa, and the Middle East, as each of these areas had its own specific involvement in colonial processes. However, these three main types do make it possible to identify forms that still today most strongly differentiate hierarchical and spatial configurations in city systems across the world.

6.3.2 *Typical Emergent Properties of the Three Major Styles of Urban System*

The three types of system described above are not differentiated by their histories alone, but also by certain features of their hierarchical and spatial configurations that are still perceptible today.

A first feature that differentiates the three styles of systems is their morphology, i.e. the pattern of their occupation of space as measured by the density of the cities and the degree of hierarchization (Table 6.2). When cities were established at a time when the means of transportation were very slow, via the spontaneous emergence of agricultural markets or by the establishment of relay posts on the main communication routes, they were very close to one another. For example, in Europe and India alike, the average distance between two towns or cities is around 15 km. The newer countries show a wider spacing between towns and cities, marked concentration in the largest cities, and an urban hierarchy that is characterized by greater variation in size. As a convenient measure of the degree of hierarchical inequality, the slope of the straight line adjusting the rank-size distribution of cities and towns can be used. This parameter value is under one for Europe and India, while it is clearly above

Table 6.2 Typical parameters of the three major styles of urban systems

	Average distance to nearest neighbor in km (2000)	Inequalities in city sizes (2000)	Average annual growth rate (%) during the urban transition (19th or 20 th century)	Macrocephaly index (2000)
Long standing urbanization: Europe	15	0.94	1–2, distributed	2 to 3
Long standing and external shock: India	16	0.99	2–3, dual	4 to 8 (regional indices)
Recent systems:				
(a) U.S.A	36	1.20	3–4, higher on frontier	1.5
(b) South Africa	32	1.19	3–4, higher on frontier	2

NB: the degree of inequality is measured by the absolute value of the adjustment slope in the rank-size graphs.

one for the USA and South Africa. Indeed, in the latter two, countries, towns and cities developed in a pioneer logic, i.e. aiming to occupy the widest space possible, even if this was in an extensive manner, and they developed with faster and more efficient means of transportation (especially the railways). This results in systems where towns and cities are less numerous, less dependent upon the initial agricultural settlement, more widely spaced, sizes more contrasted, and where the largest cities can reach sizes greater than those observed in the Old World. The automobile later enabled a small number of urban centers to have an influence over very distant outer rings: for instance in 2000, the mean radius of functional areas is 46 km for the American metropolitan statistical areas (MSAs), as compared with 13 km for French urban areas.⁴

A second feature that differentiates the three styles of urban systems concerns the regimes of city growth. By analogy with demographic transition, what is known as urban transition (Zelinski, 1971) is a period of massive urban growth in the course of which, settlements, thus far made up of villages and fairly homogeneous and scattered, became much more heterogeneous by concentration around urban centers. This transition, which began at the time of the great industrial revolution, is complete in countries where it took place early (Europe, the USA), but it is still underway in countries in Asia (China, India) and is just starting in certain African countries. While the main influx of new urban settlers was above all from rural areas in the 19th century urban transition, it also was fed by demographic growth specific to the towns themselves in the 20th century transition. This increase in population (Table 6.2) is spread across towns and cities according to three main patterns:

1. In countries where urbanization is long-standing, and where it developed in the continuous mode (European type), urban growth was *distributed*, i.e. spread across all parts of the territory in a manner that is proportional to the size of the towns and cities, even if a slightly higher relative growth rate is seen in the large cities. Average growth rates were low throughout the 19th century, at around 1 to 2% per year (1.7 for London and 1.3 for Paris), except in certain localities in industrial basins (Bradford 4%, Valenciennes 3%).
2. In new countries, urban growth followed a “wave” settlement pattern, with markedly higher growth rates in the newly created towns and cities. In the USA, for instance, Chicago saw an annual growth rate of 12% per year between 1850 and 1870, and urban growth was more intense on average during the urban transition than in older countries (around 3 to 4% a year, mainly because of immigration from abroad).

⁴ Caution is required in comparing these figures: indeed the commuting range threshold taken into account to define functional areas is not the same in the USA (15%) and France (40%), and it refers to very different dimensional grids (a ratio of 100:1 on average between the surface areas of *counties* and that of the French *communes*). But the choices made by the census bodies involved is an actual reflection of the specific nature of urban settlement patterns in each country.

3. In systems that were reorganized following the impact of colonialism (some Asian and African countries), urban growth was *dual* i.e. fairly small and late endogenous growth of markets, administrative centers, and local artisan activity is superimposed on very marked concentration in the large cities, often the capital cities, or acting as an interface with the capital. This macrocephaly (highest ratio between the populations of two cities of consecutive rank) is typical: while the ratio is 1:2 or 1:3 on average in most European countries, it is 1:6 in Ivory Coast and 1:7 in Mali. In India, several regional capitals were driven by colonialism (Bombay, Calcutta, Delhi, with 16 million, 13 million and 12.5 million inhabitants respectively in 2001), and each is far ahead of the other cities in their respective regions (the populations of these cities today are four to eight times that of the largest city in their hinterland). In these developing countries, the mean growth rates during urban transition are high, often more than twice those observed in “older” countries, especially because that can coincide with demographic transition and its reduced urban mortality (Table 6.2).

6.4 Competitive Expansion and Distributed Growth

It could be thought that such marked historical differences would result in different patterns of urban dynamics, and that each main type of city system thus identified would evolve in its own specific way. In fact, once the city system is established and integrated into a political territory, the resemblances in the way they evolve are striking. Paradoxically, it can even be said that it is because all city systems, once formed, evolve in the same manner that they continue to carry the marks of their histories. Indeed, such marks are not the indication of the “inertia” of geographical structures, but conversely of their extraordinary ability to adapt.

6.4.1 Increase in Size and Number of Cities

Whatever the level of development, and however long-standing the urbanization, city systems always have been characterized by a tendency to grow both in terms of their urban population (the maximum city size increased from one million in 1800 to three million in 1900 and 30 million in 2000) and in terms of their number (Table 6.3).

For a total surface area of around 4.5 million km², Europe today has more than 5,000 cities, but the period of greatest expansion was between 1700 and 1800 (where the number of cities increased threefold) and then between 1800 and 1900 (where the number increased 2.5 times, while it barely doubled in the following century). India, where the urban transition began much later, made up the difference in a spectacularly short period between 1950 and 2000, since for a total surface area of 3.2 million km² the total number of cities is 3300 (the same density as Europe). In the USA, the progression of the number of cities is regular through the 19th and

Table 6.3 Evolution of the number of cities from 1900 to 2000 (Europe, India, USA, South Africa)

	1900	1950	2000	Surface (million km ²)
Europe	2532	3702	5123	4.8
India	580	1095	3285	3.2
United States	382	717	934	7.8
South Africa	14	41	307	1.2

Sources and data bases: Europe: Bairoch, Batou, and Chèvre (1988), Pinol (2003); India: Census, Oliveau (2005); United States: Census of the U.S., Bretagnolle, Giraud (2006); South Africa: Davies (1972), Vacchiani-Marcuzzo (2005).

20th centuries, the decrease in the rate of growth that can be seen from 1950 being due to grouping of contiguous urban units into the SMAs. South Africa is a case apart because of its very recent urbanization, and a city system that was not mature before 1950.

The important fact is that, whatever the country, over recent centuries, urban populations do grow faster on average than the number of cities, which results in greater concentration in urban settlement. The number of cities is still increasing in developing countries, but in highly urbanized countries, the appearance of new cities with populations above a certain threshold is slowing down, and the number of urban units can even decrease on account of the fusion to form larger units as a result of urban sprawl. It is likely that in a few decades, the stabilization of the overall world population will be reflected in the stabilization of the growth of urban populations, which will not necessarily prevent their continued economic growth. It is here that the lack of reliable indicators of economic growth in cities makes itself felt. Nevertheless, recent and present-day dynamics of city systems are characterized by a historical tendency to vigorous growth. How does this growth spread over the different cities?

6.4.2 Hierarchical Differentiation and Distributed Urban Growth

The size of a city is the product of a long-term process of local accumulation. There were many fluctuations in city sizes over historical times, but once systems of cities are established in a given territory, the rapid upsurge of new cities as well as the sudden collapse of some of them become very rare, and less and less probable. As already mentioned, there are very large differences in city sizes, when measured by the total population they concentrate, their surface areas, or their economic gross product. The number of cities is in inverse geometrical proportion to the number of their inhabitants, as summarized by Zipf in his famous “rank size rule” (1941, 1949). This hierarchical differentiation within systems of cities is an emergent property that characterizes the organization of consolidated and integrated urban systems (Pumain, 2006). In Fig. 6.6, the remarkable stability of this hierarchical structure in the long term can be seen for each of the four city systems under observation, whatever the historical “style” to which they are assigned.

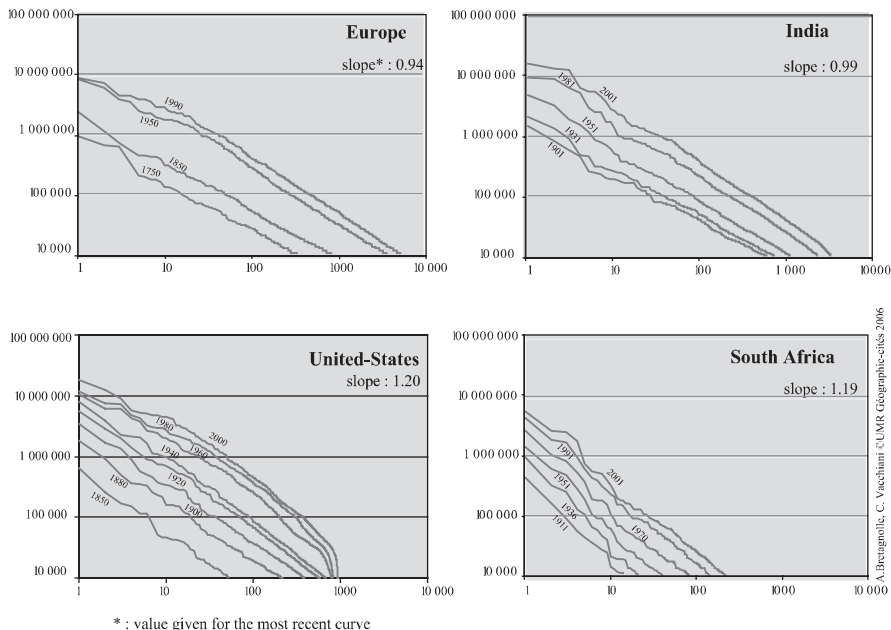
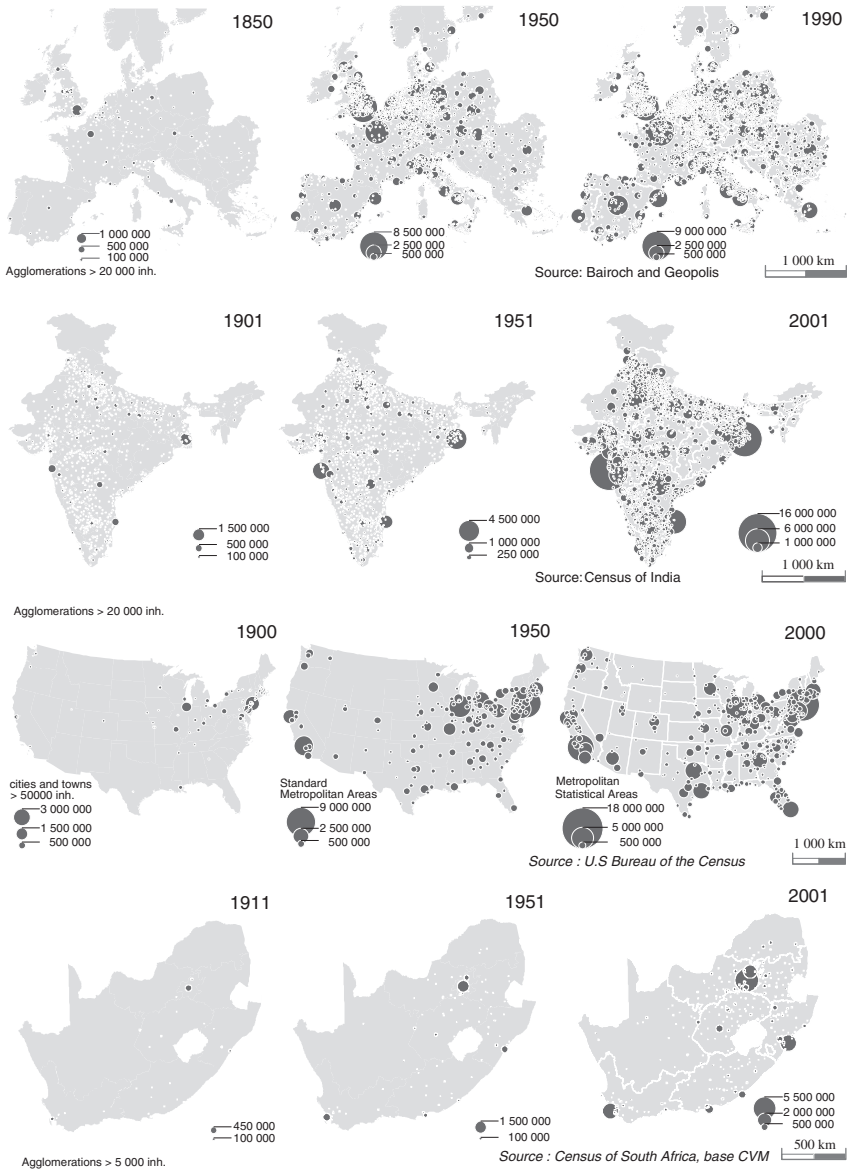


Fig. 6.6 Pareto-Zipf distribution of city size in the long-term period for the four categories of city systems

Not only does the shape of the distribution of city sizes remain very similar over decades, the hierarchical and spatial order given by the population sizes also evolves very little. The maps in Fig. 6.7 show this evolution (the surface area of the circles increases proportionately to the population of the cities between 1850 or 1900 and 2000 for the four cases under study). The similarities in evolution as represented in this manner are striking: on average, cities grow in a manner that is proportional to their size. This explains how initial differences tend to persist over very long periods. The reference model is that developed by Gibrat, explaining the non-symmetrical distribution of city sizes (a lognormal distribution that differs slightly from Zipf’s law) by way of a random growth process in which all cities, at each time interval, possess the same probability for growth rate. This model, which has been termed “distributed growth” (Pumain, 1997) because all parts of the system grow at more or less the same rate, gives a good account of growth processes in established city systems, although with certain systematic divergences that will be examined below. This statistical model merely informs us that the causes of variation in city populations are so numerous and diverse that it is sufficient to use a straightforward random process to represent this variation, and, thereby, to explain why all systems have the same general distribution pattern for city sizes. But to go further in interpretation, and in particular to understand phenomena of divergence from this model, we also need to understand what properties of city systems explain this distributed growth.



A. Bretagnolle@UMR Géographie-cités - 2006
S. Olivreau@UMR Géographie-cités - 2006
C. Vacchiant@UMR Géographie-cités - 2006

Fig. 6.7 Evolution of urban patterns in integrated city systems

Going back at least to the enlightened work by Botero in the 16th century (Pumain & Gaudin, 2002), there has been an awareness that cities are constantly *in competition* to capture resources and innovation so they can continue to make good use of what they have already acquired and maintain or increase their influence within the city systems with which they entertain relationships. This competition explains why innovation spreads quickly from one city to another and why the resulting qualitative and quantitative changes are of more or less the same magnitude, over short periods, in integrated systems. The consequence of this mode of growth is that city size hierarchy and the inequalities that arise following various types of political accidents (war, choice of capital city) or economic accidents (functional specializations, as in the industrial revolution), are maintained over periods that last a lot longer than the events that caused them.

6.5 Reinforcement of Urban Hierarchy

In the course of time (at least since the first industrial revolution), inequalities in city sizes have shown a tendency to grow. This growing inequality can be explained both by the arrival of new smaller towns in the system and by the fast growth of the largest cities. A first divergence from Gibrat's model that is often noted is that the variations in growth are not totally independent of city size: over long periods of time, the large cities grow a little faster and the smaller towns and cities a little more slowly than the average, and inequalities in size reinforce more markedly than what is expected by the model (Bretagnolle, 1999).

6.5.1 Concentration of Urban Population in the Integrated Systems

Whatever the system considered, the degree of hierarchization increases over time in integrated systems (see Table 6.4; an exception is noted for South Africa, which can be explained by the late period when the system became established).

The observation of adjustment slopes in the rank-size graphs for the last three decades could suggest that this process of historical concentration is terminating

Table 6.4 Increase in inequalities among city sizes

	1900	1950	2000
Europe	0.74 (1850)	0.91	0.94
India	0.76	0.86	0.99
United States	0.97	1.15	1.20
South Africa	1.39	1.16*	1.19

*This slope is for 1960, because from this date on the entries of small towns and creations became less frequent.

The degree of inequality is measured by the absolute value of the adjustment slope in rank-size graphs.

in countries with long-standing industrialization. However, observations once again show that the choice of delimitation criteria for urban areas has a strong influence on the results of measures of growth and concentration. Using France as an example, if evolving delimitations are adopted that take the increase in transportation speeds into account, it can be seen that concentration, far from falling off, actually increases (Bretagnolle, Paulus, & Pumain, 2002). In the USA, the results are less clear-cut: whatever the delimitations used (official SMA/MSA definitions or the harmonized database that we have developed), concentrations do indeed decrease between 1970 and 1980, and fluctuate after that time.

6.5.2 Selection Process and Hierarchical Diffusion of Innovations

The reinforcement of urban hierarchy can be explained by the process of hierarchical and selective diffusion of innovation, which combines two effects: (1) a growth advantage for the largest cities in a system by way of early adoption of innovations and (2) a tendency of relative decline for smaller towns and cities short circuited by these innovations. This process will be examined in more detail in Chapter 8. The capture of innovations by the large cities is explained by the complexity and diversity of their functions and infrastructures (the result of a long history of successive adaptations), which provide better access to information and a greater capacity to carry the high cost and risk associated with innovation. Over the long term, the most marked growth trends that arise from the initial advantage benefit the larger cities and result in “top-down” hierarchization within city systems.

Conversely, smaller towns and cities gain access to innovation at a later date, or not at all, which results in a “bottom up” hierarchization effect. The diffusion of innovations in transportation illustrates this selection process. Whether we consider the pre-industrial era with its mail-coach network or the later railway network, motorways, or airline networks, the large cities and the smaller towns possessing an attractive specialization first profited from such services. From as early as 1820, the quest for increased speed, which is not solely a modern concern, led to a reduction in the number of intermediate stops throughout the 19th and 20th centuries. Figure 6.8 shows how, over and above specific savings brought about by any given technical innovation (for instance the braking power of locomotives or the electrification of the railways), there is a constant and regular increase in the average speed from Paris to the main cities in the French urban system.

For the towns or cities that lost their “nodal” position of access to the fastest networks of their time, effects are visible in the long term, even if demographic variables alone are taken into consideration (see Bretagnolle, 2003). Thus, in disagreement with certain widespread assumptions about the growing universal availability of information, in history the reverse has occurred. Innovation in transport and communication networks does indeed spread among towns and cities overall but with time lapses in relation to city size, which leads to greater concentration of the channels for the circulation of information that is essential at a given time.

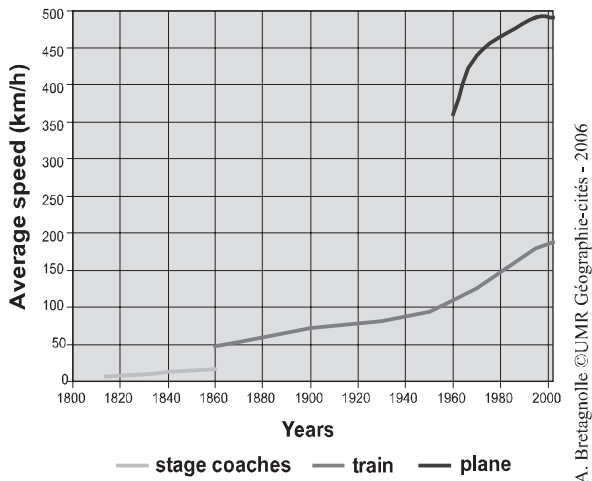
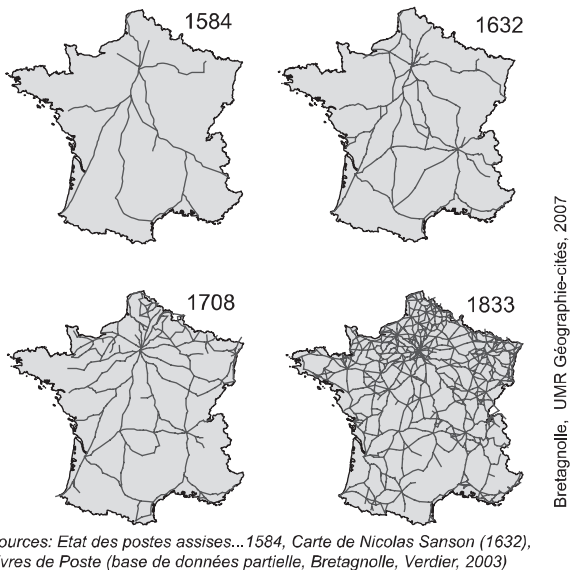


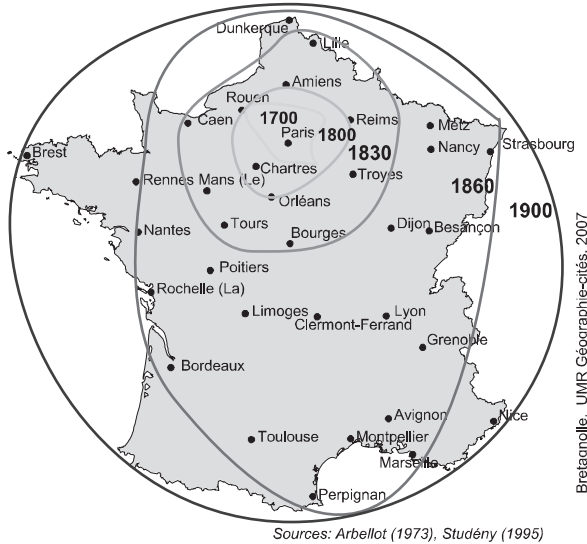
Fig. 6.8 Increasing speed of transportation, between Paris and main French metropolises (1800–2000)

As an example of the selective diffusion process, we have mapped the successive extensions of the postal road network in France from 1584 to 1833 (Bretagnolle & Verdier, 2005, 2007). As an innovation, these postal roads connected first the cities that were already in the upper part of the urban hierarchy, following a dual process of hierarchical diffusion and space-filling (Fig. 6.9).



Sources: *Etat des postes assises...1584, Carte de Nicolas Sanson (1632), Livres de Poste (base de données partielle, Bretagnolle, Verdier, 2003)*

Fig. 6.9 The evolution of the postal roads between 1584 and 1833 (current boundaries)



Bretagnolle, UMR Géographie-cités, 2007

Sources: Arbellot (1973), Studény (1995)

Fig. 6.10 Isochronic map: cities that are located at a one day distance from Paris (1700–1900)

As an example of the acceleration in the evolution of social interaction space from postal network to railway network, we have mapped the French cities that are located at a one-day distance from Paris between 1700 and 1900. These were the places that had the benefit of rapid interaction and that were not accessible to non-connected cities (Fig. 6.10). The consequence of this “space-time contraction” (Janelle, 1969) is that the smaller towns and cities grow on average at lower rates than the largest (Table 6.5).

Table 6.5 Average annual growth rates (%) of cities according to their size

Size class	Europe (since 1850)	India (since 1901)	USA* (since 1940)	South Africa (since 1951)
>100 000	1.38	2.45	1.79	3.21
50–100 000	1.13	1.95	0.92	3.11
< 50 000	0.99	1.71	*	3.10

*Values are computed for SMAs, i.e. since 1940 and with a minimal size of 50,000 inh.

6.6 Conclusion: Toward an Evolutionary Theory of Urban Systems

Through observation of progress through history, and comparing different regions in the world, we have identified several patterns of city system dynamics. We have underlined the importance of identifying urban and territorial entities that are geographically relevant to comparisons in time and space. We went on to distinguish phases in the establishment of these systems, or in their de-structuring

and re-structuring, and phases of “normal” evolution in consolidated systems, when interactions among the towns and cities are controlled and regulated in a more homogeneous manner, for instance inside national frontiers. We have shown that the structural properties of these integrated urban systems are partly similar, insofar as they result from these dynamics, and partly dependent on the history of the political territories to which they belong. The main difference observed between the “older” city systems and those in “new” countries can thus be explained by the evolution of the material conditions in which individuals, goods and information were circulating at the time when they became established. The second difference, distinguishing more “monogenetic” systems from “dual” systems, can be explained by the interference between dynamics belonging inside the system and the dynamics of relationships outside their territory. Thus, there is indeed a correlation between the structure of city systems and the terms of exchanges between cities, over time.

These exchanges and interactions among cities, on all scales, are of the center-periphery type, which generates asymmetry, and, which in turn, enables the accumulation of population and activity in certain places and some redistribution via diffusion. Exchanges among cities are multiform, and their range is highly varied; they produce complex networks where patterns of hierarchical structuring can be detected, but it is very difficult to observe such interaction in a direct manner. Just as we had to infer the dynamics of city size from the observation of the evolution of population sizes, in Chapter 8 we will attempt to deduce the logic of the ability of cities to generate and adopt innovation from observations of changes in their economic activities and in their social composition. Understanding these processes is essential to analyze how urbanization actually contributes to innovation and to globalization, so as to go on to identify levers for action.

References

- Bairoch, P., Batou, J., & Chèvre, P. (1988). *La Population des Villes Européennes de 800 à 1950, Banque de Données et Analyse Sommaire des Résultats* (Vol. 2). Publications du Centre d'Histoire Economique Internationale de l'Université de Genève. Geneva, Switzerland: Librairie Droz.
- Berry, B. J. L. (1964). Cities as systems within systems of cities. *Papers of the Regional Science Association*, 13, 147–163.
- Berry, B. J. L. (1967). *Geography of market centers and retail distribution*. Englewood Cliffs, NJ: Prentice Hall.
- Braudel, F. (1967). *Civilisation matérielle et Capitalisme*. Paris, France: Colin.
- Bretagnolle, A. (1999). *Les systèmes de villes dans l'espace-temps: effets de l'accroissement de la vitesse des déplacements sur la taille et l'espacement des villes*. Thèse de doctorat, Paris, France: Université Paris I.
- Bretagnolle, A. (2003). Vitesse et processus de sélection hiérarchique dans le système des villes françaises. In D. Pumain, M.-F. Mattéi (Eds.), *Données Urbaines* (Vol. 4). Paris, France: Anthropos.
- Bretagnolle, A., & Verdier, N. (2005). Images d'un réseau en évolution: les routes de poste dans la France pré-industrielle (XVII^{ème} - début XIX^{ème}s.). In *Mappemonde*, No. 79 (2005-3) <http://mappemonde.mgm.fr>.

- Bretagnolle, A., & Verdier, N. (2007). Expanding the network of postal routes in France (1708–1833). In Le Roux (Ed.), *Post offices of Europe 18th–21th century. A comparative history* (pp. 155–173). Paris, France: Comité pour l'Histoire de la Poste. <http://halshs.archives-ouvertes.fr/halshs-00144669/en/>. Accessed 5 November 2007.
- Bretagnolle, A., Paulus, F., & Pumain, D. (2002). Time and space scales for measuring urban growth. *Cybergeo*, 219.
- Bretagnolle, A., Pumain, D., & Vacchiani-Marcuzzo, C. (2007). Les formes des systèmes de villes dans le monde. In D. Pumain, & M.-F. Mattéi (Eds.), *Données Urbaines*, 5, Paris, France: Anthropos.
- Bretagnolle, A., Giraud, T., & Mathian, H. (2007). L'urbanisation des Etats-Unis, des premiers comptoirs coloniaux aux Metropolitan Standard Areas, *Cybergeo*.
- Christaller, W. (1933). *Die Zentralen Orte in Süddeutschland*. Jena, Germany: Fisher.
- Crozet, Y., & Joly, I. (2004). Budget-temps de transport: les sociétés tertiaires confrontées à la gestion paradoxale du bien le plus rare. *Cahiers Scientifiques du Transport*, 45, 27–48.
- Davies, R. J. (1972). *The urban geography of South Africa*. Institute for Social Research, Durban, South Africa: University of Natal.
- de Vries, J. (1984). *European urbanisation: 1500–1800*, London, UK: Methuen.
- Gibrat, R. (1931). *Les Inégalités Economiques*. Paris, Sirey.
- Hall, P., Gracey, P., & Drewry, R. (1973). *The containment of urban England*. London, UK: Sage Publications.
- Hollar, W. (1667). *Hollar's Exact Survey of the City of London*. Reproduction, Edward Stanford Ltd., N.D. 83 × 54 cm. De la Feuille James (1690), Londini angliae regni metropolis novissima and accuratissima, <http://www.british-history.ac.uk>.
- Janelle, D.G. (1969). Spatial reorganisation: a model and concept. *Annals of the Association of American Geographers*, 59, 348–364.
- Pinol, J. L. (2003). *Histoire de l'Europe Urbaine. 1 – De l'Antiquité au XVIII^e siècle*. Paris, France: Le Seuil.
- Pirenne, H. (1927). *Les Villes au Moyen Age: Essai d'Histoire Economique et Sociale*. Brussels, Belgium: Lamertin.
- Pred, A. (1977). *City systems in advanced societies*. London, UK: Hutchison.
- Pumain, D. (1997). Vers une théorie évolutive des villes. *L'Espace Géographique*, 2, 119–134.
- Pumain, D. (Ed.). (2006). *Hierarchy in natural and social sciences*. Methodos Series Vol. 3, Dordrecht, The Netherlands: Springer.
- Pumain, D., & Gaudin, J.-P. (2002). Systèmes de villes et pouvoir. L'analyse de Giovanni Botero à l'époque de la Renaissance. *Cybergeo*, 227.
- Pumain, D., Saint-Julien, T., Cattan, N., & Rozenblat, C. (1991). *The statistical concept of city in Europe*. Eurostat.
- Reynaud, J. (1841) Villes. *Encyclopédie Nouvelle* (Vol. 8, pp. 670–687). Paris, France: Gosselin.
- Vacchiani-Marcuzzo, C. (2005). *Mondialisation et système de villes: les entreprises étrangères et l'évolution des agglomérations sud-africaines*. Thèse de Doctorat, Université Paris 1, tel.archives-ouvertes.fr/tel-00011351/en/.
- Zelinski, W. (1971). The hypothesis of the mobility transition. *Geographical Review*, 61, 219–249.
- Zipf, G. K. (1941). *National unity and disunity*. Bloomington, IN: Principia Press.
- Zipf, G. K. (1949). *Human behaviour and the principle of least effort*. Cambridge, MA: Addison-Wesley Press.

Chapter 7

The Self Similarity of Human Social Organization and Dynamics in Cities

Luis M.A. Bettencourt, José Lobo and Geoffrey B. West

7.1 Introduction

The issue of whether quantitative and predictive theories, so successful in the natural sciences, can be constructed to describe human social organization has been a theme of inquiry throughout the entire history of science. Aristotle was perhaps the first to write on the subject in ways that present clear scientific challenges that are still with us today. In *Politics* (Book I), he wrote

[. . .] it is evident that the state [polis] is a creation of nature, and that man is by nature a political animal. The proof that the state is a creation of nature and prior to the individual is that the individual, when isolated, is not self-sufficing; and therefore he is like a part in relation to the whole.

The idea that cities (the “State” for Aristotle) are natural inevitable structures on which humans coalesce and thrive has suggested many metaphors for cities as natural organisms (Miller, 1978; Girardet, 1992; Graedel & Allenby, 1995; Botkin & Beveridge, 1997; Decker Elliott, Smith, Blake, & Rowland, 2000), or ecologies (Macionis & Parrillo, 1998). Modern sociological thought about the nature of urban life (Durkheim, 1964; Simmel, 1964; Macionis & Parrillo, 1998), especially in the United States (Wirth, 1938), was born largely out these analogies. Cities as consumers of resources and energy, and producers of organizational structures and waste have a clear counterpart in biological organisms. Therefore, it is interesting to ask to what extent these analogies are more than anecdotal. For instance, are analogies of cities as organisms, with specific metabolism, useful to establish quantitative expectations for their resource demands, environmental impacts and growth trajectories?

The idea that cities are emerging natural structures, that in some sense (to be demonstrated below) are independent of culture, geography or time, also suggests that there should be universal features common to all urban agglomerations (Wirth, 1938), which, in turn, may define an *average idealized city*. Such a city would be characterized in terms of quantitative indicators and would constitute the benchmark

L.M.A. Bettencourt (✉)

Theoretical Division, Los Alamos National Laboratory and Santa Fe Institute, Santa Fe, NM, USA

against which real cities should be measured, both in their successes and in their shortcomings. This concept of a quantifiable average city may be natural, even intuitive, to anyone who seeks synthesis among the huge diversity of urban social life and is implicit in many policy considerations (Wirth, 1938; Macionis & Parrillo, 1998). But its pursuit is often at odds with traditions in the social sciences that emphasize instead the richness and differentiation (Macionis & Parrillo, 1998; Durkheim, 1964; Simmel, 1964) of human social expression. The tension between these two approaches can only be diffused, in our opinion, by empirical investigations determining if average idealized characterizations of urban organization are supported by data, and can be synthesized as predictive theories of certain key features of human dynamics and organization.

Below, we will pursue and partially accomplish some of these goals. We will show that, when observed from the point of view of their rates of change, cities are approximately self-similar entities across entire urban systems (usually taken to be nations), and that these properties scale with population size in a manner that is independent of any particular reference scale. In this sense, knowledge of urban indicators for a city of a given size implies predictions for those of another, given only their population ratio. Because quantitatively similar scaling laws are a property of biological organisms (West, Brown, & Enquist, 1997; West, Brown, & Enquist, 1999; West, Brown, & Enquist, 2001; Enquist, Brown, & West, 1998), we will also be able to establish to what extent cities can indeed be understood in terms of biological organization, and specifically how human societies differ and transcend these structures.

The remaining of this chapter is organized as follows. We start with basic expectations for the behavior of cities, resulting from analogies to biological scaling. We then give a brief account of the initial quantitative studies that suggested that the framework in urban organization would be more complex. Following on these hints we give an overview of other urban indicators and how they scale with population. We then discuss the implications for growth from urban scaling relations and conclude with some speculations and directions for future work.

7.2 General Expectations from Biological Scaling

Before we started our empirical investigation on urban scaling we attempted to translate the metaphor for cities as biological organisms, in terms of quantitative relations for human organizations.

The most fundamental quantity characterizing any physical system is its energy, which in turn sets dynamical time scales. For many complex systems, particularly those in biology and society, energy consumption is the key to identify leading rhythms of internal organization, growth and information creation. For a biological organism energy consumption per unit time is a measure of its metabolism. Remarkably metabolic rates Y , scale with body size M (mass) according to a simple power relation:

$$Y = Y_0 M^\beta, \quad \beta = 3/4, \quad (7.1)$$

which holds across 21 orders of magnitude in body mass, and all different species (West et al., 1997, 1999). The exponent $\beta = 3/4$ can be understood in terms of the networks of resource distribution inside organisms (e.g., the vascular system of animals and plants). These networks are hierarchical; distributing resources from central points (e.g., heart) to every component of the system (e.g., cells). In this way, they carry a conserved fluid to every part of a volume of tissue that constitutes the organism. These networks have been optimized by natural selection to be efficient, in the sense of dissipating the least amount of energy possible. Under these conditions, the networks can be abstracted as hierarchical branching processes with a non-trivial fractal dimension. In d dimensions, it can be shown that the exponent $\beta = d/(d + 1)$, which becomes $3/4$ for $d = 3$. Because the scaling law (7.1) has the dimension of a rate, it also predicts how characteristic times, characterizing the organism's behavior, scale with size. Specifically, characteristic times (e.g., lifespan) per unit mass scale as $M^{1/(d+1)}$, while rates (such as heart or respiratory rhythms) scale inversely to times, as $M^{-1/(d+1)}$ (West et al., 1997, 1999).

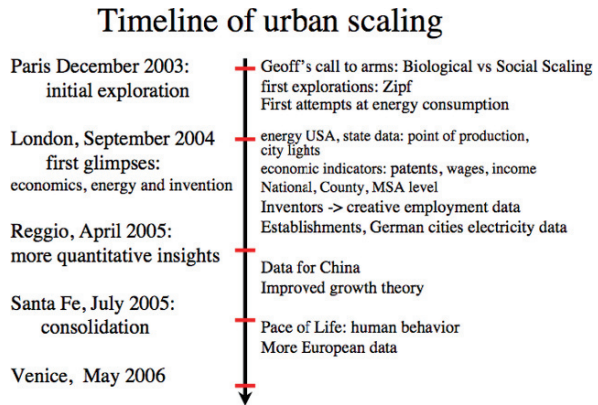
At least at the superficial qualitative level, cities can, likewise, be thought of in terms of idealized networks of distribution that supply people, households and institutions with water, power, etc., and remove unwanted byproducts. It is less clear, however, what quantity may play the role of scale. The most natural is population, but other units, such as households or firms, are conceivable. Below, we show that adopting population as a measure of size of a city does indeed produce clear scaling relations for many urban quantities.

Another issue that arises when translating expectations from biological scaling, is the dimensionality of the system. While the natural dimension of a city may be $d = 2$, dense cities can also show growth in height producing structures that may have $2 < d < 3$. One last point concerns the definition of city itself (see Chapter 6). Many studies in urban geography have struggled with creating good definitions of the spatial limits of a city. This is difficult to do in terms of population density or built up area, since these quantities vary continuously from the city core to the periphery (increasingly in the US, city cores show decreases in population density, which has been moving to lower suburban locations). Because we are primarily interested here in human social activity, we adopted a definition of city that, as much as possible, reflects its economic character as an integrated labor market, comprised of a city core and all surrounding areas where substantial fractions of the population work within the city limits. These are Metropolitan Statistical Areas (MSA) in the USA, Larger Urban Zones (LUZ) in the European Union and Urban Administrative Units (UAU) in China.

7.2.1 A Timeline of Urban Scaling Results in ISCOM

As we briefly discussed above, our initial investigation and pursuit of data was guided by the presentation by one of us (West) at the ISCOM meeting in Paris 2003, drawing the analogy between biological organisms and human social organizations as a working hypothesis towards building a scaling theory of cities.

Fig. 7.1 Time line of our investigation in urban scaling. The initial motivation from Biology grew into a more complex and detailed picture as data came in. In retrospect, it is easy to recognize consecutive ISCOM meetings as major milestones in our progress



In the period between the ISCOM Paris meeting in December 2003 and the subsequent working group meeting in London in September 2004, our efforts concentrated on finding datasets to empirically test the expectations from Biological scaling for cities. A time line of our investigation in urban scaling is shown in Fig. 7.1. Energy consumption at the metropolitan level turned out to be difficult to measure in the USA, because most data are proprietary and fragmented. Information about points of production and about the networks of distribution is available but requires a large amount of reconstruction (and extrapolation) work to be mapped into city consumption patterns. We had more success with Germany (through our collaborators Dirk Helbing and Christian Kuehnert), where electrical production is tied in with individual cities for historical reasons. We also explored definitions of city, as it was unclear how best to aggregate socioeconomic data. We investigated scaling for counties, cities, and metropolitan areas. We found that, although indications of scaling existed at different aggregation levels, MSAs provided the most persuasive and consistent statistical signatures, and we adopted these units for subsequent studies in the USA and abroad. The rationale of the definition of MSAs, as the set definition of city that is as much as possible devoid of arbitrary administrative units and is instead an integrated economic and social unit, makes the most sense with our findings shown below. This does not preclude, of course, that even better definitions of city limits, an important problem in urban geography, can be defined. It is, in fact, plausible that greater understanding of urban functionalities, captured in a set of scaling relations to be discussed here, will aid guide such constructions.

7.2.2 Energy Consumption and Invention Rates vs. Urban Population Size

Our first results involved data on general socioeconomic activity, such as wages (Fig. 7.2), as well as on electrical energy consumption in German cities, see Table 7.1. Both data sets pointed immediately to fundamental differences in Biology, in that patterns of energy consumption and wealth generation did not show economies of scale ($\beta < 1$). In fact, if anything, these results indicated that scaling

Fig. 7.2 Superlinear ($\beta = 1.12$) scaling of wages with metropolitan population for the USA in 2000

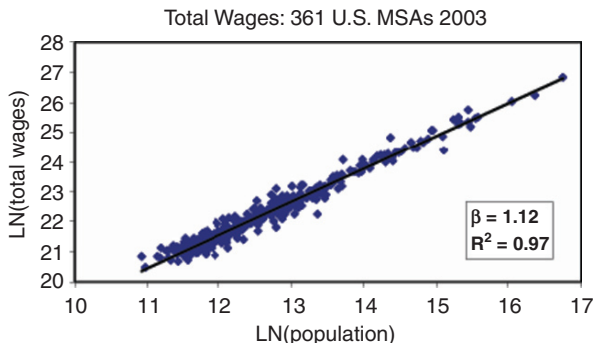


Table 7.1 Scaling exponents for electrical energy quantities for German cities in 2002

Variable	exponent \pm standard deviation
Usable energy	1.09 \pm 0.03
Household supply	1.00 \pm 0.03
Length of cables	0.88 \pm 0.03
Resistive losses	1.10 \pm 0.03

is superlinear ($\beta > 1$), but that the effects were small, with $\beta \sim 1.1-1.2$, but statistically significant. We also found at this point – from Helbing et al. (Chapter 16 of this volume) – that certain amenities, such as numbers of restaurants, scaled with clear superlinear exponents, whereas others such as hospital beds, were approximately linear.

Although showing a mixed picture, these results pointed to several features that were confirmed by subsequent data. First, urban indicators, from infrastructure to socioeconomic characteristics, show clear and manifest scaling with city population size, characterized by exponents that deviate from unity by relatively small but statistically significant margins. Secondly, exponents for socio-economic quantities scale with $\beta > 1$, individual needs with $\beta \sim 1$, and material infrastructure (such as the length of electrical cables) scale with $\beta < 1$, see (Bettencourt, Lobo, Herbing, Kuehnert, & West, 2007).

7.2.3 Patenting Rates and Creative Employment

The next set of data we analyzed dealt with measures of innovation, as measured via patenting rates (Bettencourt et al., 2007). When organized in terms of inventor’s residential address, new patents filed in the US can be tallied up by metropolitan statistical area. It has been recognized for quite some time that cities are the primary seats of innovation. Patenting, as an admittedly limited proxy to general innovative processes, has been studied in terms of its geographic preferential location for quite some time. As a result, patenting has been identified in the US and in other countries as a primary metropolitan phenomenon. Despite this rich evidence, the systematic study of the rate of patenting with metropolitan population had not been undertaken (Bettencourt et al., 2007).

The results of plotting new patents per MSA vs. MSA population size are shown in Fig. 7.3 for 1980 and 2000. Although there is some scatter in the data (which are shown without any averaging), a clear scaling trend is present in both years, with an exponent that is statistically consistent across two decades. This statistical invariance of the scaling exponent is particularly impressive when we note that, in those two decades, patenting subjects shifted dramatically away from traditional industries to new activities in electronics, software and biotechnology.

Patenting data also offered the possibility of testing several scenarios to explain the observed superlinear scaling of innovation rates with population. Two alternative scenarios are natural and testable given available data: observed increasing returns to scale (superlinear scaling) could (1) be the result of increased individual productivity, following from greater interactions with a larger number of inventors, proportional to city size; or, alternatively, follow from (2) individual productivity that is independent of city size, but is compensated by a greater number of inventors that are disproportionately represented the larger the city.

We further hypothesized that the first scenario, increased productivity due to greater interactions, should display a signature in patent co-authorship, since contact between a number of inventors naturally scale superlinearly (with an exponent $\beta = 2$, if all authors connect to each other). Thus, scenario (1) would predict a num-

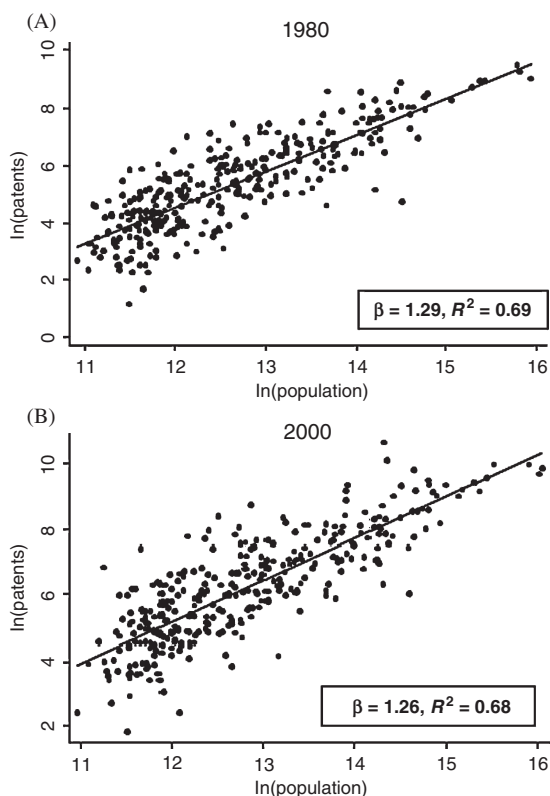


Fig. 7.3 Number of new patents per year in 1980 (A) and 2000 (B) vs. metropolitan population size. The solid line shows the result of a power law fit to the data with exponent β , shown as inset. Remarkably, despite enormous changes in technology, scaling laws stay statistically equivalent across the two decades, see (Bettencourt et al., 2007)

ber of inventors proportional to metropolitan population, but a superlinear scaling of inventor connectivity, which would yield the overall observed gains in productivity.

Conversely, scenario (2) predicted simply that productivity per inventor (average number of patents per author) would stay constant across city size, but that inventors would disproportionately be located in larger cities, thus accounting for greater rates of patenting. Table 7.2 and Fig. 7.4 show how empirical evidence settles the case in favor of scenario (2).

Table 7.2 Scaling of inventor connectivity, measured via patent co-authorship, and of number of inventors with metropolitan population size in the USA. These results, taken together with the evidence of Fig. 7.4, indicate that superlinear scaling in invention is primarily the result of the presence of a disproportionate number of inventors in larger cities and not due to superlinear increases in individual inventor productivity

Variable vs. # of	Scaling exponent
Connectivity	$\beta = 0.823 \pm 0.001$
Inventors	$\beta = 0.981 \pm 0.002$

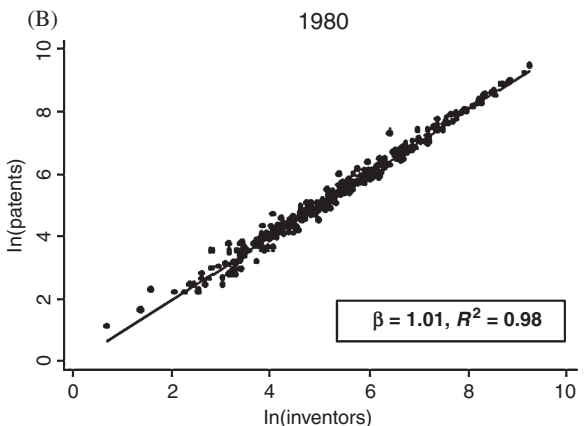
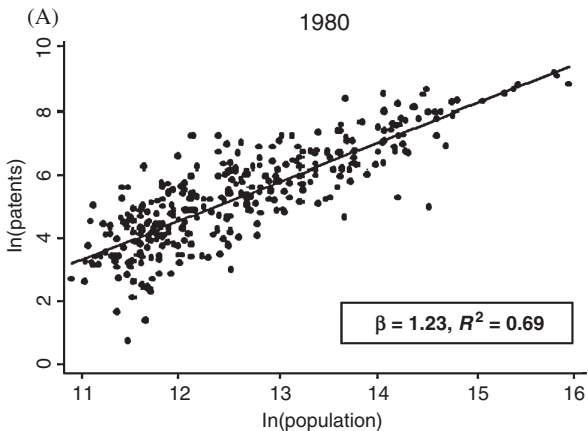
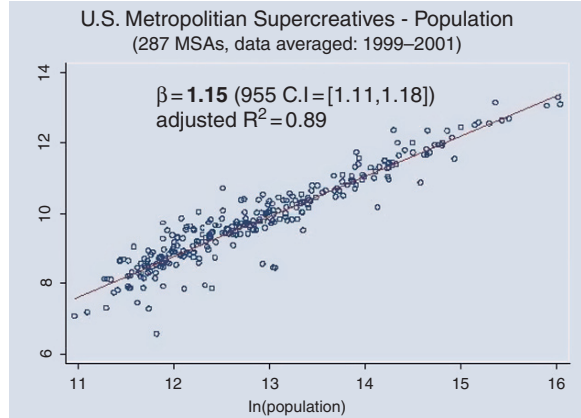


Fig. 7.4 Scaling of number of inventors with metropolitan population (A) and of patents with number of inventors (B). Taken together with the results of Table 7.2, these data indicate that inventors are disproportionately represented in the larger cities, but that individuals do not become more productive in larger populations

Fig. 7.5 Scaling of supercreative professionals with metropolitan city size



Finally, we asked whether the phenomenon observed here for inventors was in fact much more general, indicating a disproportionate number of inventive and creative activities the larger the city. Indeed, this expectation was confirmed by analyzing numbers of professionals in specific activities vs. city size, both in France and in the US (Chapter 8, this volume). As a summary, we show the scaling of numbers of “super-creative” professionals (Florida, 2004) with metropolitan population size (Bettencourt et al., 2007) in Fig. 7.5, indicating that scientific, technical, artistic, media and management activities scale superlinearly with city size, with an exponent $\beta = 1.15$. These results indicate that larger cities are not functionally scaled-up versions of smaller towns but, rather, are different in their relative activity breakdown, with more people disproportionately occupied in innovation and invention.

7.3 Urban Scaling and the Interplay Between Social Processes and Infrastructure

The set of results discussed above set the stage for a taxonomy of quantities that characterize cities as self-similar structures. Cities, in fact, realize both certain economies of scale in resource networks, typical of biology, and enable superlinear processes that are unique to human social organization.

7.4 Economies of Scale and Material Infrastructure

As we have already seen for electricity, certain aspects of infrastructure benefit from higher population density to realize economies of scale (see Table 7.3). Note that most of these quantities (length of cables, road surface) refer to material infrastructure networks, but not to resource consumption rates, which may scale superlinearly. Note also that the expectation for $1/3$ power laws for infrastructure networks in a

Table 7.3 Scaling of infrastructural quantities with city size realizes economies of scale, analogous to those in Biology. These economies of scale appear as sublinear scaling ($b < 1$) laws

Y	β	95% CI	adj.- R ²	Observations	Country/year
Gasoline Stations	0.77	[0.74,0.81]	0.93	318	USA/2001
Gasoline Sales	0.79	[0.73,0.80]	0.94	318	USA/2002
Length of electrical cables	0.88	[0.82,0.94]	0.82	387	Germany/2001
Road surface	0.83	[0.74,0.92]	0.87	29	Germany/2001

2-dimensional city is not borne out by data. This may be a consequence of several factors, including gradients in population density, the not purely two dimensional character of cities, and the fact that resource delivery requirements may drive these networks to inefficiency.

This last point is important as it raises the question of cause and effect, namely whether infrastructure is the driver of human social behavior or if the converse is true. Although this question cannot be settled satisfactorily with the present evidence, the case of electrical consumption in German cities may be paradigmatic. Although economies of scale are certainly realized in cabling, total consumption scales superlinearly. This can only be achieved at the cost of rising inefficiency, which is manifested as a superlinear scaling in resistive losses (see Table 7.1). Thus, at least in this case, it is suggestive that human social needs drive infrastructure, rather than the other way around, as happens in biological organisms.

7.4.1 Individual Needs

Another interesting instance of urban scaling is that certain quantities are neither directly related to social behavior or to material infrastructure, but simply reflect individual needs that, once satisfied, cannot be easily expanded. For example, typically each person needs one job, one dwelling, and a typical amount of water and electricity at home. These quantities scale linearly with city size as shown in Table 7.4.

Note that although electrical consumption increased with city size, household consumption increases only linearly. Thus, it is the energy used to enable social

Table 7.4 Individual needs, such as household utility consumption, numbers of jobs, and dwellings, scale linearly with metropolitan population

Y	β	95% CI	adj. R ²	observations	Country/year
Total establishments	0.98	[0.95,1.02]	0.95	331	USA/2001
Total employment	1.01	[0.99,1.02]	0.98	331	USA/2001
Total household electrical Consumption	1.00	[0.94,1.06]	0.70	387	Germany/2001
Total Household electrical Consumption	1.05	[0.89,1.22]	0.91	295	China/2002
Total Household water Consumption	1.01	[0.89,1.11]	0.96	295	China/2002

productive activity – devoted to work rather than maintenance – ranging in scope from industry, to culture and learning, and street lighting that accounts for the superlinear character of the total consumption.

7.4.2 The Urban Economic Miracle

One of the most important characteristics of cities is that they are the primary centers for wealth creation in every human society. Although urban economists have established a positive relationship between urban size and productivity (Sveikauskas, 1975; Segal, 1976; Henderson, 2003) the identification of these statistical regularities in terms of scaling laws is new and extremely important for the understanding of the self-similar social processes that enable prosperity and economic growth. Measures of wealth creation or productivity follow exquisite superlinear scaling relation, across time and for different nations, with adjusted R^2 very close to unity, see Table 7.5, indicating nearly perfect fits.

Table 7.5 Wealth creation and productivity follow exquisite scaling laws with metropolitan population, with exponents $b \sim 1.10\text{--}1.15$, and adjusted R^2 close to unity. Data aggregated at the level of the European Union encompasses several loosely connected urban systems and gives a poorer fit

Y	β	95% CI	adj.- R^2	observations	Country/year
Total Wages/yr	1.12	[1.09,1.13]	0.96	361	USA/2002
GDP/yr	1.15	[1.06,1.23]	0.96	295	China/2002
GDP/yr	1.13	[1.03,1.23]	0.94	37	Germany/2003
GDP/yr	1.26	[1.03,1.46]	0.64	196	EU/2003

Clearly, economic growth is strongly correlated to innovation and fast adaptation to new opportunities (Romer, 1986, 1990; Lucas, 1988, Glaeser, Kolko, & Saiz, 2001). Table 7.6 shows a summary of measures of innovation and employment in creative activities, which all show strong superlinear scaling.

Table 7.6 Superlinear scaling exponents for innovation and employment in innovative sectors

Y	β	95% CI	adj. R^2	observations	Country/year
New Patents/yr	1.27	[1.25,1.29]	0.72	331	USA/2001
Inventors/yr	1.25	[1.22,1.27]	0.76	331	USA/2001
Private R & D employment	1.34	[1.29,1.39]	0.92	266	USA/2002
“Supercreative” professionals	1.15	[1.11,1.18]	0.89	287	USA/2003
R & D employment	1.67	[1.54,1.80]	0.64	354	France/1999*
R & D employment	1.26	[1.18,1.43]	0.93	295	China/2002

7.4.3 The Darker Side of Cities: Costs, Crime and Disease

If the possibility of a larger and richer set of human contacts, made possible in a larger city, enables the creation of ideas and wealth, then they may also encourage

Table 7.7 Costs, the incidence of certain transmissible diseases, crime and other patterns of human behavior, such as walking speed (Bornstein & Bornstein, 1976), are also superlinear scaling laws with city size

Y	β	95% CI	Adj. R^2	Observations	Country/year
Cost of housing (per capita)	0.09	[0.07,1.27]	0.21	240	USA/2003
New AIDS cases	1.23	[1.18,1.29]	0.76	93	USA/2002
Violent crime	1.16	[1.11,1.18]	0.89	287	USA/2003
Walking Speed (per capita)	0.09	[0.07,0.11]	0.79	21	Several/1979

other types of less benign social activities (Milgram, 1970) such as those involved in crime (Glaeser & Sacerdote, 1999) and disease transmission. Thus, we may expect, at least in the absence of strong intervention, that crime and disease incidence (both temporal rates, analogous to idea or wealth creation) also scale superlinearly with city size. These expectations are well borne out by data as shown in Table 7.7, and Fig. 7.6.

These results highlight an important feature of scaling laws for quantities that are time dependent. Rates of per capita behavior scale with $N^{\beta-1}$, thus, under

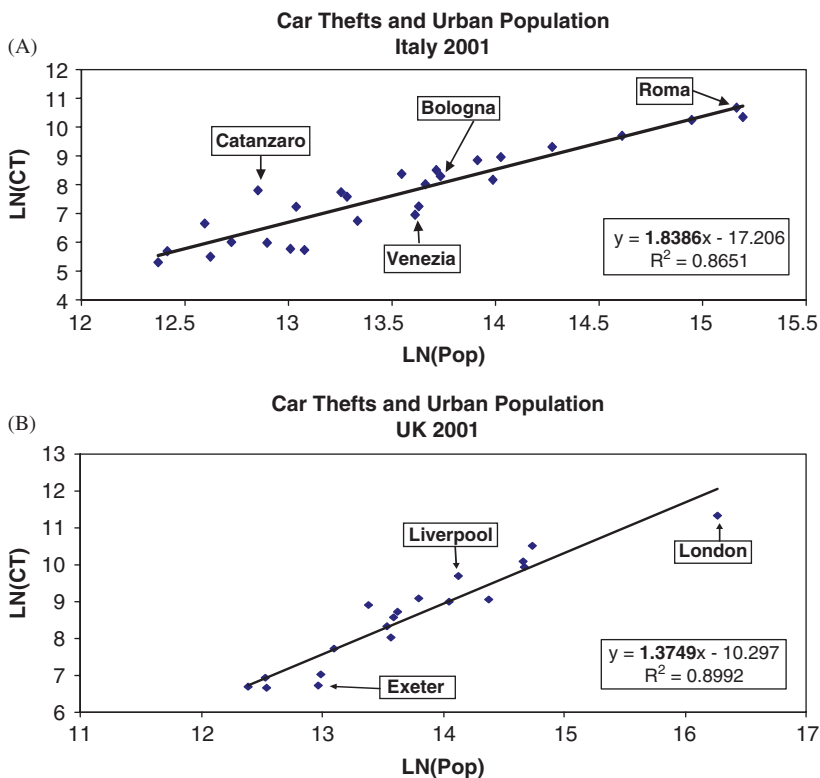


Fig. 7.6 Car thefts (per year) in Italian (A) and British (B) cities show superlinear scaling with metropolitan size. Cities above the scaling law are more theft prone than expected for their size

superlinear scaling, contrary to biology, the pace of social life (measured in disease incidence crime or indeed walking speed (Bornstein & Bornstein, 1976) increases with city size. Life is indeed faster in the big city.

7.5 Implications of Urban Scaling for Growth and Development

Growth is constrained by the availability of resources and their rates of consumption. Resources, Y , are utilized for both maintenance and growth. If, on average, it requires a quantity R per unit time to maintain an individual, and a quantity E to add a new one to the population, then this is expressed as

$$Y = R N + E (dN/dt), \quad (7.2)$$

where dN/dt is the population growth rate. This leads to the general growth equation:

$$\frac{dN}{dt} = \frac{Y_0}{E} N(t)^\beta - \frac{R}{E} N(t). \quad (7.3)$$

Its generic structure captures the essential features contributing to growth. Although additional contributions can be made explicit, they can typically be incorporated by a suitable interpretation of the parameters Y_0 , R and E , leaving the general form of the equation unchanged. For simplicity, we assume that R and E are approximate constants, independent of N . The solution of (7.3) is given by

$$N(t) = \left[\frac{Y_0}{R} + \left(N^{1-\beta}(0) - \frac{Y_0}{R} \right) \exp\left[-\frac{R}{E}(1-\beta)t\right] \right]^{\frac{1}{1-\beta}} \quad (7.4)$$

This equation exhibits strikingly different behaviors depending on whether $\beta < 1$, > 1 or $= 1$. When $\beta = 1$, the solution reduces to classic exponential growth:

$$N(t) = N(0)e^{(Y_0-R)t/E}, \quad (7.5)$$

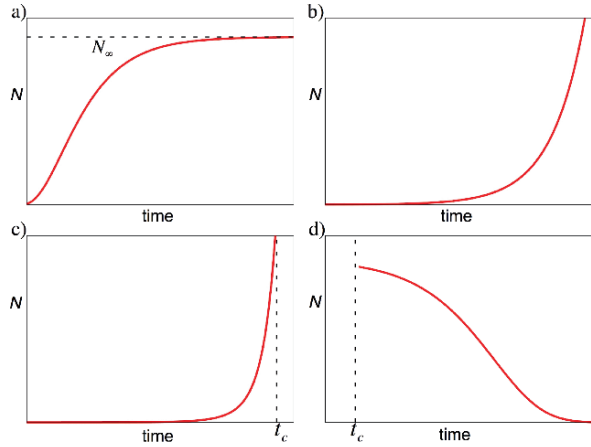
as shown in Fig. 7.7B, while for $\beta < 1$ it leads to a sigmoidal growth curve in which growth ceases at large times ($dN/dt = 0$), where the population approaches a finite carrying capacity given by

$$N(\infty) = (Y_0/R)^{1/(1-\beta)}, \quad (7.6)$$

as shown in Fig. 7.7A. This is characteristic of biological systems where the predictions of (7.3) are in excellent agreement with data. Thus, cities driven by economies of scale are destined to eventually stop growing.

The character of the solution changes dramatically when $\beta > 1$. If $N(0) < (R/Y_0)^{1/(\beta-1)}$, then (7.3) leads to unbounded growth for $N(t)$ (Fig. 7.7C). Growth becomes faster than exponential eventually leading to an *infinite* population in a *finite* amount of time given by

Fig. 7.7 Regimes of urban growth. Plots of size, N , vs. time t : **(A)** Growth driven by sublinear scaling eventually converges to the carrying capacity N_∞ . **(B)** Growth driven by linear scaling is exponential. **(C)** Growth driven by superlinear scaling diverges within a finite time t_c (dashed vertical line) **(D)** Collapse characterizes superlinear dynamics when resources are scarce



$$t_c = -\frac{E}{(\beta - 1)R} \ln \left[1 - \frac{R}{Y_0} N^{1-\beta}(0) \right] \approx \left[\frac{E}{(\beta - 1)R} \right] \frac{1}{N^{\beta-1}(0)}. \quad (7.7)$$

For a city of about a million, t_c is in the order of a few decades. These results highlight an important characteristic of our social mechanisms to generate innovation and wealth. Even as we strive to accelerate prosperity and creativity, we sow the seeds for a crisis, manifested by a finite time singularity, where adaptation processes in society will break down. These crises can be avoided if major adaptations reset the dynamics to generate successive cycles of superlinearly driven growth as shown in Fig. 7.8. These expectations are borne out by data on the population growth of New York City or for the entire world population (Kremer, 1993; Cohen, 1995; Kurzweil, 2005).

Population

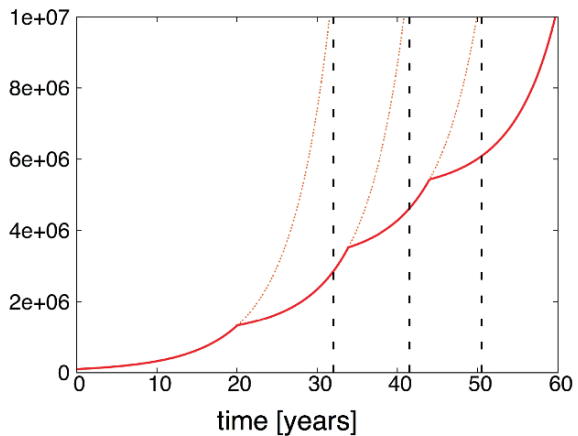


Fig. 7.8 Successive cycles of superlinear innovation reset the singularity and postpone instability and subsequent collapse. The vertical dash lines indicate the location of the sequence of potential singularities. Equation (7.7), with populations of the order of a million, predicts t_c in decades

7.6 Summary and Discussion

We have shown that power law scaling is a pervasive property of human social organization and dynamics in cities and holds across time and for different nations with very different levels of development, economic sector distribution, and with different cultural norms and geographic location. The existence of scaling laws signifies that cities within the same urban system (usually a nation) are self-similar. This is an extraordinary assertion indicating that, on average, different cities are scaled up versions of each other, particularly in terms of rhythms of social activity – such the creation of wealth and ideas, infectious contacts and crime, and patterns of human behavior – even if individual cities vary enormously in terms of their population constitution (e.g., age, race, ethnicity), geographic characteristics and countless other factors.

Urban scaling reveals a tension between quantities that constitute material infrastructure (length of cabling, road surface, etc.) and those that are eminently social (wealth, idea creation, etc.). Larger population densities allow for economies of scale in terms of infrastructure, but these may be driven to less than optimal operation by the requirements of social activity. A theory that encapsulates these compromises and is predictive of scaling exponents is a central objective for future research.

Particularly important are the consequences for growth of urban resource availability driven by innovation ($\beta > 1$) or economies of scale ($\beta < 1$), see summary in Table 7.8. The latter implies growth that eventually slows down, and an ultimate limit to the size of a city, in analogy to growth in biological organisms. The former is radically different, and probably unique to human social organization. It implies accelerating growth, towards a finite time singularity, thus linking inextricably the desired properties of fast economic and technological development to crises of adaptation. Growth in the superlinear regime never converges to a static equilibrium, defying common theoretical assumptions in economics. Instead, it requires constant adaptation to complex new situations created by faster and more efficient human social contact, both desirable and pathological. In particular, major adaptations must occur to reset growth under superlinear scaling to manageable levels, possibly explaining the cyclic nature of most instances of population and economic growth, as well as of technological development.

In closing, we would like to stress in this volume, in the spirit of the contribution by Sander van der Leeuw and collaborators, that cities can be seen as very

Table 7.8 Classification of scaling exponents for urban properties and implications for growth

Scaling Exponent	Driving Force	Organization	Growth
$\beta < 1$	Optimization, Efficiency Creation of Information,	Biological	Sigmoidal, Long term stagnation
$\beta > 1$	Wealth and Resources	Sociological	Boom/Collapse, Finite time singularity, Increasing acceleration/discontinuities
$\beta = 1$	Individual Maintenance	Individual	Exponential

large-scale social information engines, producing open ended innovation and wealth (as well as waste and pollution) out of incoming population, energy, and other resources. As cities grow, disproportionately large numbers of their parts – in terms of population and institutions – are dedicated to innovation, forcing their population either into cycling out of the city or towards adaptation to new roles and behaviors. It is perhaps this necessity for the city as the engine of human social development that makes *man a political animal by nature*. It may well be that the self-similarity revealed by urban scaling laws is the clearest quantitative expression of our unique human social nature, and its understanding the key to a future where sustainability and creativity can coexist.

References

- Bettencourt, L. M. A., Lobo, J., Herbing, D., Kuehnert, C., & West, G.B. (2007). Growth, innovation, scaling, and the pace of life in cities. *Proceedings of the National Academy Sciences of the USA*, 104, 7301–7306.
- Bornstein, M. H., & Bornstein, H. G. (1976). The pace of life. *Nature*, 259, 557–559.
- Botkin, D. B., & Beveridge, C. E. (1997). Cities as environments. *Urban Ecosystems*, 1(3), 3–19.
- Cohen, J.E. (1995). Population growth and earth's human carrying capacity. *Science*, 269, 341–346.
- Decker, E.H., Elliott, S., Smith, F.A., Blake, D.R., & Rowland, F.S. (2000). Energy and material flows through the urban ecosystem. *Annual Review of Energy and the Environment*, 25, 685–740.
- Durkheim, E. (1964). *The division of labor in society*, New York: Free Press.
- Enquist, B. J., Brown, J. H., & West, G. B. (1998). Allometric scaling of plant energetics and population density. *Nature*, 395, 163–166.
- Florida, R. (2004). *Cities and the creative class*, New York: Routledge.
- Girardet, H. (1992). *The Gaia atlas of cities: New directions for sustainable urban living*, London: Gaia Books.
- Glaeser, E. L., & Sacerdote, B. (1999). Why is there more crime in cities? *Journal of Political Economy*, 107, S225–S258.
- Glaeser, E. L., Kolko, J., & Saiz, A. (2001). Consumer city. *Journal of Economic Geography*, 1, 27–50.
- Graedel, T. E., & Allenby, B. R. (1995). *Industrial Ecology*, Englewood Cliffs, NJ: Prentice Hall.
- Henderson, V. (2003). The urbanization process and economic growth: The so—what question. *Journal of Economic Growth*, 8, 47–71.
- Kremer, M. (1993). Population growth and technological change: 1,000,000 B.C. to 1990. *Quarterly Journal of Economics*, 108(3), 681–716.
- Kurzweil, R. (2005). *The Singularity is Near*, New York: Viking.
- Lucas, R. (1988). On the mechanics of economic development. *Journal of Monetary Economics*, 22, 3–42.
- Macionis, J. J., & Parillo, V. N. (1998). *Cities and urban life*, Upper Saddle River, NJ: Pearson Education Inc.
- Milgram, S. (1970). The experience of living in cities. *Science*, 167, 1461–1468.
- Miller, J. G. (1978). *Living systems*, New York: McGraw-Hill.
- Romer, P. (1986). Increasing returns and long-run growth. *Journal of Political Economy*, 94(5), 1002–1037.
- Romer, P. (1990). Endogenous technological change. *Journal of Political Economy*, 98, S71–S102.
- Segal, D. (1976). Are there returns to scale in city size? *Review of Economics and Statistics*, 58, 339–350.

- Simmel, G. (1964). The metropolis and mental life. In K. Wolff (Ed.), *The sociology of George Simmel* (pp. 409–424). New York: Free Press.
- Sveikauskas, L. (1975). The productivity of cities. *Quarterly Journal of Economics*, 89, 393–413.
- West, G. B., Brown, H. H., & Enquist, B. J. (1997). A general model for the origin of allometric scaling laws in biology. *Science* 276, 122–126.
- West, G. B., Brown, H. H., & Enquist, B. J. (1999). The fourth dimension of life: Fractal geometry and allometric scaling of organisms. *Science* 284, 1677–1679.
- West, G. B., Brown, H. H., & Enquist, B. J. (2001). A general model for ontogenetic growth. *Nature*, 413, 628–631.
- Wirth, L. (1938). Urbanisation as a way of life. *American Journal of Sociology*, 44, 1–24.

Chapter 8

Innovation Cycles and Urban Dynamics

Denise Pumain, Fabien Paulus and Céline Vacchiani-Marcuzzo

8.1 Introduction: Urban Systems are Adaptive Systems

Urban systems are adaptive systems, in the sense that they continuously renew their structure while fulfilling very different functionalities. Many examples of adaptation in city size, spacing, and their social and functional components have been given in Chapter 6 of this book. There, we defined the structure of urban systems as a rather persistent configuration of relative and relational properties differentiating cities, which, over long periods, maintains the same cities in categories of size or socio-economic specialization. The content of these categories changes in terms of the quantitative thresholds or the qualitative attributes used for defining them at each date, but they retain the same meaning in terms of the relative situation of cities in the urban systems. Hierarchical differentiation and socio-economic specialization are the major structural features shared by all city systems. On the scale of national, continental, or world urban systems, the structures result mainly from self-organization processes, even if intentional decisions made by individuals or institutions (for instance, the choice of Brussels for the seat of many European Union institutions) may sometimes influence the general configuration.

In the present chapter, we emphasize how the process of innovation is essential in shaping the structure and dynamics of urban systems. Feedback processes can be observed, through which social and technological change occurs in every town and city, while the particular features of this propagation of innovation determine functional and size differentiation among cities. In addition, the spontaneous organization of systems of cities encourages further production of innovation. There is an incentive to innovate that stems from the very structure of urban systems. Urban systems are viewed as subsets of cities involved in a multiplicity of exchanges, through different networks that use these exchanges for a variety of economic, political and social functions relating to operation, management, or control. The exchanges that take place in these networks also convey information about innovation. While most innovations induce smooth change, without any deep structural transformation and

D. Pumain (✉)
Université Paris I Panthéon Sorbonne, UFR de Géographie, Paris, France

only slightly affect the urban hierarchy, some of them emerge in correlated bundles, which can accelerate the hierarchization process, or even lead to the emergence of new types of cities, via specialization.

Using different examples at different times in history and in different parts of the world, we demonstrate how the urban hierarchy is linked to the hierarchical and selective process of diffusion of innovation, as well as to the improvement in transportation technologies. The dynamics of competition inherent in urban systems at once activates, and, is reactivated by, the innovation process. We also discuss the vulnerability of specialized cities and the conditions for their resilience.

8.2 Innovation and Hierarchical Structure of Urban Systems

Here, we discuss the feedback between the innovation process and the hierarchical structure of urban systems. First, we analyze how this structure constrains the propagation of innovation that in most cases takes the form of a hierarchical diffusion process. Second, we recall how this diffusion process and the correlated distributed growth process in systems of cities shape their hierarchical structure. Third, we show that the asymmetries and staggered time-lapses in this process reinforce the urban hierarchy over time by introducing a hierarchical selection within urban systems.

8.2.1 Innovation Propagation in Urban Systems: Hierarchical Diffusion

We define innovation as an invention that has become socially accepted. Innovation can be of various kinds and includes new products, or new technology, as well as new social practices, which in general are more long-lasting than mere fashions. More specifically, we define an innovation cycle as a bundle of new products, new economic activities, new professions, and the accompanying new social practices that are created more or less simultaneously over a rather short period, because they rely on the same kind of technology (such as the type of energy used or a production process). Because of the correlation between the multiple features of change, these innovation cycles may have a large impact in differentiating cities in urban systems (Hall & Preston, 1988).

Spatial diffusion of innovation is the term used to refer to the process of adopting this type of grouped innovation by cities. This process is by no means passive and spatially homogeneous. On the contrary, cities (i.e. the social entities that are stakeholders that were or are investors in these urban places, such as firms, local authorities, or private individuals) are engaged in permanent competition with other cities (i.e. with stakeholders investing or living there) for the capture of as many innovations as early as they can. Indeed, there is an economic advantage attached to the early adoption of innovation because the profits are maximum, not at the very beginning of its existence (the risks of failure and the costs of testing various options are still high then) but in the early stages of diffusion, when the prices are still high (before production processes and consumption become widespread) (Pred, 1966).

Hägerstrand (1952) was the first to formalize the propagation of innovation among towns and cities as a hierarchical diffusion process: the largest cities are the first to capture the benefit of the innovation, then the innovation filters down the urban hierarchy, according to urban size, through imitative or competitive processes: the larger cities adopting first, then the medium size cities, and later the smallest towns. The early adoption in largest cities is easily explained by the high levels of information and skilled labor and the diversity and capacity of infrastructures, that are the distinctive attributes of large cities (these attributes being themselves the result of previous successful adaptation to previous cycles of innovation, i.e. the consequence of an intrinsic historical path dependency in the dynamics of urban systems). The largest cities are those that have benefited from their adaptations to many successive innovation cycles, which explain their large sizes. As a consequence, they have also developed broader diversity of activities and attained higher levels of social and organizational complexity. These characteristics explain why they have a greater probability of developing further innovation at an early stage.

8.2.2 *Innovation, Distributed Growth and Hierarchical Structure*

When a city adopts an innovation (or is selected as a place to produce the corresponding goods and services), there is generally a return, including profits that are generated by the new activities, as well as indirect benefits. This process was modelled a long time ago under the economic base theory (Ullman, Dacey, & Brodsky, 1971). Urban growth is greatest in the emergence stage (because of the initial advantage associated with a new production), so that at each time period, advanced cities keep pace with innovation, draw returns from it and grow, whereas non-adapting cities grow less (or not at all). Urban growth may be translated, in variable proportions, as an increase in population or general wealth, but it can also include more qualitative aspects such as changes in human capital, acquisition of knowledge, and diversification of local resources.

In Chapter 6 of this book, we described the dynamics of urban growth in a variety of urban systems. We demonstrated that, once a territory has stabilized under political control, and towns and cities have been established until they completely fill the geographical space under consideration, the systems of towns and cities, whatever their former history and corresponding morphological features, evolve according to a common process we call *distributed growth*. Over short time intervals, there is a wide variation in city growth rates, but many seemingly random fluctuations between time intervals, so that over long periods, all cities grow at the same rate on average (Pumain, 2000). In Fig. 8.1, we illustrate how this distributed growth process, which progressively adapts the system of cities to a larger size without changing its basic hierarchical structure, is linked to the qualitative changes that occur in towns and cities because of the spatial diffusion of innovations. Paulus (2004) performed multivariate analyses on the distribution of the labour force among economic activities for 354 French *aires urbaines* observed at five points in time between 1962 and 1999. The trajectories in Fig. 8.1 connect the successive positions of the same city in the plane defined by the first two factors of

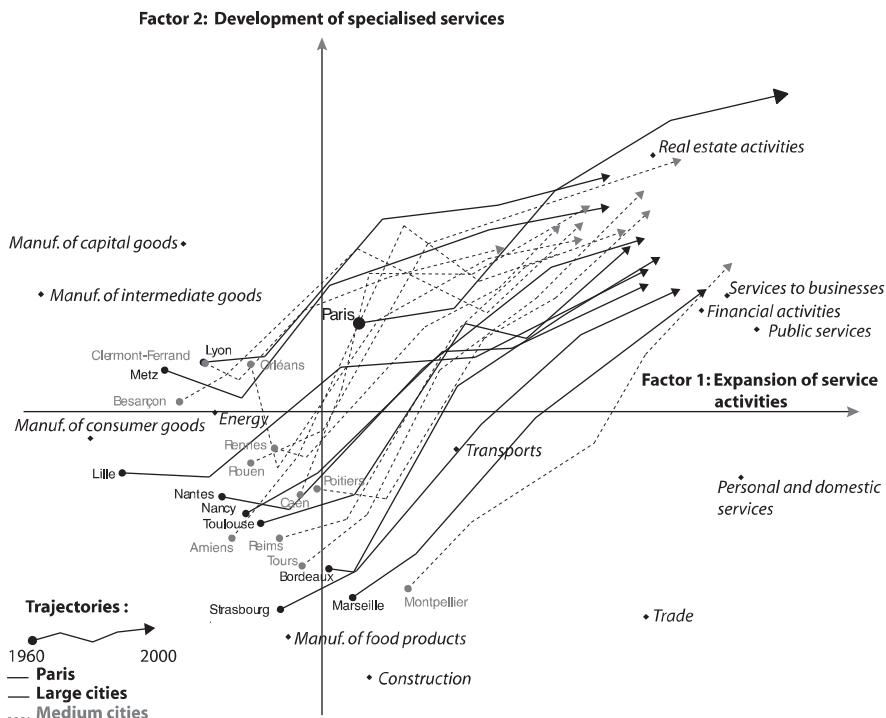


Fig. 8.1 Co-evolution of cities in socio-economic space

Source: Fabien Paulus (2004)

a principal component analysis for that period of time (for the sake of readability, only the largest cities are represented in the figure). The parallelism of the curves is striking. It shows that all cities have registered the same trends in the transformation of their socio-economic structure over that interval of fifty years, reflecting a general decrease in manufacturing activities and a transition from traditional retail and services towards more specialized business, financial, and administrative services.

Figure 8.1 thus illustrates a wide diffusion of the innovations of the time in all parts of the city system. These innovations are expressed here for economic sectors but they follow similar trends when professions or skills are considered. It is because of this general adaptation to socio-economic change that the fundamental structure of the system is maintained: as changes diffuse rapidly everywhere, the initial relative differences are not greatly modified. The diffusion of innovation not only contributes to maintaining the differentiations among cities, but it also explains why all urban systems have a highly skewed distribution of city sizes. The French statistician Gibrat (1931) was the first to demonstrate how a stochastic process of urban growth generates the hierarchical distribution of city sizes (as a lognormal distribution or “law of proportional effect”). This model has been tested many times on a number of urban systems (Robson, 1973; Pumain, 1982; Guérin-Pace, 1995; Moriconi-Ebrard, 1993; Bretagnolle, 1999; Gabaix & Ioannides, 2004) and can be

accepted as a rough but rather robust first approximation of the growth distribution in urban systems. However, the most frequently observed deviations go against the hypothesis of independence of cities, since they include periods of temporal auto-correlation between growth rates and a slight trend towards positive correlation (or negative in some urban systems, especially in the US) between growth and city size. The challenge now is to find a better model for describing urban change, which includes interactions between cities, and to explore further the connection with the innovation process (Favaro, 2007).

It has also been observed in empirical studies that urban growth rates are linked with innovation cycles (Berry, 1991), involving a changing relationship with city size: at the beginning of a cycle, larger cities tend to grow faster, then growth rates tend to equalize, then small towns tend to grow faster (Robson, 1973). This last stage was even interpreted as a “counterurbanization” trend (predicting a decline or even the “end” of the largest metropolises) during the years 1970–1980, while in fact it simply marked the declining stage of the post Second World War innovation cycle (Cattan, Pumain, Rozenblat, & Saint-Julien, 1994). However, in the long run, and partly because of these time lapses and partly because of selection processes during the diffusion of innovations, the urban growth process is not purely stochastic. Some cities grow significantly faster, and others undergo relative or even absolute decay, to an extent that cannot be predicted from a homogenous stochastic model.

8.2.3 Hierarchical Selection and Reinforcement of Urban Hierarchies

The consequence of the early adoption of innovation by large cities is that they draw greater benefit from the innovation (the initial advantage), and this is translated into a persistent tendency for their growth rates to be slightly above the average of towns and cities overall. We have seen in Chapter 6 that the inequalities in city sizes increase over time, especially during periods of fast growth. This *hierarchical selection* is reinforced by the logically correlated trend of smaller towns having growth rates below the mean, either because they adopt innovation when the associated benefits are becoming much smaller or because they are never reached by the innovation. The latter is especially likely when rapid transportation modes or infrastructures are considered: it was observed for the diffusion of railways, free-ways, and, more recently, high-speed trains and airports (Bretagnolle, Paulus, & Pumain, 2002). But it was also observed in the case of much older networks, using less sophisticated transportation technologies.

8.3 Innovation and Specialization: Emergence and Persistence of Urban Geodiversity

Besides the effects of hierarchical selection, a second type of asymmetry is created in urban systems by the innovation process. Sometimes, the resources for which exploitation becomes profitable are not available in every location; they give rise

to *urban specialization* because the related economic activities can only develop in a few urban sites. Usually, the development of a new urban specialization gives rise to exceptionally high growth rates, as the booming cities attract migrants and profit-oriented investments. This has been the case when certain mineral resources became exploitable, such as in 19th century coalmines in the British Midlands or Belgian Wallonia, gold in California, or diamonds at Kimberley, South Africa. In this line of thinking, we should remember that several location factors that, in a deeper past, explained the emergence and success of many towns and cities are also geographical “accidents” of another kind, those influencing the layout of long distance trade routes, for instance wide valleys or topographic corridors, estuaries and bays for maritime routes, major crossroads, or contact points between different regions. Among the more recently exploited spatially concentrated resources are mountain slopes for skiing or coastlines for tourism. The places where public or private funds have been injected to create a local concentration of skill and knowledge can also be considered as nodes of possible concentration of investment and urban specialization, for instance the large universities that have generated “technopoles.” The recurrence of specialization processes and their effect on cities’ development explains that, over long periods, the positive correlation between city sizes, diversification of urban activities, and urban growth tends to decrease with the length of the time interval under consideration (Shearmur & Polese, 2005).

Since medieval times, less than a dozen large innovation cycles have left traces still visible today. An example in the European urban system is the existence of specialized cities (Fig. 8.2). Even if the economic cycle that gave rise to these urban specializations has been over for a long time, cities carry the marks of the momentary intensity of investments that were made at the relevant time. This is visible in their architecture and in the collective representations, what is called the “city image.” Urban marketing experts use architectural and urban heritage, as well as mental connotations that are associated with the city, as a resource for developing the city’s

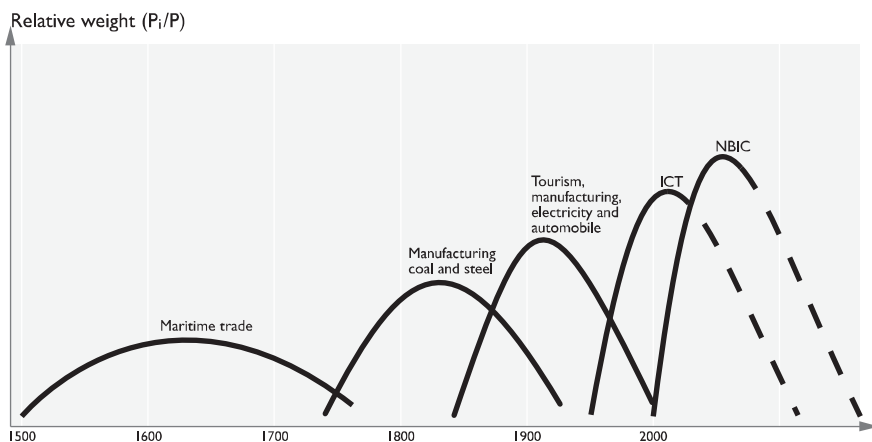


Fig. 8.2 Main innovation cycles having generated urban specialization in Europe

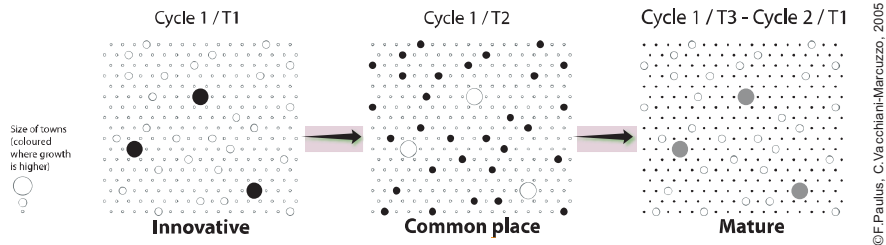


Fig. 8.3 Diffusion and substitution in the urban system

attractiveness for tourism purposes, for new types of economic investors, and for consolidating a social consensus around an urban development project (Markusen & Schrock, 2006).

Thus, urban specializations are explained partly by the unequal diffusion of some innovation cycles that are linked to spatially concentrated resources. But they may also result from the hierarchical diffusion process itself. For stakeholders, the interpretation of the hierarchical diffusion of innovation can be considered as a rational choice (in the sense of economic game theory), which is constrained by the hierarchical configuration of urban systems at a higher level. Duranton and Puga (2001) have recently explored the linkages between initial location and subsequent relocation of firms according to the industrial cycle to which they belong. Firms can only find the urban diversity that favors their learning of technologies and procedures in large cities, and they subsequently relocate to places where they can develop their activity. The costs of investment in a new product or services, which are higher in large cities (higher rents and wages), can only be borne in the first stages of innovation diffusion, when expected profits are still high; in later stages, the activity becomes profitable only in smaller towns where the costs are lower. However, this purely economic reasoning can hold only if the accompanying social potential can follow, in other words, when the required knowledge has percolated through the spatially differentiated social system from large metropolises towards smaller towns. Fujita and Hamaguchi (2001) thus explain how large cities favor innovation while relocation leads to more specialized urban places. Figure 8.3 gives a diagrammatic representation of how innovative activities corresponding to a given innovation cycle locate at first (time T1) in the largest cities, then diffuse to medium size towns (time T2), and become restricted to certain specialised small towns (time T3), while activities of another innovation cycle emerge in the largest cities.

8.4 Innovation Cycles and Scaling Properties of Urban Systems

The largest cities become larger because they were successful in adopting many successive innovations. Many of these innovations later become part of the activity of all towns and cities, since they meet needs that become commonplace (for instance, the primary and secondary education and health services in cities of the developed world today). But the functioning costs in these large urban areas are also much

higher, and many activities are forced to migrate out to smaller settlements where they can sustain their economy. So at each time period, the activities belonging to a new cycle of innovation remain circumscribed for a while at that upper level of the urban hierarchy, then diffuse among other cities, then become more restricted, first escaping from the largest cities, and finally remaining concentrated in only a few small towns.

Thus at a given moment, it can be expected that the most advanced technologies concentrate in the largest cities, while current technologies are ubiquitous and outdated technologies remain only in small towns. The corresponding activities can then exhibit three different scaling parameters for a general model

$$S_{ij} = P_i^\beta \quad (8.1)$$

where the importance of economic sector j in city i (measured by employment or production) is expressed as function of the city size P_i :

Leading technologies (top of current innovation cycle): $\beta > 1$

Commonplace widespread technologies (diffusion stage): $\beta = 1$

Mature technologies (decay or substitution stage): $\beta < 1$

This model was applied to three urban systems. Our first test of the theory is based on the distribution of the labor force in 276 French urban areas (the largest “*aires urbaines*”,¹ which are roughly equivalent to the American Metropolitan Standard Areas). Inevitably, the official nomenclatures used for economic activity (32 categories of the NES – *Nomenclature Economique de Synthèse*) do not always correspond exactly to historical innovation cycles: they were not designed for this purpose, even if they are revised from time to time to provide a more apt description of current economic activities (Desrosières, 1993). Because of the somewhat arbitrary aggregation of activity sectors that they give, it would be hazardous to interpret the value of scaling parameters in absolute terms. Moreover, it is obvious that the content of activities that we classify as “mature” at the level of aggregated economic sectors can be just as up-to-date, in terms of technological and managerial processes, as the diffusing or even leading activities, at the level of individual firms. What the model seeks to express is an aggregated spatio-temporal view of the whole system of cities and over very long periods.

8.4.1 *Scaling and Diversity of Urban Functions*

According to the theory above, the activity profile of the largest cities is expected to be more diversified than that of the smallest cities: if large cities successfully adopt

¹ There are 354 “aires urbaines.” We selected the 276 largest in order to have the same number of urban units as in the case of US. The minimal size is then 17,000 inhabitants.

many innovation cycles, they will carry traces of past cycles, so that their functional profile is likely to be more diverse and more complex. In support of this, the number of employees in 20 economic sectors was collected in order to calculate a diversification index (Paulus, 2004). This diversification index D is based on Isard’s specialization coefficient I , that is,

$$D = 1 - I, \tag{8.2}$$

where

$$I = \frac{1}{2} \sum \left| \frac{x_{ij}}{(x_{i.} - x_{.j})/x_{..}} \right|, \tag{8.3}$$

and x_{ij} = employment of activity j in city i , $x_{i.}$ = total employment of city i , $x_{.j}$ = total employment in activity j , and $x_{..}$ total urban employment – I corresponds to the Euclidian distance between the economic profile of the town and the mean profile). When D is close to 1, it indicates that the city’s economic profile is diversified. On the other hand, a city with a diversity index close to 0 has most of its employment concentrated in a single economic sector.

The relationship between city size and economic diversity is clearly visible from the graph in Fig. 8.4. The correlation is strong, with a coefficient of determination equal to 50%, even if variations remain. All “aires urbaines” larger than 200,000 inhabitants belong to the most diversified group of cities. The less diversified cities are only the small ones.

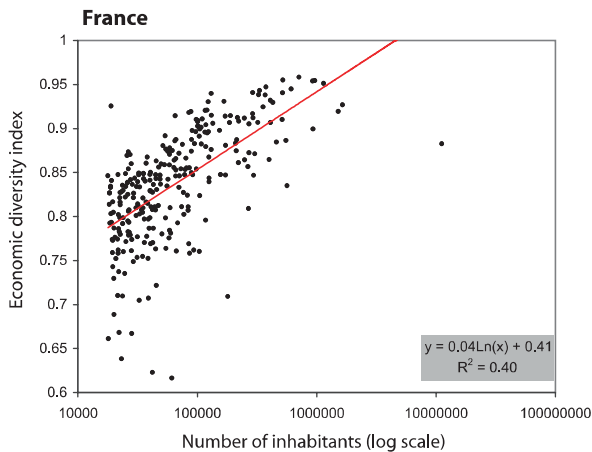


Fig. 8.4 Economic diversity and urban size.
Source: INSEE – Recensement de la population, 1990

8.4.2 *Scaling Parameters and Stage of Activity Sectors in the Innovation Cycle*

Not only is this global indicator in agreement with our interpretation, but the same is true for the results of more detailed investigation about scaling parameters, using data on employment according to economic sector. We plotted cities according to their size (logarithm of the number of inhabitants) on the X-axis and the logarithm of the number of employees in a given economic sector on the Y-axis. To calculate the scaling exponent (β), using the least squares technique, we estimated the slope of the line that fits the set of points. This data set is provided by the last French census, in 1999.

Table 8.1 shows scaling exponents for certain economic sectors classified according to their approximate stage in the innovation cycle. Consultancy and assistance activities, as well as financial activities (Fig. 8.5) are representative of economic sectors which became leaders during the current innovation cycle and emblematic of the “knowledge society.” The β exponents are well above 1. This result confirms that these economic sectors are much better developed in the largest cities and absent or tiny in the smallest ones.

Employment levels in hotels and restaurants can be interpreted as a proxy for measuring the impact of the tourism innovation cycle. Tourism emerged at the end of the 19th century, as long distance travel became faster via railway networks. This activity spread widely during the 1960s and can now be considered as a diffusing activity. The β exponent is close to 1 and the quality of fit is very good. Just a few small towns have many more employees in hotels and restaurants than the average in the urban system. These cities remain specialized.

The manufacture of food products, as a mature industry, scales sublinearly with city size (Table 8.1). This activity was an innovation a long time ago, when it replaced domestic production. Today it remains important in small towns only, and tends to account for smaller proportions in Paris and other large cities, which have

Table 8.1 Scaling parameters and stage of economic sectors in the innovation cycle (France)

Stages in technological development innovation cycle	Economic sector	Power-law exponent (β)	(β) 95% Confidence limits	R ²
Innovative	- Consultancy and assistance activities	1.21	1.17–1.26	0.92
	- Financial activities	1.16	1.11–1.21	0.91
Common place (adapting)	- Hotels and Restaurants (tourism)	1.03	0.99–1.07	0.90
	- Health and social services	0.96	0.93–1.00	0.92
	- Education	0.98	0.96–1.01	0.96
	- Manufacture of food products, beverages and tobacco	0.90	0.83–0.97	0.70

Source: INSEE, Recensement de la population, 1999, 350 aires urbaines

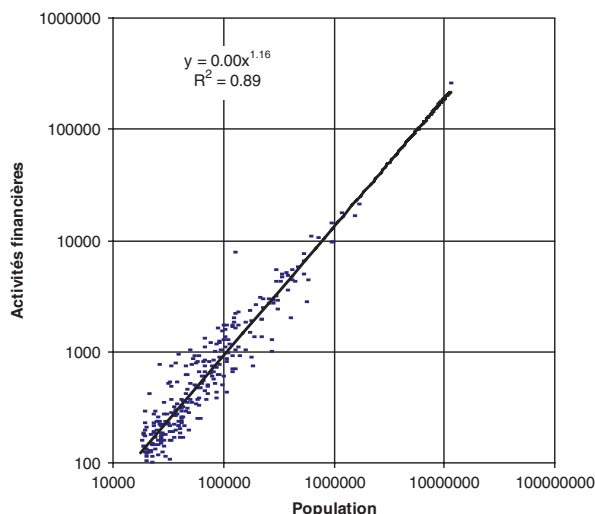


Fig. 8.5 Employment in financial activities as a function of city size (France)
Source: INSEE, Recensement de la population, 1999

been deserted by manufacturing activities since the 1960s (there can of course be other reasons, such as the proximity with places of production, for locating food production in the countryside). We could not identify any other sector scaling sub-linearly with city size, but a much more detailed analysis of the manufacturing sector would be required: we recognize here the inadequacy of the existing nomenclature for the purpose of our study.

We compared the scaling exponents that have been calculated for the United States urban system, using, as urban units, the 276 MAs (258 MSAs and 18 CM-SAs, Porto Rico is excluded – the smallest size is 57,800 inhabitants). Employment data per sector is derived from the Census 2000, using the 20 sectors of the North American Industry Classification System (NAICS) nomenclature.² There is no exact match between this economic nomenclature and that used by the French statistical institute, but reasonable comparisons are possible (Table 8.2). In particular, the number of urban units and the level of disaggregating of activities are close and introduce no bias. Globally, the values of scaling exponents for similar activities are close. The economic sectors that belong to the most recent innovation cycle, including all modern business services like finance and insurance, real estate or scientific services (generally summarized as APS and FIRE), scale superlinearly in both countries with city size. The beta exponent is 1.21 for the professional,

² Actually, NAICS (US nomenclature) and NES (French nomenclature) are both related to the International Standard Industrial Classification of All Economic Activities (ISIC) proposed as guidelines for national classifications by the United Nations Statistical Commission.

Table 8.2 Comparison of scaling parameters for similarly defined economic sectors in US and France

	France	United States
Professional; scientific and technical services/Conseils et assistance et Recherche et Développement	1.21	1.21
Finance and insurance/Activités financières	1.15	1.14
Wholesale trade/Commerce de gros	1.11	1.09
Administrative and support and waste management services/Services opérationnels	1.07	1.11
Accommodation and food services/Hôtels et restaurants	1.04	0.98
Construction	0.99	1.01
Retail trade/Commerce de détail, réparations	0.97	0.98
Health care and social assistance/Santé, action sociale	0.96	0.96
Manufacturing/ensemble des industries	0.92	1.0
Manufacturing/ensemble des industries sans IAA	0.95	

scientific and technical services and 1.15 for finance and insurance sector in both countries (Fig. 8.6). Common activities such as utilities, accommodation, food services or retail trade scale almost linearly with size (for instance, retail trade: 0.97 in both cases). Sublinear scaling would probably characterize some subdivisions of the manufacturing sector if details were provided for the US, as is the case using the French nomenclature of activities.

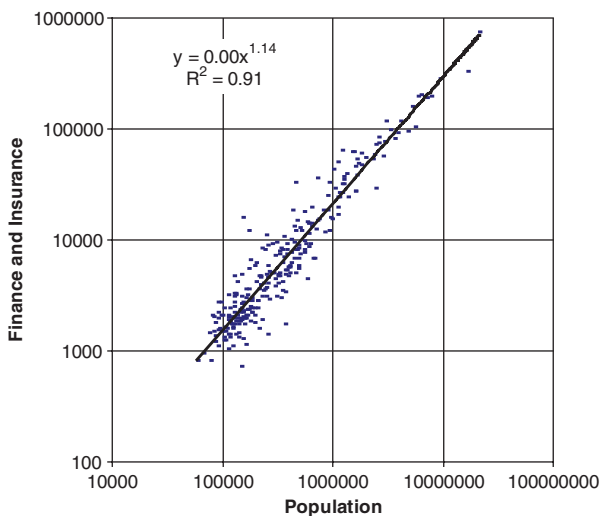


Fig. 8.6 Employment in financial activities as a function of city size (United States)
Source: US Census.

8.4.3 *Scaling Parameters and Hierarchy Among Occupational Groups*

We also applied the same method to occupational groups, as they are described by the French census. Following this nomenclature, it is quite easy to classify the labor force according to average skill levels. We find that the highly skilled jobs scale superlinearly with city size, whereas unskilled jobs (mainly employed workers) do scale sublinearly, and standard skills (such as teachers) are simply proportional to city size (Table 8.3).

Once again, we compared those results on French urban system with the US urban system. We used the SOC nomenclature which is quite different from the French one (PCS) as it takes more into account the hierarchy among the occupation in firms. The French one emphasizes instead the skill criterion and the homogeneity of individuals and households behaviours among social “milieux” (Desrosières, 1993). Nevertheless, we found corroborating results (Table 8.4).

We can display a synthetic view of the strong relationship between professions and city size by using factor analysis. The input tables are fairly simple: for France, it includes eight entries for urban size and five columns for occupational groups; correspondingly eight and seven for US. The less detailed classifications for professional groups are used here and for both cases. On both plots (Fig. 8.7 parts A and B), the distribution of occupations underlines the transition from small towns to the largest cities (Paris alone in the French case). The society of small towns has a concentration of many workers, whom we can consider as unskilled in terms of current technological development. Medium size towns concentrate relatively more skilled employees, such as technicians, clerks and salesmen. The largest cities, especially Paris and MAs from 5 to 10 million inhabitants in the US, house a larger proportion of highly skilled people (executives and professionals). This cross-sectional view corresponds to the historical process of emergence of more skilled activities that

Table 8.3 Scaling parameters and hierarchy of skill among occupational groups in French urban areas

Stages in technological development innovation cycle	Occupational group	Power-law exponent (β)	(β) 95% Confidence limits	R^2
Highly skilled	- Civil servant executives	1.21	1.15–1.26	0.84
	- Management and business executives	1.15	1.11–1.18	0.86
Skilled	- Technicians	1.06	1.03–1.13	0.86
	- Teachers	0.95	0.93–0.97	0.96
Unskilled	- Skilled workers	0.87	0.82–0.92	0.74
	- Unskilled workers	0.80	0.75–0.86	0.70

Source: INSEE, Recensement de la population, 1999

Table 8.4 Scaling parameters of main US occupational groups

Main US occupational groups	Power-law exponent (β)	(β) 95% Confidence limits	R ²
Management, business and financial	1.11	1.09–1.14	0.97
Professional A (1)	1.16	1.12–1.20	0.92
Professional B (2)	0.96	0.94–0.98	0.97
Service (3)	0.97	0.96–0.98	0.99
Sales	1.01	1.00–1.02	0.99
Office and Administrative	1.04	1.02–1.06	0.98
Working-class (4)	0.96	0.94–0.98	0.97

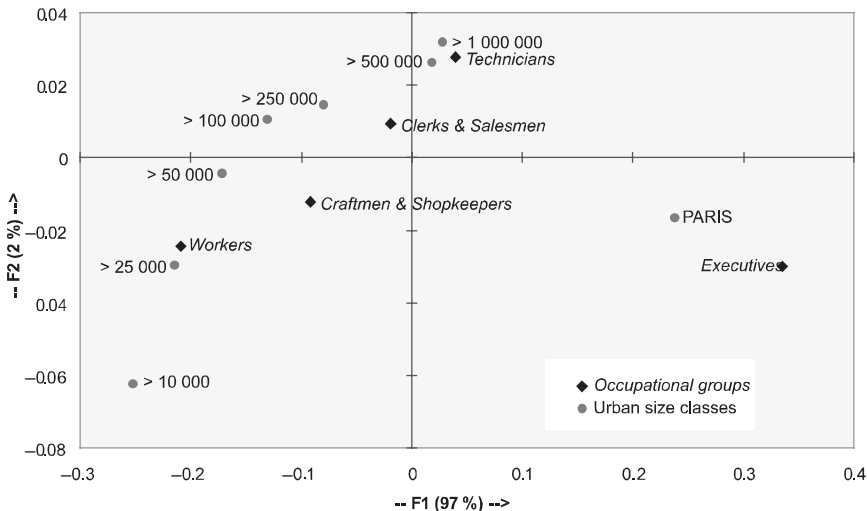
(1) Computer and mathematical – Architecture and engineering – Life; physical; and social sciences – Legal – Arts; design; entertainment; sports; and media; (2) Community and social services – Education; training; and library – Healthcare practitioners and technical; (3) Healthcare support – Protective service – Food preparation and serving – Building and grounds cleaning and maintenance – Personal care and service; (4) Construction and extraction – Installation; maintenance; and repair – Production – Transportation and material moving.

Source: US Census, 2000.

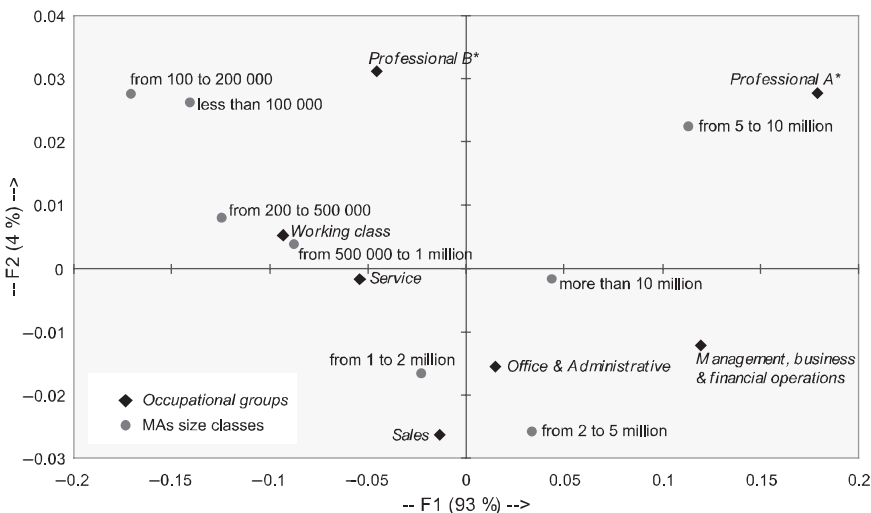
develop first in the largest cities. It represents in a synthetic way the evolutionary process of divergence of human capital levels across cities (Glaeser & Berry, 2005). Unlike the economic nomenclature where types of products change over time almost every couple of decades, more or less mirroring innovation cycles, the nomenclature of professions evolves at a slower pace. It approximates the hierarchy of social status, which in modern societies is linked to the level of skill, even if only loosely. The actual content of skill may change rapidly, while the identification of the corresponding social categories remains almost the same. That is why, unlike economic sectors, the aggregate categories corresponding to the highest skill (executives, intellectual professions), which contribute intensely to innovation, are not likely to diffuse through all cities over time, but remain concentrated in the largest. But if we envisaged professions in a more detailed way, we might find that for instance the highly skilled mechanic who constructed automobiles one by one in the very center of Paris at the beginning of the 20th century corresponds today to a worker in a decentralized plant in a much smaller town. Of course, the equivalence is not easy to establish, and the evolution is not so simple.

8.4.4 A Further Test: Evolution of Scaling Exponents Through Time

Another test of the theory consists in observing how the scaling parameters evolve over time. It can be expected that the now leading technologies can still increase their parameter value, while the activities of older cycles should have decreasing values. Using our historical database on employment per economic sector in French urban areas from 1962 to 1999 (Paulus, 2004), we explored the evolution of the scaling parameters (Fig. 8.8).



A: France; Source: INSEE, Recensement de la population, 1999.

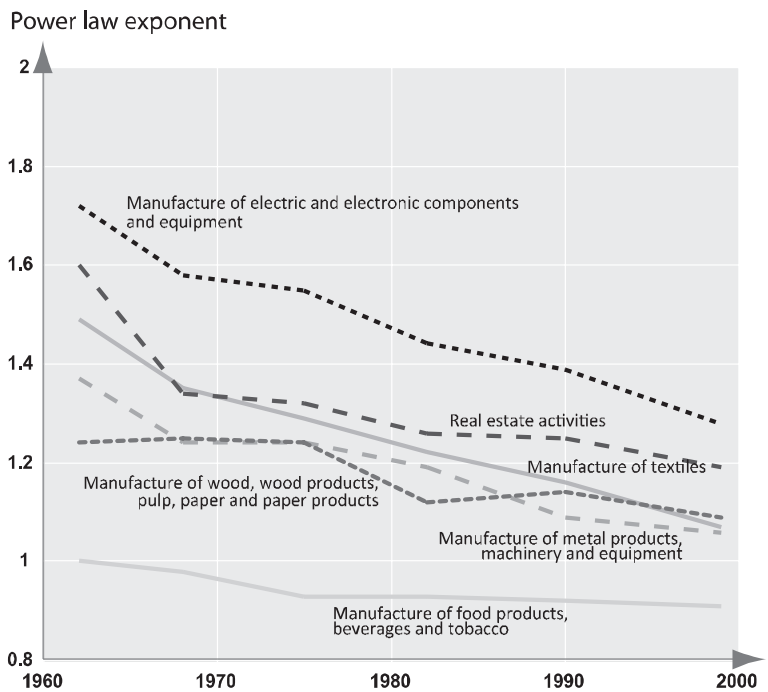
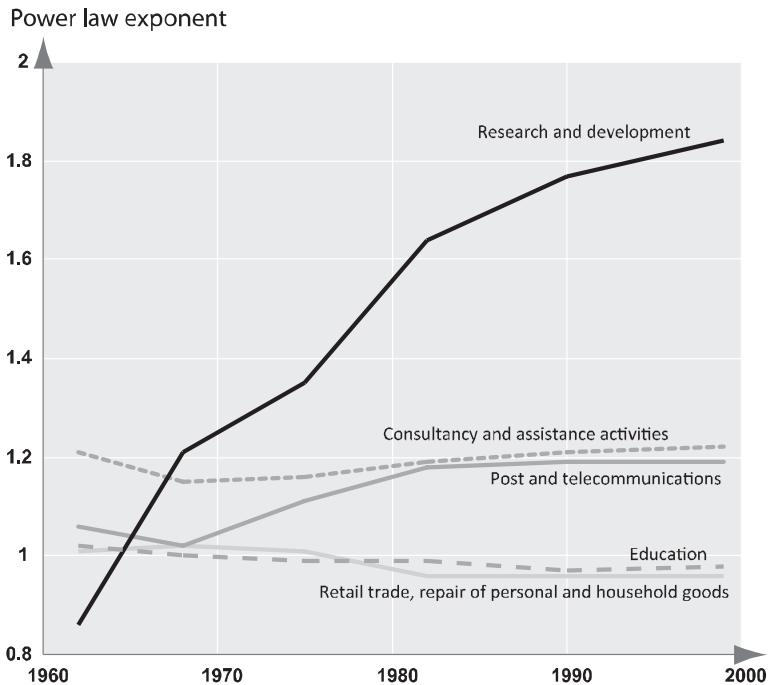


* Professional A : Computer and Mathematical ; Architecture and Engineering ; Life, Physical and Social Sciences ; Legal ; Arts, Design, Entertainment, Sports, and Media.
 Professional B : Community and Social Services ; Education, Training and Library ; Healthcare Practitioners and Technical.

B: US; Source: US Census 2000.

Fig. 8.7 City size classes and occupational groups (Multivariate analysis)

On the upper graph in Fig. 8.8 are represented the sectors which have an increasing or stable β exponent from 1960 to 2000. They are the three economic sectors which are involved in the current innovation cycle: a good example is research and development, where the β exponent was about 0.9 in 1962 and rose progressively to 1.2 in 1968 and 1.84 in 1999. It may seem surprising that the β exponent is



FP © ULP - Image et Ville - 2006

Fig. 8.8 Evolution of scaling parameters (France)
 Source: INSEE, Recensement de la populations, 1999

less than 1 at the beginning, at a time when this sector began to acquire a crucial role in the production system. This can be understood if we keep in mind that, at this time, the total number of employees in this sector was very small and 90% of cities had no employees at all in this economic sector. In this context, some small towns hosting a research center could have as many researchers as larger cities, Paris being an exception. But rapidly, as this research activity became prominent in the economy, with a high rate of employment growth, the largest cities assumed a leading position. In the same context of the growing importance of the society of communication, we also notice that post and telecommunications scale at 1 in 1968 and at 1.2 in 1999. Consultancy and technical assistance activities are more stable, with a β exponent always above 1. On the same graph, education and retail trade exhibit β exponents which are very close to 1 at each date. This can be explained by their fairly stabilized function at this stage of development (although the function has frequently been renewed in its qualitative content), even if it is tempting to interpret their spatially ubiquitous distribution in a static functional perspective, as expressing the satisfaction of basic universal and elementary needs of resident populations. Such activities would probably scale non linearly in countries with a lower level of development or in earlier historical periods.

The lower graph in Fig. 8.8 represents the decreasing values of β exponents for certain economic sectors over time. Most of these sectors are manufacturing industries. This decrease can be understood as a process of hierarchical diffusion of the technological development of these industries, leading towards a higher relative concentration in the lower part of the urban hierarchy. Nevertheless, it should be noted that all these activities retain β exponents above 1. Manufacturing industries are not all mature. Diversity within these manufacturing industries is wide, with some plants that are on the forefront of the technology and some others that belong to an older stage. For example, while numbers of employees in the textile industry are falling, the value of the β exponent is also decreasing but remains above 1 (in 1962, it was equal to 1.5 and at the end of the century, its value is 1.1).

We see here the difficulty in considering that economic activities or urban functions directly reflect different stages in innovation cycles. Economic data provided by the statistical offices are not well suited to pinpointing innovations. Nomenclatures of economic activity sectors are constructed to identify products, and are periodically revised in order to categorize new, innovative productions. But this process of categorization is not systematic: some new sectors have been identified, while old sectors may remain under the same name in the nomenclature with completely new content. A good example would be the automobile industry: at first the small innovative workshops where automobiles were invented did not appear in the nomenclature under any other name than "mechanics." Later, the category "automobile industry" was invented. Today, the category still exists, but the content of the activity has changed, involving robots and all sorts of technological improvements. Nowadays, its content covers both innovation in the production process and a long standing invention in the field of transportation technology. This can explain the poor quality of fit and the fact that the value of the β exponent remains close to 1 after having decreased.

The above point helps to decide between two alternative interpretations of the urban scaling laws. The first is longitudinal, and it considers that the model represents the relationships between the size of a typical town or city and parts of its activities at different stages in its development. The second interpretation of the model is transversal, considering it to represent the distribution of different activities among cities of different sizes at a given period. The first interpretation does not consider the diversity in functional specialization among cities, but interprets the differences in the scaling parameters as reflecting the ability of different activity sectors to adopt a spatial organization that optimizes the trade-off between the advantages and the costs of location in a city of a given size. The second interpretation acknowledges the functional diversity of the system of cities and the progressive substitution of activities at the different levels of the urban hierarchy over time. Our last result showing the evolution of the scaling parameters of certain economic sectors over time seems to support this interpretative framework. This would also remain consistent if we enlarge our observation of urban system to extend beyond national boundaries, including international division of labor and delocalization of activities with lower technological levels and fewer requirements for the skill of labor force towards countries where the production costs are lower. (But in these countries, the multinational firms, which represent an “innovation” there, locate according to the hierarchical diffusion process that we mentioned above, see for instance Vacchiani-Marcuzzo, 2005).

At first sight, our theory may appear counter-intuitive, since urban activities which scale sublinearly with city size should, as is the case among biological species for metabolic rates and size of organisms, illustrate a state of better efficiency or adaptation: they have found a mode of organization that provides scale economies. This may be true for some specific urban services such as water or energy supply. To maintain a static interpretation of this type, it must be concluded that innovative activities waste resources, since they are more abundant in the largest places. In the early stages of emergence of a new market system, many resources, both human and material, are not used in an efficient way. Nor do they need to be, given the possibility of monopoly profits. Efficiency is gained over time. Thus we prefer to support our evolutionary view, even if seemingly more complicated, rather than the more general static explanation, which cannot explain in a satisfying way all the empirical evidence we have found. Of course, more empirical testing is required to consolidate our hypothesis.

8.4.5 Innovation in Emerging Countries: The Example of South Africa

Another illustration of the generality of the theory can be provided by exploring countries that are totally different in their economic structure. South Africa is an interesting example because, as an emergent country, its economy is somehow “between” industrialized and developing worlds. There is a huge volume of literature on new patterns of international trade and division of labour between developed

and developing countries, which we will not detail here. Our hypothesis is that, even if the temporality of innovation cycles is different (and probably delayed) when compared to developed countries, any economic “innovation” would adapt to the existing structure of the system of cities and diffuse downwards through the urban hierarchy of city sizes. Of course, testing this hypothesis is not easy because the existing databases do not include in their nomenclatures typologies that would be directly relevant for our theory, in terms of duration of economic cycles.

Nevertheless, Vacchiani-Marcuzzo (2005) analysed the repartition of economic sectors in South African urban agglomerations and built a database that includes 90 urban agglomerations defined according to our criteria of consistent geographical entities (morphological and functional, see Chapter 6). The South African nomenclature of economic activities for the year 2001 is aggregated in ten sectors, which cannot be perfectly compared with those used for French and American cities. Moreover, we know that the informal sector, which is not taken into account in the official nomenclature, represents a non-negligible part of the total employment (one third at least).

Despite these differences in conception, which also reveal qualitative differences in the economic structure of the country, the results of scaling measurements are interesting (Fig. 8.9). They clearly differentiate a very dynamic sector, regrouping

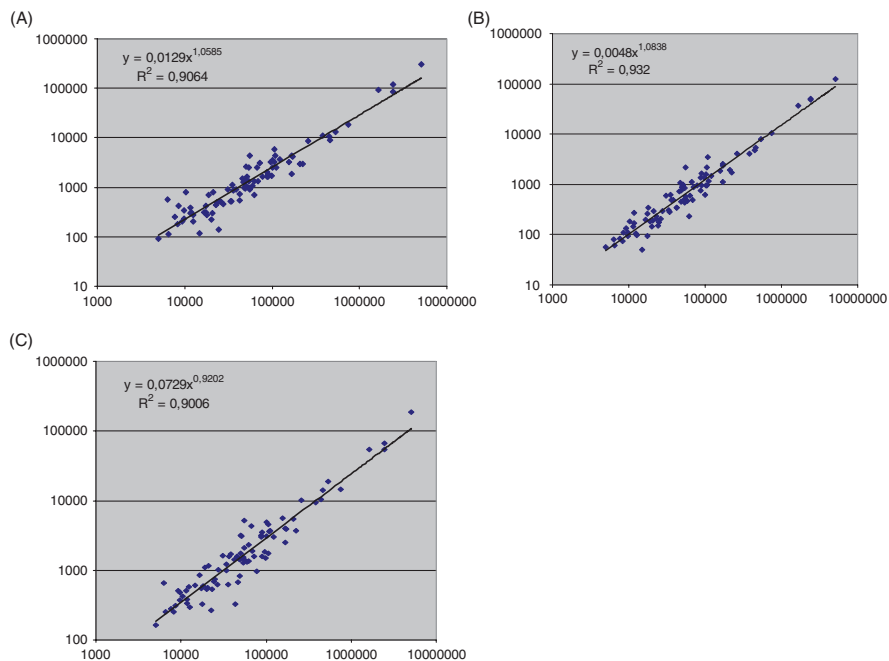


Fig. 8.9 Employment in economic sectors and city size in South Africa. **A:** Leading economic sectors – Finance, Insurance, Real Estate (FIRE); $\beta = 1.06$; 95% CL: 0.98–1.13; $R^2 = 91\%$; **B:** Leading economic sectors – Transportation; $\beta = 1.08$; 95% CL: 1.02–1.15; $R^2 = 93\%$; **C:** Mature economic sectors – Domestic services; $\beta = 0.92$; 95% CL: 0.86–0.98; $R^2 = 90\%$.

Source: Data base CIS-CVM/Census of Population 2001 and Data base CVM

activities of finance, insurance and real estate (FIRE), that scales supralinearly with city size ($\beta = 1.06$), and another one, the domestic services, which scales sublinearly ($\beta = 0.92$) and represents a more obsolete part of the economy. Whereas the FIRE sector reveals that the largest South African cities participate in the economic globalization, the domestic services express a form of social organization that was diffused broadly over all the South African territory, since it corresponds to a very old type of division of labour that was effective in this country. The first process clearly selects the upper level of the urban hierarchy, while smaller towns can still be chosen as locations for very specific economic sectors such as mining industries, which in the past contributed in a very significant way to the creation of the Northern part of the South African urban system. However, the transportation sector represents a specific case: it scales supralinearly with city size with the highest slope ($\beta = 1.08$). This result may appear surprising compared to other countries such as France or USA where this sector is part of mature activities, but on the contrary, in South Africa it reveals the effects of emerging new transport networks and still booming logistic activities.

8.5 Conclusion: Innovation and Sustainability of Urban Development

We shall discuss the relationship between innovation and sustainability of urban systems in more detail in Chapter 12. We want to recall here that innovation is an essential driving force in urban dynamics. Knowledge and information, reflexivity, and the ability to learn and invent provide the impetus for urban development. The crucial role that cities have played in generating innovations – intellectual and material, cultural and political, institutional and organizational – is well documented (e.g. Bairoch, 1988; Braudel, 1992; Hall, 1998; Landes, 1999). The role of cities as centers for the integration of human capital and as incubators of invention was rediscovered by the “new” economic growth theory, which posits that knowledge spillovers among individuals and firms are the necessary underpinnings of growth (Romer, 1986; Lucas, 1988). As Glaeser (1994) points out, the idea that growth hinges on the flow and exchange of ideas naturally leads to recognition of the social and economic role of urban centres in furthering intellectual cross-fertilization. Moreover the creation and repositioning of knowledge in cities increases their attractive pull for educated, highly skilled, entrepreneurial and creative individuals who, by locating in urban centers, contribute in turn to the generation of further knowledge spillovers (Feldman & Florida, 1994; Florida, 2002; Glaeser & Saiz, 2003; Bouinot, 2004). This seemingly spontaneous process, whereby knowledge produces growth and growth attracts knowledge, as the driving force enabling urban centers to sustain their development through unfolding innovation, actually is the result of their adaptive organization. As stressed by David Lane, it is the organization of cities, that provide scaffolding structures where *knowledge* can be generated, developed, stored and accessed, and economic organizations – firms and networks

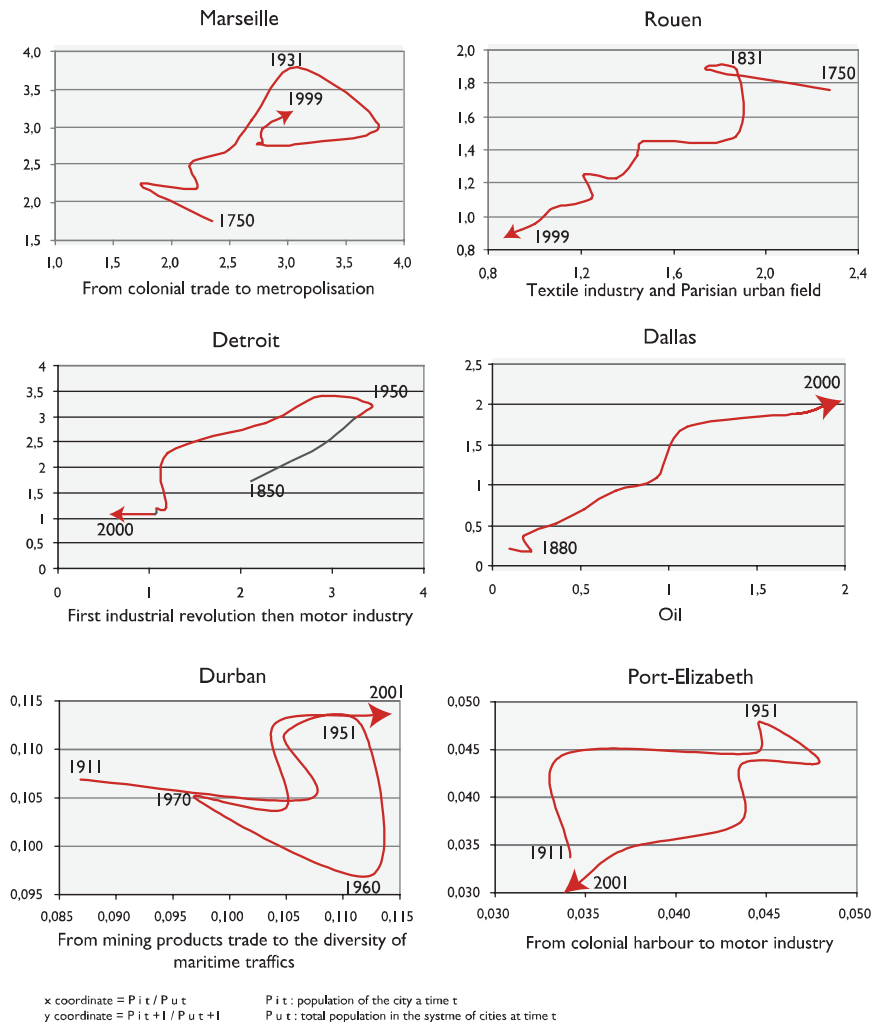
of firms, as well as development agencies etc. – that carry out economic activities – production, exchange, finance and so on – which generate growth. The conceptual definition of a city is a center of organizational activities, and this fundamental functionality is generally undervalued by the advocacy of “agglomeration economies.”

The two processes of diffusion of innovation and specialization have consequences on the dynamics of systems of cities. The activities that can diffuse widely through the system tend to reinforce the relative weight of the large cities, because of the growth advantage and attractivity characterizing the earlier stages of innovation, whereas the activities that focus on a few specialized towns because of some specific location factors, after boosting the development of these towns, sometimes with spectacular growth rates, then tend to hamper their further development by weakening their ability to adapt. Regarding that, predictions could be made about the threatened future of some highly specialized and valued cities, as the tourism centers of our consumption economy, which may encounter future difficulties, although in a lapse of time difficult to predict now.

The many contemporary studies on so-called “metropolization” rediscover a process that has been constitutive of the dynamics of systems of cities (Pred, 1977; Pumain, 1982) at a time when the globalization trends and the general conversion to the “information society” are generating a new broad cycle of innovations. There is an obvious relationship between the maximum possible city size of a metropolis in a given country and the population size of that country. Even if it has not yet been tested because of a lack of relevant data, the same scaling effect exists for the global urban product. Therefore, one must consider the impact of changes in the world economy at the level of nation-states when predicting the future trajectories of large cities.

The question of sustainability of urban development holds for large cities as well as for specialized ones. Watching the past, again in order to think about the future, we can illustrate how the relative weight of a city in the urban system is related to its participation, successful or not, in successive waves of innovation. On Fig. 8.10, the relative size of one city in the system at time t (on X axis) is compared to its relative size at time $t + 1$ (y axis). Examples were chosen for their clear connection between a city relative growth and an innovation: the trajectory of Marseilles is a growing influence linked with the development of colonial trade, then a recession followed by a recent recovery; Rouen is in continuous relative decline for the past two centuries, because of its closeness to Paris and the decline of the textile industry and harbour activities; the whole cycle of growth and decline of the motor industry explains by itself the trajectories of Detroit, and, later, Port Elizabeth, as they are both examples of extreme specialization; Durban has a more complex trajectory including, first, a prosperity stage with colonial trade (in particular, mining products from Johannesburg were exported through Durban), followed by a recession, then a recent new development of harbour activities; Dallas has been driven upward by the oil industry, whose cycle is not yet finished.

It is probable that these different evolutions will continue, within the contemporary context of globalization, no longer at the scale of national systems, but in global urban networks. The colonial period already introduced durable asymmetries



A. Breagnolle, C. Vacchiani-Marcuzzo@UMR Géographie-cités 2006

Fig. 8.10 A few examples of diverging trajectories of cities over time

in urban growth and perturbed the organization of urban systems in colonized countries. The foreign investments and the redistribution of economic activities that characterize the globalization of economies will reproduce or reinforce some of these effects. The dynamics of systems of cities will henceforth have to be analyzed at that scale.

A last important methodological and theoretical remark is about our interpretation of scaling laws in the field of urban geography. We insist that it is not a matter of disciplinary shyness or backwardness if we are reluctant to consider that there would theoretically exist something like an “average ideal city” and that “on average different cities are scaled up versions of each other” where “major adaptations must

occur to reset growth under superlinear scaling to manageable levels” as assumed in Chapter 7. If we resist this “monocity” interpretation of urban scaling laws, it is not because “social sciences emphasize the richness and differentiation of human social expression” as presented in Chapter 7. Our scientific approach is not idiographic but is nomothetic, as well as the physical or biological ones. And according to our view, to be scientifically understood in their development, *cities have to be conceptually represented as elements of a differentiated system of cities*. Cities are not living only on their own resources, but from the valorization of information about distant resources that are more and more located in other cities, enhancing the social and economic power of networks. This is demonstrated by the fact that whatever their cultural, political, economic or historical context *there always are simultaneously in any given territory of the world* towns and cities of different sizes that are also functionally differentiated. Moving investment as well as social value from one city to the next is an essential part of the urban dynamics. That is why our interpretation integrates the hierarchical and functional differentiation of cities within systems of cities as a fundamental part of the explanation of urban scaling laws.

Moreover, we think that the conclusion of Chapter 7 that “it is perhaps this necessity for the city as the engine of human social development that makes man by nature a political animal” has to be reversed. It is because human beings are social animals who developed politics in increasingly large organizations that cities were invented, and that is why they remain the evolutionary expression of the political order of societies. No doubt this discussion between theoreticians of complex systems will continue well beyond the ISCOM project where it was especially fruitful. Arbitration may be provided by a closer observation of cities’ past and future evolutions, using the best quality data as well as the more sophisticated methodologies of complex systems sciences.

References

- Bairoch, P. (1988). *Taille des villes, conditions de vie et développement économique*. Paris, France: EHESS.
- Berry, B. (1991). *Long-wave rhythms in economic development and political behavior*. Baltimore, MD: John Hopkins University.
- Bouinot, J. (2004). Des évolutions dans les comportements spatiaux des entreprises en 2003, Point Chaud, *Cybergeo*.
- Braudel, F. (1992). *Civilisation Matérielle, Economie et Capitalisme, XVe-XVIIIe siècle*. Paris, France: A. Colin.
- Bretagnolle, A. (1999). Les systèmes de villes dans l’espace-temps: Effets de l’accroissement des vitesses de déplacement sur la taille et l’espacement des villes, Thèse de doctorat. Paris, France: Université Paris I.
- Bretagnolle, A., Paulus, F., & Pumain, D. (2002). Time and space scales for measuring urban growth. *Cybergeo*, 219.
- Cattan, N., Pumain, D., Rozenblat, C., & Saint-Julien, T. (1994). *Le système des villes européennes*. Paris, France: Anthropos.
- Desrosières, A. (1993). *La Politique des Grands Nombres. Histoire de la Statistique*. Paris, France: la Découverte.

- Duranton G., & Puga, D. (2001). Nursery cities: urban diversity, process innovation and the life cycle of products. *American Economic Review*, 91(5), 1454–1477.
- Favaro, J. M. (2007). *Croissance urbaine et cycles d'innovation dans les systèmes de villes: une modélisation par les interactions spatiales*. Paris, France: Université Paris I, thèse de doctorat.
- Feldman, M.P., & Florida R. (1994). The geographic sources of innovation: Technological infrastructures and product innovation in the United States. *Annals of the Association of American Geographers*, 84(2), 210–229.
- Florida, R. (2002). *The rise of the creative class*. New York: Basic Books.
- Fujita, M., & Hamaguchi, N. (2001). Intermediate goods and the spatial structure of an economy. *Regional Science and Urban Economics*, 31(1), 79–109.
- Gabaix, X., & Ioannides, Y.M. (2004). The evolution of city size distributions. In V. Henderson & J-F. Thisse (Eds.), *Handbook of regional and urban economics* (Vol. 4, Chapter 53, pp. 2341–2378), Amsterdam, The Netherlands: North-Holland.
- Gibrat, R. (1931). *Les Inégalités Economiques*. Paris, France: Sirey.
- Glaeser, E. L. (1994). Cities, information, and economic growth. *Cityscape 1*, 9–47.
- Glaeser, E. L., & Berry, C. R. (2005). *The divergence of human capital levels across cities*. Harvard Institute of Economic Research, Discussion paper 2091.
- Glaeser, E.L., & Saiz A. (2003). *The rise of the skilled city*. NBER Working Paper Series, 10191.
- Guérin-Pace, F. (1995). Rank-size distribution and the process of urban growth. *Urban studies*, 32(3), 551–562.
- Hägerstrand, T. (1952). *The propagation of innovation waves*, Lund, Sweden: The Royal University of Lund.
- Hall, P. (1998). *Cities in civilization: culture, innovation, and urban order*. London, UK: Weidenfeld and Nicolson.
- Hall, P., & Preston, P. (1988). *The carrier wave: New information technology and the geography of innovation 1846–2003*. London, UK: Unwin Hyman.
- Landes, D. (1999). *The wealth and poverty of nations: Why some are so rich and some so poor*. (2nd ed.). New York: W W Norton
- Lucas, R. (1988). On the mechanics of economic development. *Journal of Monetary Economics*, 22, 3–42.
- Markusen, A., & Schrock, G. (2006). The distinctive city: divergent patterns in growth, hierarchy and specialization. *Urban Studies*, 43(8), 1301–1323.
- Moriconi-Ebrard, F. (1993). *L'urbanisation du Monde*. Paris, France: Anthropos, Economica, Collection Villes.
- Paulus, F. (2004). *Coévolution dans les systèmes de villes: croissance et spécialisation des aires urbaines françaises de 1950 à 2000*. Paris, France: Université Paris 1, thèse de doctorat.
- Pred, A. (1966). *The spatial dynamics of U.S. urban-industrial growth, 1800–1914*. Cambridge, MA: MIT Press.
- Pred, A. (1977). *City systems in Advanced Economies*. London, UK: Hutchinson.
- Pumain, D. (1982). *La Dynamique des Villes*. Paris, France: Economica.
- Pumain, D. (2000). Settlement systems in the evolution. *Geografiska Annaler*, 82B(2), 73–87.
- Robson, B. T. (1973). *Urban growth, an approach*. London, UK: Methuen.
- Romer, P. (1986). Increasing returns and long-run growth. *Journal of Political Economy*, 94, 1002–10037.
- Shearmur, R., & Polese, M. (2005). Diversity and employment growth in Canada 1971–2001: can diversification policies succeed? *The Canadian Geographer*, 49(3), 272–290.
- Ullman, E. L., Dacey, M. H., & Brodsky, H. (1971). *The economic base of american cities*. Seattle, WA: University of Washington Press.
- Vacchiani-Marcuzzo, C. (2005). *Mondialisation et système de villes: les entreprises étrangères et l'évolution des agglomérations sud-africaines*. Thèse, Paris, France: Université Paris 1.

Part III
Innovation and Market Systems

Chapter 9

Building a New Market System: Effective Action, Redirection and Generative Relationships

David Lane and Robert Maxfield

9.1 Introduction

In this chapter, we describe some episodes in the early history of LonWorks, a technology for distributed control networks.¹ LonWorks was introduced in December, 1990 by Echelon, a Silicon Valley start-up company. The launch generated considerable enthusiasm: the editor of one leading control engineering trade journal even predicted that LonWorks would do for control what the microprocessor did for computing. While that prediction still remains to be realized, LonWorks has found many applications in the past 18 years, and Echelon has enjoyed a steady if not spectacular growth over the same period. By end of LonWork's first decade, over 3000 companies had purchased LonWorks development tools, and more than 5000 products incorporating LonWorks technology had been brought to market. Annual sales of these products now probably exceed \$1.5 billion. Since August 1998, Echelon has been a public company, whose shares are traded on NASDAQ. Many leading producers in the building automation and some in the process and discrete control industries now manufacture devices and systems incorporating LonWorks technology. But the adoption of LonWorks technology is by no means limited to companies that belong to these three traditional control industries. Rather, LonWorks has been successfully applied in a wide variety of industries, from transportation² to semiconductor manufacturing³ to food services⁴ to intelligent meter

D. Lane (✉)

Università di Modena e Reggio Emilia, Modena, Reggio Emilia, Italy
e-mail: lane@unimo.it

¹ See Lane and Maxfield (2005) and Chapter 10 of this volume for treatments of later episodes in this story.

² For example, the NY Transit Authority specifies LonWorks for train braking systems, LonWorks has been declared the standard for European gas station forecourts, and Raytheon has developed a fiber optic airplane control system based upon LonWorks.

³ For several years, Echelon's largest-volume customer was a manufacturer of vacuum pumps for semiconductor manufacturing. In 1996, LonWorks was declared a standard for control networks by SEMI, a semiconductor manufacturing trade organization.

⁴ One of the earliest LonWorks applications was by TruMeasur, who used the technology in 1991 to control a liquor dispensing-and-billing system installed in several Las Vegas casino-hotels.

reading.⁵ Partly directing this movement of LonWorks across vast tracts of agent-artifact space,⁶ and partly in response to exigencies that this movement has generated, a new *market system* is emerging, organized around LonWorks technology itself. This chapter is about Echelon's early approach to building this market system.

What is a market system? By a market system, we mean a set of agents that engage with one another in recurring patterns of interaction. These interactions are organized around an evolving family of artifacts. Through their interactions, the agents produce, buy and sell, deliver, install, commission, use and maintain artifacts in the family; generate new attributions about functionality for these artifacts; develop new artifacts to deliver the attributed functionality; and construct, augment and maintain new agents and patterns of agent interaction, to ensure that all these processes continue to be carried out, over time, even as the circumstances in which they take place are changing in response to perturbations from inside and outside the market system itself.

A central theme in our story is the distinction between a market system and "the market," that abstract entity defined formally by economists and employed informally in the popular vernacular. "The market" is a locus of impersonal exchange activities, where agents buy and sell products with defined characteristics, at prices that – according to standard economic theory – reflect supply-and-demand induced equilibria. Economic theory accords these prices the role of principal communication media between agents, who use the information prices convey to drive the actions they take in the economy. Relationships between agents do not count for much in "the market". What matters is how each agent separately values each artifact in "the market", values that "the market" then aggregates into prices for these artifacts. Frequently, in popular narratives about developments taking place in business and the economy, "the market" is assigned the central causal role. Indeed, many of the people we interviewed and the documents we read to construct the story we tell here claim that it is "the market" that will determine LonWorks' destiny.

Certainly, "the market" is always lurking behind all of the activities we relate in our story. In the end, the success of LonWorks will be measured by the number of purchases of devices, machines and control systems that use this technology – and the success of Echelon, by the profits it generates from the sales of the LonWorks products and services that it provides. Yet it is striking how relatively minor a role "the market" actually plays in the story we have to tell. It is frequently being upstaged by activities and processes that bear little resemblance to that abstract entity that we described in the preceding paragraph. We will find participants in the emerging LonWorks market system constantly negotiating the *meaning* of artifacts,

McDonald's specified LonWorks for its experimental "Kitchen of the Future" project in 1993. Recently, LonWorks has been adopted by several large European manufacturers of catering and food dispensing systems.

⁵ See Chapter 10 for an account of the Echelon-ENEL joint venture in this area, the largest single application in LonWork history. Over the past two years, Echelon has repositioned itself as an intelligent meter reader company.

⁶ For a discussion of agent-artifact space, see Lane and Maxfield (2005), and Russo (2000).

both those currently in use and those that are being conceived or designed. As a consequence of these negotiations, LonWorks and its associated artifacts *take on value* – not just inside the heads of individual agents, but through social processes that require concrete social settings, which the emerging market system has to provide. Moreover, the agents in the LonWorks market system learn much more from each other than they do from prices, and beyond learning *information* that other agents already know, these agents jointly construct new *interpretations* of themselves and the artifacts around which their activities revolve, interpretations that drive action in hitherto unexplored, even unimagined, directions. Finally, the agents carry out these activities in the context of *relationships*, the most important of which are those we call generative (Lane & Maxfield, 2005).

The relationships between agents are constructed upon *scaffolding* that the emerging market system must provide: *agent-space structures* like standards organizations, users groups and trade fairs; *communication media* like journals, newsletters, and web sites; *rules* that govern both “market” and non-market interactions between the agents and artifacts that participate in the market system; and *shared attributions* about agent *roles* and artifact *functionality*. Without such scaffolding, to speak of “the market” deciding for or against LonWorks makes little sense. And it will not be “the market”, but a complex network of agent interactions that will bring the scaffolding and the system it supports into existence.

How are market systems constructed? The problem we have to consider in our story is not just what a market system *is*, but how it is *constructed* – or, better, *emerges*, since the structure it assumes is frequently quite different from what the agents whose interactions are directed towards bringing it into being intend. In this chapter, we seek to understand some of the difficulties in generating *effective action* for agents that participate in the construction of the market system. To study effective action, we concentrate on the actions and attributions of the people who work for Echelon. We do this for two reasons. First, we have access to considerable information about what Echelon actors did and what they said about why they did it. We have interviewed many people associated with Echelon intensively and repeatedly, from February 1996 through November 1999. In addition, our interviewees have provided us with various internal documents, by means of which we have been able to monitor some processes that leave no trace in the public record. Second, Echelon not only developed the core LonWorks technology, but also has taken a position of leadership in building the market system around LonWorks. Hence, it has played an important role in many, though by no means all, of the processes we seek to understand.

In our discussion of effective action, we highlight three concepts: the role of *generative relationships* in constructing new attributions about the identity of agents and the functionality of artifacts; *redirection*, the process by which an agent changes orientation in mid-course, on the basis of attributions about where current interaction streams are flowing; and *aligning attributions*, processes in which particular agents attempt to bring everyone involved in the market system to share certain attributions about the identity of key system agents or artifacts. Section 9.4 describes an important Echelon redirection episode. Section 9.5 tells about some difficulties

Echelon encountered in establishing generative relationships with other key agents in the LonWorks market system.

Why the LonWorks story? The LonWorks story has two features that make it a particularly interesting case study of the processes through which market systems come into being. First, LonWorks is a network technology, and market systems organized around networked artifacts play an increasingly important role in modern economies. Some of the most interesting features of the LonWorks story involve the construction of the material and cognitive infrastructures that all market systems organized around networked artifacts require. Second, LonWorks is about distributed control, and in following its story we may gain insights into how control is distributed: in artifact space, in the construction of ever more complex local operating networks; and in agent space, in the construction of a market system by the heterogeneous agents that will populate it.

Network infrastructure Like railroads, electric power, telephony, and data communications, LonWorks is a network technology. Like these other network technologies, LonWorks has the potential to change the way many different economic activities and processes interact with one another, and in so doing to generate new kinds of economic processes and functionality. However, this potential can only be fully realized after a great deal of network infrastructure is already in place, and the costs of putting all this infrastructure in place must be borne by someone, in some way.

Unlike railroads, electric power and telephony, distributed control technology does not require a huge upfront investment in physical infrastructure, like railroad tracks, power grids and telephone lines. In fact, the costs of wiring and control modules in a distributed control system may be lower than those for alternative technologies that lack the potential to deliver the extra functionality that distributed control makes possible.

But there are two other types of infrastructure that are necessary for distributed control technology and are very costly to produce. First, there is the core technology on which distributed control depends: communication protocols, integrated circuit designs, network operating systems, application interfaces, media-specific communication hardware, general and application-specific control algorithms, and tools for developing, installing, debugging, monitoring and maintaining control networks. This core technology constitutes infrastructure, in the sense that it provides the essential building blocks for distributed control networks, but unless embedded in such networks it has no value in itself. Someone must bear the cost of developing this technology before it can be instantiated in physical artifacts, which may be combined into control networks that deliver functionality that others may come to value.

The second kind of network infrastructure is cognitive, and its associated costs may be substantial. What allows control to be *distributed* is that there are many different control points – sensors and actuators – in most modern control systems. In principle, at least with LonWorks, control can be distributed down to the level of the individual sensors and actuators. Herein lies one of the greatest advantages of distributed control: the possibility of enabling previously distinct devices, perhaps components of separate systems, to communicate directly with one other and to control directly one another's operations. In this way, whole new levels of system

functionality may be achieved. This possibility can only be realized if all the devices to be combined are already equipped to function as nodes on a distributed control network. The simplest and cheapest way to guarantee this is to internalize the control potential of each device – that is, to design and produce the devices with all the necessary control hardware and software already in place. As we shall see, this is the strategy adopted by Echelon, based upon the Neuron chip.

But here we encounter a serious chicken-and-egg dilemma. To create the potential for their artifacts to generate new kinds of functionality through connection to a distributed control network of other artifacts, producers must be already convinced that this possibility warrants the expense of incorporating the necessary equipment. Clearly, the strength of this conviction depends upon how many other producers of artifacts share it. Thus, the more producers that incorporate the distributed control technology, the easier it becomes for the next producer to decide to incorporate it as well. For the technology to be widely adopted, a kind of *belief infrastructure* has to be constructed: enough producers have to believe that their products' value will be sufficiently enhanced by incorporating the necessary equipment to justify that belief; and then the success of the technology in generating new functionality by combining products produced by the believers must convince enough other producers to equip their product for the technology, until an avalanche of adoptions is generated, and the technology takes off. LonWorks' promoters must figure out how to construct this belief infrastructure – and how much it will cost to do so.

Distributed control in agent-artifact space We shall argue that control of the LonWorks market system is in fact distributed among a number of agents and structures in agent space. In 1988, when Echelon was formed, the company constituted the entirety of the proto-LonWorks market system and hence, tautologically, “controlled” it. Building the market system entailed distributing this control, a process fraught with difficult choices for Echelon executives, with consequences that, we shall argue, would have been impossible to foresee *a priori*. Discovering where the loci of control of the market system lie, and how control is exercised from and among these loci, are recurrent problems in our narrative – for us, as well as for the agents about whom we write. In our story, there will be a continuing counterpoint between distributed control in artifact space, achieved by LonWorks technology – and distributing control in agent space, to construct a functioning LonWorks market system.

9.2 Planning for Control: 1985–1988

Our story begins in the early 80's, with an idea. A.C. “Mike” Markkula, co-founder and then president of Apple Computers, contemplating the decreasing size and cost of computers, asked himself, “What will happen when a computer costs only \$1 and weighs only a few grams?” At that size and price, he thought, most things can have their own embedded computer. But what would these little computers *do*?

Markkula imagined two broad classes of applications. First, by adding intelligence and memory to the things in which the computers are embedded, these things may be able to perform their accustomed functions better and faster. A thermostat, for

example, could be programmed to switch automatically between different heating and cooling modes for different times of day, days of the week, and months of the year. Second, and more important, if the computers in *different* things could *communicate* with one another, the things might in combination provide new functions that neither could provide alone. For example, an occupancy sensor, a light switch and a thermostat might “cooperate” to begin heating or cooling a room as soon as someone enters it – and to turn off the lights after everyone leaves. Or when an alarm clock rings to wake a person up, it might also send a message to a coffee pot to begin brewing that person a cup of his favorite morning beverage. In general, by *combining* different kinds of things, and allowing them to *control* what each other is doing, a *new kind* of thing, delivering a *new kind* of functionality, can be created.

In 1985, Markkula set up a private company, called ACM Research. The company hired a small group of engineers from Apple to begin developing a prototype for a small, embeddable computer and a communications protocol that would allow these computers to talk to one another and control aspects of one another’s operations. By mid-1988, the ACM Research engineers had designed an integrated circuit, called the Neuron, which could be connected together by means of a variety of different media to form a Local Operating Network (LON), which they believed was capable of handling any problem involving “on/level/off” control. At that point, Markkula decided to form a new company, to be called EcheLON,⁷ which would develop the Neuron, the communications protocol and the other technology needed to implement the LON idea, and would then bring LONs to the market.

Markkula and his associates drafted a business plan for Echelon, which they circulated to prospective investors in November, 1988. The plan began by making explicit the analogy between the LON idea and the on-going revolution in data communications, which was then in the process of creating the new workplace paradigm of distributed computing, based on Local Area Networks (LANs). “Personal computers are tools for our intellectual being,” the plan proclaimed. “They extend our ability to think. Echelon networks are tools for our physical being. They extend our ability to do.”

Through LONs, society would have available to it “millions of tiny computers that help us do things.” These “things” fell into three conceptual categories: communication, identification, control. The networks would allow the artifacts in which the computers were embedded to communicate to each other what they sensed of their own states and their local environments. Since each computer would have a unique identification number, any artifact with a computer inside attached to a communication network could be “located” and distinguished from otherwise apparently fungible artifacts. These artifacts could then be armed with appropriate programs stored in their computers’ memories. Finally, the computers could execute control actions with respect to the artifacts to which they were attached – for example, if the artifacts were light switches or motors, by turning them on or off or changing their levels of brightness or speeds of rotation.

⁷ Soon to be simplified to Echelon, as we do hereafter in our narrative.

The core LON technology would consist of four elements: the *Neuron chip*, an integrated circuit with four processing units, shared memory (RAM, ROM and EEPROM) and three input-output ports;⁸ a standardized communications protocol, subsequently called *LonTalk*, by means of which Neurons could formulate and interpret messages for and from other Neurons; *transceivers*, by means of which communication channels in various media (twisted-pair wire, coaxial cable, radio frequency, infrared – even a.c. power lines) could connect Neurons into a network; and a *development tool*, with which networks could be designed and nodes programmed and commissioned. A first essential task of the new company would be to develop this core technology, building on the work already accomplished by ACM Research.

According to the Business Plan, the primary advantage of LON, as compared to existing models of control, was a trio of related concepts: *universality*, *standardization*, and *interoperability*. The artifacts around which control systems were constructed differed markedly among the various control industries of this period: pneumatic and electrical signaling; PLC's, FCU's, and centralized computers; as well as unlimited numbers of "dumb" sensor and actuator devices. The universality of the LON solution provided the way out from this Babel in the world of control: all control systems – whether from the control industries, buried inside machines,⁹ in cars, airplanes, homes or factories – would consist just of networks of sensors and actuators, each equipped with a low-cost Neuron chip.

Universality in turn implied the possibility of *standardization*: if every control system were a LON, they could all speak the same language and share the same basic architectural principles – and even, if it were desirable, be connected onto the same network. Just as Ethernet was rapidly becoming the standard for computer LANs, permitting different kinds of computers to share data, programs and peripheral resources, LonTalk as a universal standard would allow virtually any kind of device to share status and control information with any other.

Finally, standardization implied the possibility of *interoperability*. Devices belonging to different systems, manufactured by different companies, would be able to work together to deliver new functionality. For example, if a building's security system and lighting system were on the same LON, an occupancy sensor from a security system could notify all the light switches in a room when the room was empty, and the switches could be programmed to turn off their lights in response to this message, even if the systems were manufactured and installed by companies that knew nothing of each other's products.

To the authors of the Business Plan, the advantages to be gained from universality, standardization and interoperability were as substantial as they were self-evident. They predicted a very rapid and widespread adoption of LON technology. The markets they expected to penetrate with LON technology included residential

⁸ Which included programmable analog-to-digital and digital-to-analog converters.

⁹ Or standing adjacent to the machine it controlled, as a separate "box," as in computerized numerical control (CNC) for machine tools.

construction, vehicles, retail and distribution, building automation, factory control,¹⁰ aircraft, agriculture and medicine. The most important functions to which the technology would be applied included security, identification, HVAC (that is, heating, ventilation and air conditioning), inventory, illumination, electronic control, and recording and monitoring.

To recruit companies to adopt LON technology, the Business Plan foresaw two stages of activity. First, Echelon would target a small number of important players in the key industries listed above. About a half dozen of these companies would become “alpha” partners, who together with Echelon would develop first-generation LON products, which would “PROVE the technology and showcase real products for real markets.” These partners would benefit from the usual first-mover advantages in their markets. Next, Echelon would carry out a “blitzkrieg” strategy, with a highly public launch, featuring testimonials from their early adopters, followed by seminars for engineers and end-users held around the world. These seminars would be essentially open to all-comers, and it was anticipated that what would emerge, would emerge: better not to try to foresee all the ways in which ingenious designers might use LONs to solve their industries’ problems.

How could Echelon position itself so that it could stay on top of the distributed control wave? The Business Plan raised three problems that Echelon must solve. First, it must take care that LON and not some competitor became the standard distributed control technology. Second, it must guarantee that interoperability could be achieved among all the products from every supplier who adopted the LON standard. Third, Echelon must earn substantial profits from growth in the markets for products that used LON technology.

With respect to the first problem, one of Echelon’s first priorities must be to find *allies* to quickly capture the market for LON and the Neuron. Echelon, after all, had to start small, and yet the battle for LON would have to be fought over a huge swath of agent space. The most important allies – even more important than the early adopters – would be semiconductor manufacturers, to whom Echelon would license the right to produce Neurons. Semiconductor manufacturer partners were important for two reasons. First, they already had the equipment and the competence to produce integrated circuits, neither of which Echelon could hope to achieve without a huge infusion of initial capital. Second, semiconductor manufacturers were large companies, with extensive marketing and sales resources they could devote to LON technology, in which they, as well as Echelon, would have a stake. To make an Echelon alliance more attractive for semiconductor manufacturers, Echelon would charge a very low royalty rate for the Neuron chip. This would also accelerate the rate of adoptions of LON, since the cheaper the chip, the greater the sales – and the greater the sales, the more the semiconductor manufacturers would be able to drive production costs even further down. As a consequence of this policy, though, Echelon would need to find its own profit sources elsewhere than from Neuron sales.

¹⁰ Little was made in the Business Plan of the distinction between discrete and process control, since from the abstract LON point of view, the differences seemed only questions of detail.

The semiconductor partners were also the key to the proposed solution to the problem of guaranteeing interoperability. Every standard, at least in information technology, leaves some design leeway for manufacturers to differentiate their products from competitors'. The problem is to make sure that these differentiating features do not compromise the ability of different products based on the standard to interoperate with one another. The Business Plan's solution to this problem lay in its licensing agreements with semiconductor manufacturers and with original equipment manufacturers (OEMs) – the companies that produced artifacts that incorporated LON technology. The OEM license would require that the artifacts they produce “comply with Echelon specifications and that new features, etc., be approved by the LON Systems Standards Committee (LSSC) before they are released.” The Plan did not spell out the constitution of the LSSC. It was clearly meant to be a standards committee with a mandate broad enough to cover all the industries to which LON technology might be applied. Its task would be “to maintain standards, arbitrate conflicts, and coordinate network issues so that all users are assured of compatibility.”

Of course, without some sort of teeth, merely signing a license agreement could hardly ensure that OEMs would comply with LSSC interoperability standards. The Business Plan proposed to supply these teeth through the license that the semiconductor partners sign with Echelon. This license would require the semiconductor manufacturers to sell Neurons only to customers approved by Echelon. If that provision were enforced, any rogue company producing non-interoperable LON artifacts would be unable to purchase Neurons.

Why should the semiconductor manufacturers allow Echelon to veto potential Neuron customers? Two reasons: the semiconductor manufacturers would be Echelon's partners, tightly bound to them through engineering, marketing and sales co-ventures, and so could be won over to act in what Echelon regarded as the interests of the whole LON community; and, guaranteed interoperability would be in the semiconductor partners' *own* interest, since interoperability was the key to the continuing expansion of LON functionality and hence the market demand for Neurons. In contrast, Echelon could hardly expect to exercise much direct control over the anticipated horde of OEMs, some of whom might well prefer to produce systems that did *not* interoperate and hence whose components could not be substituted with cheaper or better products made by competitors.

Where would Echelon's profits come from? The Business Plan put forward three possibilities. First, Echelon would develop several product lines, consisting of devices and tools for constructing LON networks: the Plan mentioned explicitly development tools and “subassembly modules,” for example Neurons mounted together with transceivers. Second, “at some time in the future, the company may elect to enter one or more of the end markets that have developed around Echelon technology.” Third, Echelon may invest in companies that plan to develop LON applications. In this way, the company can hold a “small percentage stake in the entire LON industry.”

As we have already noted, the Business Plan is a very optimistic document. Its market analysis predicted that Echelon revenues would exceed \$1 billion by 1995. The next step was to make Echelon happen – and start building.

9.3 The Birth of Echelon – Raising Capital, Recruiting People: 1988-June 1989

In mid-1988, Echelon was just an idea, a handful of employees, and designs for a chip and a communication protocol. To build a company that could develop and market control network components and systems, Markkula thought that he needed to raise about \$15 million in capital and recruit a management team. Given his current commitments at Apple and elsewhere, he had no intention of leading the company himself. His choice to head Echelon management was another successful Silicon Valley veteran, M. Kenneth Oshman, founder and ex-CEO of ROLM Corporation, the company that introduced digital switching and computers to telephone PBX's and successfully challenged the American telephonic giant AT&T for leadership in that business (Lane & Maxfield, 1997).

Like Markkula himself, Oshman was frequently described in the popular and business press as a “Silicon Valley legend.” ROLM had been the first Silicon Valley start-up to take on an established, traditional monopolist like AT&T in one of its key product lines and win. Though he had proved to be an astute businessman, Oshman was by training an engineer. A native Texan, he graduated from Rice University, and then earned a Ph.D. in electrical engineering at Stanford before founding ROLM in 1969.

In 1984, IBM acquired ROLM. Less than two years after the sale, Oshman left IBM. While he was serving on the board of several companies, including Sun Microsystems, he had no current full-time commitments in 1988. Markkula approached Oshman before the business plan was completely drafted. The two men came to an agreement: Oshman would purchase a major stake in the company from Markkula and would become CEO of Echelon, while Markkula would serve as Chairman of the Board.¹¹ The two would collaborate on the final revisions of the business plan and would raise the rest of the capital they thought Echelon would need before the company could support itself from its own revenues.

Many a hopeful entrepreneur with an idea and a prototype faces the step of raising capital, necessary though it may be, with considerable trepidation. After all, most venture capitalists turn down many more proposals than they accept, and even when they do provide funds, the entrepreneur must as a consequence accept a diminished share in whatever financial rewards his ideas eventually generate and in the control over the company he has set up to bring his ideas to market. More than one entrepreneur has even found himself ousted as head of his own company by a Board of Directors controlled – or at least strongly influenced – by representatives

¹¹ Within a year, Oshman became Chairman and Markkula Vice-Chairman of the Echelon Board.

of his venture capitalists, anxious to salvage their investments with a timely and profitable exit strategy. Even when the entrepreneur keeps his position, he may have to abandon a cherished development strategy if it comes into conflict with the venture capitalist's exit strategy for the company.

Markkula and Oshman had no such problems. Between them, they had close relationships, both personal and professional, formed over many years, with most of the leading Silicon Valley venture capitalists. Based on their respective past business successes and the intrinsic merits of their project, they neither foresaw nor encountered any difficulty in obtaining the amount they sought. In fact, so as not to disappoint some of their prospective investors, they actually raised more capital, around \$25 million, than they had originally intended.

Several venture capitalists were allowed to invest in the new company. Among them were Kleiner Perkins Caufield and Byers, probably the premier Silicon Valley venture capital firm; Arthur Rock, a pioneer Silicon Valley venture capitalist, who had been the lead investor in all three of the start-ups for whom Markkula had previously worked – Fairchild, Intel and Apple; and Venrock Capital Partners, the Rockefeller brothers' venture capital arm, which had important connections among financial and industrial concerns outside the Valley that would later prove valuable to Echelon. In fact, in 1994, Peter Crisp of Venrock, who was the one representative of the venture capitalists that Oshman and Markkula placed on the Echelon Board, was instrumental in obtaining an additional \$10 million in capital from George Soros' Quantum Fund, at a time at which both the capital and the prestige associated with its source were important for Echelon's continuing development.¹² In addition to the venture capitalists, Markkula and Oshman gave a few Silicon Valley friends the opportunity to invest in Echelon.

By mid-1989, the members of Echelon's Board were Markkula, Oshman, Rock, Crisp and Robert Maxfield (a ROLM co-founder and Oshman's closest collaborator there). Larry Sonsini, head of the powerful Silicon Valley law firm Wilson Sonsini Goodrich and Rosati, joined the Board in the spring of 1993. As in most Silicon Valley start-up companies, the Board played an important role in Echelon's development. Rock, Markkula, Oshman, Maxfield and Sonsini were all linked to the network of successful Silicon Valley entrepreneurs, and thus they were connected to leading executives in most of the important Valley semiconductor and computer companies. These connections frequently vastly amplified Echelon's actual "market power", and provided the still-small and profit-less start-up access to key people in these companies, through which negotiations could be undertaken and, occasionally, partnerships attained. For example, in 1993 Rock arranged a meeting between Oshman and the CEO of PG&E, whom Rock had interested in LonWorks technology. As a member of Intel's Board of Directors, Rock also helped to catalyze a joint Intel-Echelon demonstration project in home automation for the Comdex trade show in 1994. Sonsini helped initiate discussions between Echelon and Packard Bell over

¹² In 1995, Crisp also helped convince the president of Otis Elevators to adopt LonWorks technology.

home automation strategies in 1995. Oshman, a member of Sun's Board, initiated several cooperative projects between Echelon and Sun and various Sun initiatives, like GINI and JavaSoft. In addition, Sun CEO Scott McNealy presented a rousing keynote address at the May 1996 LonUsers show, in which he launched the concept of JavaLon, the union of LonWorks control networks with TCP/IP-based communication networks. Later, this idea was further developed in partnership with Cisco, the leading manufacturer of hardware for the internet. The link with Cisco was first established through an ex-ROLM employee, Steve Behm, Cisco Vice President for Global Business Alliances. On several occasions, Behm helped push along temporarily stalled negotiations involving Echelon and Cisco employees lower down than he in the company hierarchy. A final example: in 1995, lobbyists of the EIA (Electronic Industries Association), a powerful trade association, attempted to insert a clause into the massive federal bill de-regulating the telecommunications industry that would have effectively mandated a competitive standard (CEBus) for home automation. Echelon was able, with the active participation of its Board members, to put together very quickly a powerful cross-industry coalition, including Apple, Intel, Sun, Stratacom,¹³ Detroit Edison,¹⁴ Scientific Atlanta, Motorola¹⁵ and others, which blocked the initiative – an action unthinkable for other companies of Echelon's small size (about 120 employees in 1995) and relatively minor market presence (1995 annual revenues about \$20 million).

For Board members to exercise their personal networks in behalf of Echelon, they had to be kept informed about Echelon activities, projects and strategies on a regular basis. Oshman sent a letter to Board members every month, in which he summarized recent performance and prospects for the company and occasionally raised deeper strategic questions that he wanted to discuss at the next Board meeting. Moreover, from the company's inception, the Board met regularly, usually every other month. At these meetings, they heard reports from Oshman as well as other key executives, in particular the Vice-Presidents of Marketing, Engineering and Finance. In addition, Oshman frequently solicited their opinions, advice and concurrence, especially on financial and key strategic questions. So Board members knew a lot about Echelon's activities, projects and strategies, at least from the point of view of the company's top management. But because of their ties of friendship and collaboration with Oshman, much preceding their work together at Echelon, and their genuine deep respect for his intelligence and management ability, the Board rarely challenged Oshman's ultimate authority in the company: they were his friends, advisers, helpers, not impersonal skeptical overseers.

After Oshman became CEO in November 1988, he started hiring his new management and engineering teams. ACM Research had been dominated by ex-Apple employees; Echelon quickly took on a new, ROLM flavor. The Vice President

¹³ Stratacom's CEO had been a ROLM Vice President, who later accepted a seat on Echelon's Board of Directors.

¹⁴ Detroit Edison had invested \$10 Million in Echelon.

¹⁵ Motorola was a semiconductor partner and investor in Echelon, with a seat on the Echelon Board of Directors.

for Finance, Oliver (“Chris”) Stanfield, had worked for ROLM, as had two of the four people that comprised Echelon’s marketing department in 1989. But the real focus of new hires, and of ROLM influence, was in Engineering. By the beginning of March 1989, Echelon had twenty-seven full-time and two part-time engineers on staff. The four key engineers – the Vice-President of Engineering, the System Architect, and the Directors of Hardware and Software Development – had all been ROLM employees. Every other engineer, many of them also ex-ROLMans, reported to one of these four men. The ex-Apple engineer who had directed engineering for ACM Research stayed on awhile as head of project planning, but he left Echelon in 1990. Several others from the ACM Research days stayed longer. In 1996, a recently hired new Vice President of Engineering (another ex-ROLMan) reported to us that there was still some resentment among these people about alleged “favoritism” granted to ex-ROLMans. Favored or no, the ROLM engineers who took over Echelon Engineering had excellent academic backgrounds, a record of superb engineering accomplishments at ROLM and later ports of call – and a measure of financial security, obtained from their ROLM stock options. It was a group accustomed to one another and to success.

One thing, though, that the group lacked was any prior experience in the techniques and problems of bread-and-butter control engineering, as practiced in the existing control industries. Not a single one of the new Echelon employees was hired from these industries. This pattern endured. As it expanded, Echelon continued to hire from Silicon Valley. In marketing as well as engineering, Echelon sought smart engineers who knew a lot about computers, digital logic and analog electronics. Given that most engineers who worked for the mainly Midwestern-based control companies had neither the academic backgrounds, computer experience nor Valley style that Echelon sought, they would have been unlikely to get jobs at Echelon – or to prosper, had they somehow been transplanted there. Not until fall, 1994 did Echelon hire its first employee with a controls industry background (see Lane & Maxfield, 2005 for the story of how he landed there and what happened as a result).

In this respect, Echelon was carrying on a Silicon Valley tradition, and specifically following ROLM’s example. ROLM started as a militarized computer company. Its founders understood computers and what one could do with digital logic. When they looked about for a business into which they could expand, they were intrigued by the possibilities that digital switching and computer control might open up for telephony. But they did not know very much about telephones. Even when they hired an engineer to design their first PBX, they recruited Jim Kasson from Hewlett Packard, the Silicon Valley digitizing company par excellence, rather than seeking someone from AT&T, Western Electric, or Northern Telecom, who knew something about PBX’s and what they did. In the Valley in general, and ROLM and, later, Echelon in particular, there was the sense that if you had a group of smart, technologically avant-garde engineers, you could solve any problem. Too much “application area” experience often turned out to be more a hindrance to creativity than an advantage in itself.

The first important task facing the new Echelon engineers was to review the existing designs for the Neuron chip and the LonTalk protocols. At the same time,

the small marketing group began making contact with potential early adopters in a variety of industries, sounding them out about their interest in applying LonWorks networks in their products. Through the first months of 1989, a picture began to emerge, considerably at variance with the rosy vista of the Business Plan.

9.4 Crisis and Redirection: June 1989–1992

By June 1989, the engineering and marketing reviews were nearly complete, with discouraging results. The logic of universality, standardization, and interoperability, leading to steadily decreasing prices for Neurons and increasing numbers of new markets for LonWorks, no longer seemed compelling to the Echelon management team.

To begin with, different applications seemed to have very different control requirements. For example, HVAC (heating, ventilation and air conditioning) control algorithms appeared to require more memory than the current Neuron design could provide, while many industrial control applications seemed simply too complicated to be configured in Neuron-based networks. In fact, potential customers from the process control industry seemed interested in the Neuron only as a communications device, for which they would be willing to pay only \$1–\$2 per node. Moreover, at the other extreme of application complexity, most embedded control problems could be solved more cheaply with application specific integrated circuits, and so the additional but unused capabilities of the Neuron would price it out of this market. And there were other, more specific problems with respect to particular applications. For example, home security companies were less interested in the possibility of integration of security systems with other home automation systems than they were in the costs of power. Since security systems are always on, customers wanted them to consume as little power as possible, and the Neuron's power requirements were several times higher than the current industry norm.

Some of these problems were tied to specific Neuron design decisions that could be corrected in the design revision that now seemed inevitable. However, the deepest problem seemed not to be about the Neuron design, but with the very concept of distributed control networks as a universal control solution. Distributed control networks might just be too complicated to be economical for embedded control, while a hierarchical control architecture might actually function better than a flat network architecture for controlling complex systems. If this were indeed the case, the potential market for distributed control networks had still to be defined, and would certainly be much smaller than envisioned in the Business Plan.

As if all this were not enough, the Echelon marketing people had encountered two other unanticipated difficulties in their interviews with potential customers. First, they found that initial enthusiasm about the possible advantages of distributed control networks all too frequently dissolved into skepticism as soon as the discussion turned to implementation details for particular product lines. Second, when the potential customer's technical staff was brought into the discussion, symptoms of the dreaded NIH syndrome ("not invented here") appeared, and the temperature

would begin to cool considerably. Forging partnerships with customers to develop LonWorks-based products was beginning to appear substantially more difficult than it had originally seemed to the authors of the Business Plan.

All these problems meant serious trouble for Echelon. First, it was now clear that the time it would take to develop the core technology would be substantially longer than originally anticipated, since the Neuron and LonTalk would have to be redesigned from the ground up. Second, while a huge potential market might still exist for distributed control networks, it was no longer obvious that penetrating the *existing* control industries – especially industrial controls, HVAC and embedded controls – was feasible. Thus, the work of building wholly new control markets would have to assume a higher and earlier priority than had been envisioned previously. So the problem was not just what to do and when to do it, but how to pay for it as well, since the hoped-for revenue streams from Echelon products now appeared shallower and farther off in the future than they had six months before.

The Echelon management team outlined four possible courses of action to the Board. The first option was to go back to the drawing boards and invent a new solution. Second, Echelon could select and then enter a few niche system businesses for which the present solution seemed most promising. Third, the company could try in any case to create a “networking business for control systems,” with a restricted definition for the target control problems to which such networks might be applied. Fourth, the Board could decide to fold the company and return to the investors what remained of their investments.

The management team advised against folding the company. The idea of distributed control networks still seemed too promising to abandon. Moreover, semiconductor partner Motorola was becoming very enthusiastic about LonWorks’ future and was very eager to press forward, eager enough to agree to help fund Echelon research and development costs over the next three years in the form of licensing and prepaid royalty fees.

The Echelon management team’s proposal to the Board combined elements of each of the other options they had outlined. First, Echelon should narrow its focus to applications that “involve simple on-off, message-based control, and possibly ID applications as well.” Within this focus, it could then target a small number of specific markets and evaluate the possibility of creating a system business in each of them. Then, “after we are successful in establishing one or more revenue producing businesses, we will [re]examine the possibility of developing and promoting a general solution to our targeted subset of applications, and of making the establishment of this general solution a good business opportunity for Echelon.”

The management team offered examples of niche system businesses within the proposed focus area, all of which Echelon marketing people were currently investigating. These examples included efficient diagnostic services for LANs, shelf tags for retail and warehousing applications, load-shedding devices to redistribute demand for electricity and thus lower the cost of power, commercial lighting, home automation systems (to compete with X-10, the current leading network technology), transceivers, smart office furniture, a home controller with a television interface accessed via touch-tone telephone, and traffic monitoring. Echelon

management emphasized two particularly attractive features about entering one or more of these businesses as quickly as possible. First, it was the quickest path to a revenue stream that could pay for developing more general LonWorks technology. Second, wrestling with the real problems of an application might help Echelon learn how to define better just which features that general technology ought to have. For both of these reasons, finding the right niche system businesses and plunging into them could boot-strap Echelon from its present difficulties to a position from which it could regenerate its quests for standardization and interoperability.

At the same time, though, the management team advocated redesigning the core LON technology and, with better designs in hand, continuing to pursue OEM business opportunities.¹⁶ The way to manage the redesign, they argued, was to “optimize the protocol, neurons, development system and transceivers [that is, the core LON technology] for communications, simple control, and identification” applications. From this optimized core technology, Echelon could begin to produce “generalized products for OEMs.” Such products might include modules that integrated neurons and transceivers in a form factor suitable for a particular application environment, as well as network management software tools.

At this point, a redirection was certainly underway at Echelon, but it was much too early to speak of any strategy shift. Rather, the management was proposing several parallel lines of search-through-action, with a rough ordering on which lines to pursue first. In fact, several projects were already underway, each aiming to find a partner to explore a possible LonWorks-related business opportunity. Meanwhile, several different groups of Echelon engineers were simultaneously fixing problems in the existing neuron design and rethinking the neuron architecture and the LonTalk protocol.

In addition, another process was underway at Echelon. For the management team’s proposal to the Board represented not merely a redirection of Echelon’s *activities* but of its emerging attribution of its own *identity*: its collective sense of “who” Echelon was, what it did and how it did it. In particular, the critical review that Echelon had undertaken challenged two fundamental assumptions about its identity as set forth in the Business Plan. First, Echelon was abandoning its claim that LonWorks represented a universal solution to the problem of control. Second, the company was reversing its temporal business priorities from the order announced in the Business Plan: it would seek *first* to enter specific “niche” businesses, and *then* to develop a general control networking business. As such, Echelon would not be, first and foremost, a “technology provider,” the primary role the Business Plan had assigned to it.

Of course, agent identities do not change quickly. In particular, while they may be partially *shaped* by documents like the management team’s proposals to the Board, they are hardly *established* by them. They emerge from histories of interactions and

¹⁶ Both the revised LonTalk protocol and the Neuron 2 design specifications were completed by November 1989. First silicon versions of the Neuron 2 were completed late in 1990. The revised Neuron had three instead of 4 processors, two to handle communication and one for applications; RAM, ROM and EEPROM memory; and 11 I/O ports.

attempts by the participants in those interactions *collectively* to make sense out of what has happened to them, providing patterns to induce an order out of the ebb and flow of contingency.

In fact, neither of the attributional shifts about Echelon's identity described above stabilized within the company. Paradoxically, they were undone in the course of the very activities that were initiated as deliberate attempts to enact them. To see how this happened, we focus for the rest of this section on two of the processes set in motion by the crisis and ensuing redirection: the search for a *system business* to enter; and the emergence of an *identity* that made sense of, and provided a future orientation to, the activities in which Echelon found itself engaged in the course of its redirection.

Searching for a system business: the Smart Office Furniture story In the summer of 1989, the Echelon marketing group tried to find partners for several system business development possibilities. The guiding assumption was that the partner would pay a substantial portion of the development costs and would do the lion's share of the marketing, installation and maintenance services associated with the resulting product. Only one of these efforts survived as a viable action option at the end of 1989. The product was to be "smart office furniture", and the partner SCI,¹⁷ the office furniture market leader.

RK, one of the new hires in Echelon's marketing group, had conceived and developed the idea of smart office furniture in the context of conversations with a lock manufacture that eventually decided against adopting LonWorks for its own products. RK noted that competition among the leading companies that manufactured and installed office furniture was intense and that the companies' dealers were eager to find ways to differentiate their products from their competitors'. RK's idea was to target the "work area" sub-market, in which office furniture companies produced (and their distributors sold and installed) an entire suite of furnishings, from electric outlets to desks to file cabinets, for the work sites occupied by individual white-collar workers. A LonWorks network connecting such sites might allow each worker more individual control over his own working environment, while at the same time permitting monitoring and control services that could deliver more efficient management of the facility in which these sites were located. RK found a receptive audience for his idea among product development people at SCI's Grand Rapids MI headquarters.

By September, RK and his SCI allies had put together a preliminary plan for the smart office furniture project. Electric outlets that were also nodes on a LonWorks network would perform such services as turning off lights when employees left their work areas, as well as continuously monitor energy usage and load. Smart locks for desks and filing cases would help with security management. A Neuron in each article of furniture would simplify asset management. Finally, each worker could use the LonWorks network to control the lighting intensity and temperature in his own workspace.

¹⁷ We have changed the name of this company in our narrative.

Through interviews with a number of SCI employees, dealers and customers, the plan's concepts were tested and refined. By November, RK was able to assert that customers would be willing to pay for smart outlets if they could be assured of a two-year payback period. Given that one large SCI customer estimated that merely turning off their graphic workstations over the weekend would save \$400 per station per year, a short payback period seemed likely. Customers who had defense contracts were willing to pay for smart power locks, so long as they met Department of Defense specifications. SCI dealers were enthusiastic about the project, since it would provide an important product differentiator that would make their selling job easier. Moreover, SCI management was anxious to adopt at least the smart outlet idea, since the cabling scheme they had been implementing had been recently found to be in violation of another company's patents. As a result, Echelon managers had already assigned three engineers to begin working on hardware and software for the outlet nodes and developing a personal computer user interface for the system.

By May, 1990, negotiations with SCI had progressed far enough to plan to assign five engineers, full time, to the smart office furniture project. Late in June, Chris Stanfield, Echelon's lead negotiator, and SCI Executive Vice President agreed on the terms of a contract between the two companies. According to their agreement, SCI would pay Echelon \$3.4 million to design hardware and software for the facility management system. In addition, Echelon would be the sole-source supplier of system components. In return, Echelon agreed not to sell a similar system to any other furniture manufacturer in six countries for four years. Oshman wrote to the Echelon Board, "This has taken a long time but I believe it is well worth the effort. We will get experience with a real application. We will begin the process of learning how to manufacture modules. We may create a \$30-50 million profitable business."

Unfortunately for Echelon, SCI's top management team was divided over the Echelon alliance. After two weeks of internal discussion, SCI's CEO, his Vice President for Engineering, the SCI lawyer, and the Executive VP called Oshman to tell him that they had decided not to go through with the project. SCI's Executive VP, however, still believed in the concept, and after another several weeks he managed to convince the others to give the project a second chance. In August, SCI's Executive VP, CEO, and Vice President for Engineering visited Echelon for discussions and to view the demos Echelon engineers had prepared. The visit was a success: on September 15, Oshman announced that the contract with SCI was signed. "Now we get the opportunity to staff the project and really wow SCI with the results. . ."

But the opposition inside SCI was not interested in being wowed. Within a few months, it again had the upper hand, and at the end of February 1991, the SCI CEO notified Oshman that SCI was discontinuing its support for the LonWorks Facility Management System.

We do not know what actually happened within SCI to undo the smart office furniture alliance. An Echelon manager told us in February, 1996, that opposition to the project was led by an engineering manager, who was unwilling to commit to a new direction for a major SCI product line that depended on a technology from the outside – that is, "NIH." Another Echelon manager simply referred to "SCI

politics.” Of course, it is quite possible that some SCI executives were simply not convinced that a compelling business case could be made for the project, or that the LonWorks technology was really sufficiently well developed to deliver the necessary functionality for a price that SCI was willing to pay.

What we are certain about, though, was the effect that the cancellation had on some key Echelon managers. For them, the episode provided vivid proof of the danger of concentrating a lot of Echelon resources on developing a system business. Too much depended on the good will and intentions of the partners, whose business territory Echelon would be entering. The partners would have all the connections and all the competences connected with promoting, selling, distributing, installing, and maintaining the systems they produced. Inevitably, it would be the partner who would control how the system that Echelon developed would fare in the marketplace. The temptation for the partner to do it “their way” would be strong. In particular if the partner lacked sufficient networking expertise to realize just how important and how difficult would be the contribution that Echelon could bring to the project, Echelon would have a hard time obtaining its fair share of revenues from the partnership. So the SCI story endured in Echelon’s collective memory as a signpost pointing the company in another direction: away from developing system businesses, back to being a “technology provider.” All the possible advantages that had induced the management team to propose the re-ordering of company priorities in June 1989, were forgotten, or at least heavily discounted. After 1991, no-one in Echelon was looking for a system business to enter any more – at least until another, different kind of crisis, led Echelon again, temporarily, back to that particular road (see Lane & Maxfield, 2005).

Searching for an identity During the three years following the June 1989 crisis, the Echelon management team spent substantial time and effort trying to define just “who” Echelon was. One sign of this introspective activity was the plethora of explicit strategy and mission statements that were produced in this period. By the end of 1992, this spurt of identity-formation-by-formulation had run its course. After that, mission and strategy would be *enacted* more than they would be discussed and transcribed. Here, we trace three elements of the Echelon identity that emerged during these critical three years – and contrast them with what the Business Plan had envisioned and with the immediate response to the 1989 crisis.

The most striking change was the abandonment of the idea to enter system businesses. This change reflected more than the frustration associated with the rocky course and ultimate failure of the SCI partnership. The new solutions that chief architect Bob Dolin and his group had developed for the Neuron, the LonTalk protocol and network management services renewed confidence in the generality of LonWorks technology. As a result, the 1989 restriction to “simple, on-off applications” disappeared. By May 1990, Oshman was writing to his fellow investors that Echelon’s mission was “to create the framework, plus tools and components, for low cost, reliable, distributed sense and control applications” – without the modifier “simple.” The mission statement prepared by the Echelon management team in January 1991 went even further: “to establish the *de facto* worldwide standard for intelligent, distributed control.”

The strategy document that accompanied the January 1991 mission statement, though, introduced another kind of restriction: no longer on the relative *complexity* of the control problems to be solved, but on the identity of the *market system* to which the relevant applications belonged. The proposal was to focus especially on applications for buildings, industry, machines and the home – that is, the “traditional” control industries (building automation and industrial control), plus embedded control and home automation. These were the best-organized market systems with respect to control artifacts of the eleven potential LonWorks application areas mentioned in the Business Plan. The process of restriction-by-market-system continued. Later in January 1991, a business case prepared by Echelon’s marketing department declared that the primary areas of marketing activity would be buildings, industry and machines, while merely staying “active” in home automation, the least well-organized of the market systems still on Echelon’s focus list. In May 1991, Bea Yormark, Echelon’s Vice President for Marketing,¹⁸ informed the Board that her group was ready to “make the transition from broad seeding in the marketplace to a combination of broad and targeted selling,” where the targets, for now, would be “intelligent buildings” and “factory automation.”¹⁹ That is, Echelon would concentrate on penetrating the most established control market systems and downplay efforts to create a whole new control market system organized around distributed control networks. Note that Echelon’s principal focus would now be the very application areas that in June 1989 were judged to be too complicated to solve with Neuron-based distributed networks. Another key element in Echelon’s emerging identity was a new orientation to standards. The Business Plan had acknowledged the necessity of dealing with standard-setting committees that were developing control network technologies for particular market systems. These committees were initially not receptive to substituting LonWorks for their own proto-solutions. Indeed, Echelon came to regard the paper standards emerging from these committees as LonWorks’ primary competition. The novel idea that emerged in 1991 was that LonWorks might become a standard just by *declaring* itself a standard, frequently and vociferously enough to convince other companies to believe it – and then to adopt LonWorks for their products. “If we claim it, and others agree, we’ve won.” The January 1992 strategy statement prominently featured a new marketing slogan that succinctly summed up this approach to becoming a standard: “Win the air wars!”

Of course, Echelon would have to back up its claim to be “the” standard. Yormark proposed a simple strategy to accomplish this task: it might be enough, she reported to the Board, just to set up a “LonMark certification committee,” and then wait for

¹⁸ Since January 1990.

¹⁹ There was another reason for dropping embedded control from the list. The most important embedded control application area was automobiles, which had represented 25% of the total available control market estimated in the 1988 Business Plan. By December 1992, though, it was clear to Echelon management that LonWorks could not penetrate the automobile control market, which had already converged to an alternative protocol, CAN, developed within the automobile industry itself (MKO Board letter, December 1992).

prominent Echelon customers to bring maybe four or five highly visible, certified products to market. Actually, it would take much more than this (see Lane and Maxfield, 2005, for the story of ConMark International). But the insight that being *perceived* as a standard is essentially the same thing as *being* a standard is a powerful idea. Echelon's internalization of this idea in 1991–1992 helped to shape the way in which it fought standards wars in the ensuing years – and probably contributed substantially to some of the standards battles that it won.

During 1989–1992, Echelon made a significant change to its attribution about LonWorks' "core technology." Before 1990, "core technology" meant the tools and instruments necessary for design and *operational* control of a network: the Neuron, the LonTalk protocol, transceivers, and the LonBuilder development tool. In a January, 1991 strategy document, an addition to this list appears: network management tools. These tools would allow network supervisors to install the network, to monitor its performance, and to make changes when desired. Network management tools require a *system* view of the control network, which includes the *human beings* who interact with the network and who exercise *supervisory* control over it. As the product marketing group explained to the Board in November, 1991, the design of network management tools must take into account who will be using them, including field installers and repairmen, system integrators and end-users. Thus, these tools must be portable, easy-to-use, and with a friendly user interface. By March 1993, network management was perceived not only as an essential part of the LonWorks core technology, but as the key to "positioning for a bigger market." The "*control and monitoring market is much bigger than installation.*" It took quite some time before Echelon network management artifacts managed to embody the attributions that emerged in 1989–1992, but a key step had already been taken in this period, when Echelon actors began to understand that network management and network managers constituted essential parts of LonWork systems.

9.5 Who's a Customer, and How do we Relate to Them?

On December 5, 1990, "the day everything began to work together,"²⁰ Echelon introduced Local Operating Networks – LONs – in a public relations bash at the Equitable Center in New York. Over 300 business leaders and journalists attended the event. Markkula, Oshman, and Yormark explained how Lon technology worked and suggested some of the ways in which it might transform homes, buildings and factories. Spokesmen for Motorola and Toshiba lauded the new technology and predicted that their companies' Neuron chips would be available by mid-1991, with initial prices under \$10. But the real highlight of the carefully orchestrated presentation were videos in which executives from nine potential LonWorks customers²¹

²⁰ Title of the video Echelon produced of the launch.

²¹ Lighting systems manufacturers Advanced Transformers and Lithonia Lighting; switch and lighting fixture producer Leviton Manufacturing; factory controls market leader Allen-Bradley;

described the critical challenges their companies faced and explained how the LON concept might contribute to meeting these challenges in the future. Oshman emphasized the differences in the roles that Echelon and such companies would play in building the LonWorks market system: “It is not Echelon that will be making smart products or systems. Echelon’s business is to *support* the companies that make these products and systems. Echelon’s job is to do the research and to provide easy-to-use, low-cost enabling control technology, in the form of tools and components that other companies can use. . . *It’s the customer who’s the innovator.* We expect it will be our customers, our adopters in many industries, who will show us the different things you can do with a LON – and I’m sure many of those will be things we can’t even imagine today.”

Echelon executives were euphoric after the launch. As Oshman described it to the Board two weeks later, “Our early adopters were impressed. . . The audience understood our message and went away convinced we were real. . . We are swamped trying to sort out good leads. Our two salesmen are in heaven.” The articles that appeared in the technical press serving the various control industries were in general enthusiastic. The business press, though, had a more tempered reaction. In particular, Business Week was skeptical: “There’s a good chance that Echelon’s low-cost, do-it-all strategy will fall flat. Some network experts question whether its technology can be economical and effective for both a light socket and an industrial robot.” The article went on to claim that even Echelon’s early adopters were not “champing at the bit.” Johnson Controls, one of the companies represented in the launch videos, “says it has no firm plans for using Neurons in its building controls.” The article implied that the relationships between Echelon and LonWorks “early adopters” were a bit more ambiguous than the upbeat, partner-like impression that Echelon had tried to convey to the launch audience. In this section, we will trace the evolution of these relationships from 1989 to 1993.

First contacts Of course, relationships between Echelon and potential LonWorks adopters began well before December 5, 1990. Echelon had been organized two years before the launch, and by March, 1989, the new hires in engineering and marketing had already made contact with a number of potential customers. For Echelon, the main aim of these initial contacts was to evaluate where and how LONs based upon the Neuron chip design inherited from ACM Research might be applied. The companies with whom Echelon representatives talked came from a variety of different industries: home security systems, manufacturing control, process control, remote meter reading, network diagnostics, office furniture, home automation, building automation, and agricultural automation.²²

military system integrator CACI; Johnson Controls and Landis & Gyr, two of the three largest producers of control systems for building automation; the leading office furniture manufacturer, Steelcase; and Ziotech, a producer of industrial and board-level computer products.

²² Respectively: Radionics and Unity Systems; Allen-Bradley, the largest producer of discrete control devices and systems, and Square D, another large player in the manufacturing controls market system; Honeywell and Accurex; Itron and Metrocom; 3-Com; Intelock; Leviton Manufacturing,

The Echelon marketing group whose task it was to understand this heterogeneous collection of companies and businesses numbered only four people in April; two more people came on board in July. The understanding that these people obtained about “typical” applications in all these areas were discussed with members of Echelon’s engineering group, who then tried to design similar applications using Neuron-based networks. As we saw in Section 9.4, the results of this process were so discouraging that Echelon’s management team seriously considered shutting the company down.

However, not everyone inside Echelon was convinced that the discussions had yielded a sufficiently deep understanding of Echelon’s potential customers to justify such pessimistic conclusions. RK, one of the most active participants in the process, expressed such a view in a June 8 memo circulated to key Echelon managers:

“To date, our interactions with customers have centered around delivering the LON vision and attempting to gain customer buy-in with this vision. We postulate that since all control problems reduce to “on, off, or in-between,” a general solution exists; and then we assert that ours is the right general solution. The principal benefit of this process is that it tends to quickly reveal those assertions which potential customers find most offensive or which potential suppliers find most limiting. However, there is no reason to suspect that the results of this “assertion/refutation” process are the same as those which result from, say, studying a number of control systems in great detail, and then attempting to generalize from the data. It should thus come as no surprise that we don’t have a qualified understanding about what the LON System does well: we’ve never implemented a process which yields such information as its result. Instead, we’ve implemented a process that tends to tell us what we can’t do particularly well. While this information is extremely useful it tends to have a notably negative ring to it, and we tend to find ourselves on the defensive more often than we’d like. . . It is very difficult to determine the applicability of the LON System to an entire class of applications, all at once. One needs tremendous insight and experience in a field to be qualified to make such broad judgments, and even then, the validity of the conclusions is highly sensitive to inaccuracies in the simplifying assumptions.”

RK went on to suggest alternate processes, designed to “gather good data about a specific application:”

1. hire a consultant “who is an expert in the area and, after fully explaining our concept, charge him with developing specific solutions to particular industry problems using our system;”
2. “approach customers with a specific proposal to investigate a LON-based solution to a particular problem,” and then work closely with the customer “to learn about their industry and their concept of how the LON System is used and the benefits which result;

which not only manufactured switches and other lighting devices but also held a license to sell X-10 control equipment, the leading current home automation technology; and Priva.

3. propose to do a pilot development project with a customer to demonstrate “key aspects of our solution” and test and validate key assumptions.

RK concluded that all these solutions had a “common theme:” working closely with an expert or “real customers” towards a “well-defined target” to “flush out the key issues.” In this way, “we will gather the information we need to fully assess an application, we’ll identify points that may lie within the boundary of a general solution, and we just may learn enough to build a good systems business along the way.”

RK’s memo criticized Echelon’s way of relating to its potential customers. In RK’s view, Echelon representatives were proselytizing the gospel of distributed control, presenting Neuron-based LONs as a revolutionary universal solution to all control problems. But to RK, Echelon people did not yet understand in sufficient detail what concrete control problems potential customers solved, and hence the reasons behind the artifactual forms that embodied their solutions. As a result, Echelon was not yet in position to understand how Neuron-based control networks might contribute to these companies’ business. To RK, this understanding could be gained only through intense, temporally extended relationships, either with a consultant steeped in a particular control industry or directly with one or more customer-partners.

To understand how such relationships between Echelon and customer partners might have been constructed, it will be useful to analyze the five components of what Lane and Maxfield (1997) called *generative potential* for interfirm relationships:

1. *Aligned directedness*: the participants in the relationship need to be oriented towards similar transformations in the structure of agent-artifact space. We might assume that Echelon and potential partners would share an orientation toward developing, together, a new control artifact or a better solution to an existing control problem. But there might be obstacles in determining just which kind of artifact or what would constitute a control solution. Here, issues of cost, of intellectual property, of the importance of interoperability with other artifacts and systems, perhaps developed by the partners’ competitors, on attributions of artifact or system functionality might provide stumbling blocks to alignment.
2. *Attributional heterogeneity*: the participants need to have different ways of looking at the problem to be solved if the relationship is to generate solutions that neither could provide alone. Between Echelon and its potential partners, there was certainly considerable attributional heterogeneity. Echelon’s abstract concept that “the network is the controller” and the deep understanding that some of Echelon’s engineers had of communication protocols were certainly not yet shared or even grasped by people in control businesses. The regnant vision of controls inside control businesses was quite concrete, embracing not only intimate working knowledge of myriads of sensors and actuators and hydraulic, pneumatic, electrical, analog and digital control signaling technologies, but also an understanding of the expectations and desires of their customers, the users

of control systems²³ – all of which were outside the experience of Echelon’s team of Silicon Valley engineers and marketeers. The real problem for Echelon and its potential partners was to develop modes of discourse through which their attributional heterogeneity could surface and find expression in a new language that they jointly developed.

3. *Mutual directedness*: the participants need to have a positive orientation towards one another, feelings at least of joint comfort if not of trust. This requirement was particularly problematic for relationships between a small Silicon Valley start-up and large Midwestern control companies. Could they learn to appreciate each others’ better qualities?
4. *Permissions*: groups of people from each participant with the requisite interests and competences to forge working relationships must be allowed to engage in the kinds of interactions – both talk and joint action – that can lead to common and then new understandings and embed these understandings in new artifacts and system solutions. In principle, the engineers and marketing people from Echelon and potential partners could be granted the necessary permissions; in practice, problems often arose, from the delegation of people lacking necessary meshing competences and knowledge to interact, through the prohibition on the part of the partners to solutions “open” enough for Echelon, to disagreements about who must bear the cost of necessary research and engineering.
5. *Action opportunities*: a relationship cannot be generative on talk alone. At a certain point, concrete projects for new artifacts or systems must emerge. Frequently, relationships between Echelon and potential partners did lead to action opportunities, but only at the demonstration level.

It is interesting to note that these conditions for generativity are satisfied very well for the Silicon Valley insiders who initially directed the construction of Echelon, its financing and its engineering successes. RK’s memo, and the story of the relationship with SCI, indicates that the early relationship with potential customers was anything but generative. In fact, too much heterogeneity and too little aligned and mutual directedness seemed to characterize many of these relationships. We can interpret RK’s memo as a suggestion that Echelon re-orient its relationships with potential customers, away from a mode in which verbal interactions were essentially acts of proselytizing, to a mode of on-going co-involvement in artifact and system development. The key question is whether such relationships would actually turn out to be *generative*. Certainly, as we saw, the first serious attempt to establish such a relationship, with SCI, failed to generate anything except disappointment and misunderstanding. One possible reason for this is that Echelon had very limited resources to deal with a very large number of problems, and hence it would have been extremely difficult to concentrate the kind of human and material resources necessary to construct a relationship satisfying the five criteria for generative potential discussed above even with one other company.

²³ Refs to books and manuals.

But a closer reading of RK's memo reveals an attitude toward potential relationships that RK himself shared with the dominant Echelon point of view his memo criticized: the assumption that the problems to be solved all lie in *artifact space*. That is, the problem is just that some physical system needs to be controlled by some particular concatenation of control artifacts. It is as though all that Echelon needed to know about their customers' business is what type of technical problems they were trying to solve. In fact, the situation is much more complicated. To understand what customers recognize as problems, and what counts for them as solutions, it is necessary to understand how the customers are situated in *agent space*: with what other agents they interact, about what. Control problems, that is, cannot be construed either as "natural" or as purely technical. The social context in which they are embedded can determine to a large extent what gets identified as a problem and what makes a solution acceptable. A generative relationship between Echelon and a control company would have to explore how the participating agents are situated in agent space, and how their situation there informs their *identities*: their understanding of what they do, how they do it, with whom.

Shortly after RK wrote his memo, Echelon hired a consultant, TW, to comment on its engineering analysis of a possible LonWorks application in process controls. TW's report highlighted some of the social facts Echelon would have to face if it expected to play in the process control market system. TW observed that "the history [of process control systems] emphasizes the importance of having tools and/or architecture which allow for the configuration of a large control system by technicians with a low level of training and skills. This has become one of the primary elements considered by customers in choosing a control system vendor for a project. . . The three most important considerations when choosing a vendor for a new system [are] ease of maintenance, flexibility, ease of calibration and configuration. System cost is near the bottom of the list."

A LonWorks-based process control system would have to satisfy these essentially social constraints from day one, or no one would ever buy it. More generally, Echelon would have to know, for every market system in which it expected to operate, what sorts of people and organizations specified, designed, installed, used and maintained the relevant control systems, and how the experiences and preferences of all these agents affected the decisions of whoever it was who purchased the systems. Some of the difficulties involved in trying to do this, and the resistance of established market systems to radical changes in their structure, will be explored in the next chapter.

References

- D. Lane & R. Maxfield (1997), Foresight, complexity and strategy. In W. B. Arthur, S. Durlauf, & D. Lane (Eds.). *Economy as a complex, evolving system II*. (pp. 169–198). Addison-Wesley.
- D. Lane & R. Maxfield (2005), Ontological uncertainty and innovation, *Journal of Evolutionary Economics*, 15, 3–50.

Chapter 10

Incorporating a New Technology into Agent-Artifact Space: The Case of Control System Automation in Europe

Federica Rossi, Paolo Bertossi, Paolo Gurisatti and Luisa Sovieni

10.1 Introduction

We contribute to the debate on innovation theory and policy by exploring, through the interpretative framework provided by Lane and Maxfield's theory of innovation (1996, 2005, this volume), a set of case studies concerning the implementation of a new technology for system automation and its incorporation into agent-artifact space (Lane & Maxfield, 1997). Our purposes are on the one hand, to illustrate to what extent this theoretical approach can help explain innovation processes, and, on the other, to derive some general implications for innovation theory. By focusing on agents involved in different kinds of interactions around the same technology, we introduce and compare different practices according to which a new technology can be incorporated into the existing structure of agent-artifact space, and we discuss some of the complexities involved in processes of technological adoption and diffusion, which are often modelled in excessively simplistic terms.

While Lane and Maxfield (this volume) introduce their theory of innovation by referring to episodes drawn from the early years in the history of LonWorks technology, from the late 1980s to the early 1990s, we describe the activities of some European agents involved with the same technology a few years later, from the mid-1990s onwards. After a brief introduction to the technology and to the issues that we then confront (Section 10.2), we then describe the technology provider's efforts to construct a market system for its technology in Europe and the difficulties it encountered when attempting to modify established market system structures (Section 10.3). In Section 10.4, we describe a small German company that builds integrated control systems using LonWorks technology; this company adopts a project-based approach and relies on a web of interactions to carve out niches for its complex products. In Section 10.5, we describe an innovative remote metering project jointly carried out by the technology provider and Italy's largest electric utility; in this case, the technology has been incorporated into an existing market

F. Rossi (✉)

Dipartimento di Economia, University of Torino, Via, Po 53, Torino, Italy
e-mail: rossi.federica@unito.it

system without substantial changes to its organization. Section 10.6 draws several theoretical implications from the analysis of these case studies and presents some concluding remarks.

10.2 LonWorks: The ‘Rise and Fall’ of a Market System Program

The LonWorks technology, officially launched in the US in 1990, implements communication and control functionalities among individual devices and sets of devices; it enables the construction of control systems that connect different artifacts, whether in a factory, a building, or a residential environment. Its creator is Echelon, a small company based in Silicon Valley. One of LonWorks’ most innovative features is that it permits distributed control architectures: that is, it can be used to build control systems whose nodes communicate with each other as peers rather than in a master-slave configuration. The technology’s basic building block is the Neuron Chip, a microprocessor provided with a unique identification number, which is able to run control algorithms and can thus distribute the intelligence of the system to the level of individual nodes.¹ Compared with pre-existing hierarchical systems, a distributed control system is more reliable because the malfunctioning of a single node does not impair system functionality, less expensive in terms of cabling, and extremely flexible because the system’s configuration is logical, not physical. The individual nodes’ ‘intelligence’ allows the network’s topology to be free, and transmission of information can take place over a variety of media (infrared, twisted pair, power-line, and others). LonWorks is a ‘general purpose’ technology (David, 1991) since it provides, at least at an abstract level, a general solution to the problem of control; so much so that when LonWorks was first introduced and for a long time afterwards, many novel applications, unforeseen when LonWorks technology was first launched in the market, began to emerge in a variety of different contexts (Lane and Maxfield, this volume).

In principle, it might have been possible for LonWorks to catalyze the creation of a new market system for distributed control technology, crossing existing industry boundaries, where different companies would provide hardware and software for the construction of distributed networks connecting devices of the most varied kinds, from industrial machines to heating systems to home appliances. The distributed control architecture is compatible with a structure of agent roles conforming to an ‘open system’ model of integration, whereby independent system integrators build and configure networks by combining interoperable devices and network integration products (such as routers, gateways, servers) that are made available ‘off the shelf’ by numerous competing companies, thus ensuring low prices and ease of

¹ For a detailed description of the Neuron Chip and LonWorks technology see Lane and Maxfield (this volume).

installation. The early documents that describe LonWorks show that such a model was precisely Echelon's vision.

However, this process has not taken place according to Echelon's expectations. The widespread diffusion of LonWorks as a general purpose technology, able to flatten hierarchical control networks and leading to the creation of a single market system for distributed control systems in different applications, has not materialized, for reasons that we explore in Section 10.3. LonWorks has entered different market systems that largely have remained separate, each maintaining its own physical and cognitive scaffolding structures, role structures, conventions, and each involving different actors. LonWorks artifacts have found their niches in specific applications like building automation (especially integrated building automation systems where LonWorks' market share in the mid-2000s was estimated at about 65%²), industrial automation (although the technology has only had a modest presence in this sector), local control networks (we present some examples in Section 10.4), remote metering (we describe a very important application in Section 10.5), and numerous other applications, quite unrelated to each other, such as pumps, robots, and train braking systems. The use of LonWorks in different environments seems to be leading to a 'speciation' of the technology, which is developing in different directions according to the applications for which it is used, while the technology's basic building blocks remain largely the same – although in time these may change as well.

10.3 Distributing Control in Agent-Artifact Space: Echelon

From the start, the technology provider was aware that the creation of a market system around LonWorks required the involvement of many other agents. Echelon would need to interact with numerous companies in order to enlist their cooperation, so that they would agree to supply the products and components necessary to build distributed control networks. In a way, Echelon would need to 'distribute control in agent space,' almost in parallel with its efforts to distribute control in artifact space. In particular, manufacturers of control systems for different applications (the so-called Original Equipment Manufacturers, or OEMs) would have to be convinced to expand their offer of interoperable components implementing the Neuron Chip and the LonTalk communication protocol; system integrators would have to be persuaded of the superiority of LonWorks over its competitors; and hardware

² This figure was published in a report commissioned by a Japanese firm, which estimated LonWorks' worldwide market share in the area of integrated, high-end building automation applications to be approximately 65%. The figure therefore refers to a very specific market segment which includes a small share of building automation installations, many of which consist instead of simple non-integrated systems (HVAC, lighting, access control and so on). Because the building automation industry is transversal with respect to other more established industries, it is in general quite difficult to obtain precise data at industry level in each country, and even more difficult to provide international comparisons.

and software suppliers would have to be recruited in the production of network integration products.

Echelon initially thought that this process would be straightforward, and that the superiority of their technology in terms of reliability, flexibility and openness would suffice to persuade their interlocutors to share their ‘distributed control network’ vision. But innovation processes are complex and multilevel; market systems rely on specific structures (cognitive and physical scaffolding structures, competence networks) that allow interactions to take place and to continue over time, and create and spread narratives that communicate shared attributions about the technology’s identity and the roles and prerogatives of the agents involved (Lane & Maxfield, 2005). Overcoming the resistance posed by the structures supporting established market systems would prove to be far more difficult than Echelon originally expected.

When LonWorks was commercialized in Europe, in the early 1990s, different control industries were characterized by very different features and histories. Established market systems existed for control systems in building environments, mainly for the control of heating, ventilation and air conditioning (HVAC) units, and for the control of production processes in industrial settings.

As it did in the United States, Echelon immediately sought to establish contacts with the large European OEMs in the building automation sector. In the early 1990s, most large OEMs still sold systems based on proprietary control software and communication protocols, and their interests, as they perceived them, conflicted with the ‘open system’ model promised by LonWorks. These companies worked with consolidated networks of relationships and well defined roles; they had their trusted technicians and single-brand integrators, and the exclusive maintenance contracts generated large profits. The possibility to ‘lock in’ their clients for the expansion and maintenance of their systems was a privilege that they did not want to relinquish, and, consequently, they were not willing to make their products easily interoperable with products manufactured by their competitors.

However, LonWorks also had numerous advantages: it was robust, and since it comprised a powerful chip and a range of tools that simplified communication among devices, it solved numerous communication problems within the control system, which the large companies that produce controls and ‘big iron’ (chillers, boilers and so forth) would prefer not to have to confront. Since the mid-1990s, the largest OEMs, particularly in the HVAC industry, began to use LonWorks as a communication protocol within their systems, while, often, they continued to offer one or more alternative product lines that used proprietary communication protocols. Several companies even inserted LonWorks in systems that were then closed to the outside, preventing communication with other systems.³

³ As late as 1998, Echelon complained, “While major manufacturers of control systems continue to adopt LonWorks technology at an accelerating pace, many are worried about the market changes that will be brought about by adoption of a standard network protocol and implementation of truly open architectures. Open architectures are viewed as a possible Pandora’s Box to larger companies with substantial market shares; they prefer to maintain the status quo, which keeps their customers

Although sales to large OEMs were, and remain, one of the principal sources of revenue for Echelon,⁴ this was not very publicized: large companies were not particularly interested in promoting LonWorks, which they considered as a simple component, and the fact that most chips became part of closed systems contrasted with the ‘open systems’ integration model that was actively pushed by Echelon, especially during the 1990s.⁵

In order to build a market system for distributed control, Echelon could not simply passively rely on the large OEMs’ cooperation; it needed to recruit allies that were more interested in open systems – smaller OEMs and independent system integrators, in particular, that would gain the most from a transition to a market system characterized by greater competition around interoperable products. LonWorks appeared particularly suited to these companies’ needs not only because it allowed peer-to-peer communication among devices via a standardized protocol, but also because it was already implemented in a range of commercially available hardware and software tools that made it easy to design, install and configure distributed control networks.

Throughout the 1990s, other technologies had been launched that promised to facilitate communication among control systems produced by different manufacturers, and, hence, to enable the construction of ‘open’ control networks in building automation. Those that drummed up greater industry support were BacNet and EIB, both promoted by OEM consortia.⁶ However, LonWorks and BacNet or EIB were not equivalent. The last two were simple protocols that enabled communication among control subsystems, and, while allowing system architectures that were more flexible than strictly hierarchical ones, they were not designed to support peer-to-peer communication among devices, as LonWorks did. Furthermore, each subsystem in the control network still had to be configured using the proprietary software tools offered by its producer. These technologies were designed to enable communication among systems without opening up the control systems themselves to competition.

Echelon began to realize that the construction of a new market system for distributed control technology would entail something more than simply convincing

boxed in (...) Many of these manufacturers have found the use of LonWorks technology to be a cost-effective way to allow their proprietary devices to share information within their own closed system and wish to leave it at that” (Echelon, 1998a).

⁴ According to information gathered in the course of several interviews, the main manufacturers in the HVAC field all use LonWorks: the ‘big five’ sharing more than 50% of the market for HVAC installations (Honeywell, Johnson Controls, Siemens, TAC, Satchwell Invensys) but also Philips, Trend, and others. Always according to our interviews, the producers that have adopted LonWorks to a greater extent are TAC and Honeywell.

⁵ Publicizing LonWorks as the technology that makes it possible to create open, interoperable networks, constituted Echelon’s principal communication strategy toward the outside (Echelon, 1998a, b, 1999, 2000b, 2005).

⁶ BacNet, published in 1995, was promoted by the American HVAC trade association (ASHRAE), while EIB, launched in 1990, was promoted by a European consortium (EIBA) led by Siemens, the leading European HVAC controls producer.

potential users of LonWorks' superiority. Rather, it would require creation of new scaffolding structures sustaining numerous processes necessary for the market system to function and persist over time: providing training for interested users, verifying the compliance of products to LonWorks' standard specifications, lobbying for the promotion of LonWorks as an international standard and confronting other lobbies with contrasting interests, and, in general, supporting the many interactions needed to reshuffle roles in agent space and create new competence networks. To achieve this, Echelon acted on various fronts: it set up an internal standards setting organization (LonMark Interoperability Association) and an annual trade fair (LonWorld), it supported creation of European user groups (the LonUsers associations), and it opened up key elements of the technology so that other agents would be able to construct compatible devices (in particular, the LonTalk communication protocol was published as an ANSI standard, and Echelon issued royalty-free licenses to all who wished to implement it).⁷

In the early 1990s, to promote interoperability between LonWorks-based products and create consensus around the distributed control movement, Echelon set up the LonMark Interoperability Association (LIA) in the US, a voluntary standards organization whose members included Echelon and a large number of user companies, mostly OEMs, system integrators, and software developers. LIA members organized into task groups with the main objective to develop standard specifications to which products based on LonWorks technology would have to conform to ensure interoperability with other products. As a scaffolding structure, organized as a quality promotion and standard-setting association, LonMark has promoted interactions around the technology and has supported the expansion of artifact space. In the late 2000s, LIA associates about 300 firms worldwide, for the most part located in the US and Canada, while the LIA Product Database includes over 1300 products, of which about half are offered by European firms.⁸

In 1991, Echelon set up a trade fair, LonWorld, to promote interaction and exchange among LonWorks users and to strengthen the consolidation of the market system. LonWorld provided an important meeting place for those users – small OEMs and integrators in particular – that were interested in realizing open control networks able to integrate different functions. From the mid-1990s, several of these users, in different European countries, began to organize into user groups, called LonUsers.⁹ Echelon played a key role in supporting their formation, involving people who committed to its concept of distributed control rather than to its products in a

⁷ At the same time, Echelon tried to preserve exclusive rights to other key elements of the technology (more often through secrecy and technological barriers rather than enforcement of proprietary rights), to safeguard the possibility to maintain some control on the innovation process and to derive some revenue from its activities.

⁸ The database (<http://www.echelon.com/productdb/>) does not provide a complete overview of existing LonWorks-based artifacts, since it includes only products that have been certified as interoperable by LonMark, but many other noncertified products are also commercially available.

⁹ The largest users groups are the German LonNutzer Organization (LNO) founded in 1993, and LonUser Sweden, founded in 1992; a few years later LonUsers Italia (1999), LonUsers UK (2000), and LonUsers France were formed, together with other user groups in smaller countries

narrow sense. Interestingly, LonUsers were specific to Europe: we may suppose that, in the United States, the companies that were more active in promoting LonWorks relied on the US-based LIA as a scaffold for meeting and exchange; furthermore, US system integrators working with open systems already had their own professional organizations, such as CSIA (Control System Integrators Association), which were absent in Europe.

With the support of these scaffolding structures, Echelon enjoyed some success in recruiting numerous agents to the complex task of creating a market system for LonWorks, and the zone of artifact space around this technology expanded considerably. During the second half of the 1990s, the LonUser groups were the most active agents engaged in promoting distributed control and LonWorks technology in Europe. LonUser associations became the privileged channel to form expertise and forge stable relationships between Echelon, the system integrators, and the smaller OEMs.

In the meantime, however, important changes were taking place in the zones of agent-artifact space around control systems for building and industrial automation.

Let us first consider industrial automation applications, where LonWorks' penetration has been, over time, very modest. While some commentators claim that it was LonWorks' technical features that made it unsuitable for applications where precision of execution and rapid reaction times are fundamental (such as most industrial processes), other more convincing reasons have to do with the structure of the entire 'sociotechnical system' (Ropohl, 1999; Mollerman & Broekhuis, 2001) surrounding industrial automation technologies. In fact, when LonWorks arrived on the market, open system architectures for industrial control systems had already consolidated.

System integrators that were building industrial control networks by integrating control subsystems via standard protocols had started to appear in the early to mid-1980s. Independent software companies had begun to specialize in supervision and control software (SCADA) for managing complex networks of industrial machines; over time, a few widely used standards (Profibus, CAN, OPC), developed by industry consortia, had become the dominant communication protocols. In the course of the 1980s, this 'open system' model of integration had the better of the 'closed' model based on proprietary control systems.¹⁰ Because of the consolidated presence of integrators able to build control systems from commercial 'off the shelf' components using already standardized interfaces, LonWorks could more easily be accommodated into existing networks, instead of precipitating dramatic changes in agent-artifact space. According to one of our informants, nowadays 'companies that

(Finland, Denmark, Norway, the Netherlands, Belgium). LonUsers Espana and LonUsers Poland were founded more recently.

¹⁰ "The manufacturers' representatives, distributors, and electrical contractors who had been doing PLC-based integration projects were now able to compete directly with the major control system manufacturers [such as, in the US, Foxboro, Fisher Controls, Fischer and Porter, and Bailey Controls] . . . That they did so successfully is clear as evidenced by the fact that Bailey, Fischer and Porter, Foxboro and Fisher Controls are no longer independent companies, and Honeywell is recovering from a failed merger attempt with General Electric" (CSIA, 2002).

are making solutions with Profibus, Allen Bradley all these PLC manufacturers. . . have 99.9% of the market share, and the industrial Lon is probably sold only by about by 20 companies worldwide.’

The evolution of building automation systems, from proprietary to more open systems, may prove similar to that followed by industrial controls. Compared with industrial automation, however, the movement toward open systems in building automation has been extremely slow. Demand for integrated control networks in buildings has grown very slowly, due in part to the less pressing need to drive down costs and increase safety standards in a building rather than a factory environment, and in part to the different nature of the users – property owners are typically less technically informed than managers of industrial plants and rely on the advice of consultants to solve their technical needs. Established role structures in the construction industry generally assign responsibility for control systems to ‘mechanical contractors,’ who are responsible for the HVAC system and for all the systems that, in the words of one of our informants, ‘have to do with water and cooling;’ their competencies lie in heating and refrigeration rather than in electronics and informatics, and their understanding of innovative communication and control systems often is poor. The use of standard contracts and the reliance on established tender requirements, which rigidly assign responsibilities for different systems to different contractors, also imply that only a very small fraction of building automation projects explicitly require integration among control systems, and even less specifically provide for open technologies. In addition, the market for home automation systems, which has been considered very promising by many industry representatives for decades, has so far failed to materialize, for reasons mainly connected to these systems’ high cost and electricians’ lack of specific control competences.

Despite slow growth in the demand for building control systems, which is also strongly constrained by the vicissitudes of the construction market, some important changes have taken place in the industry since the early 2000s. These are connected, in the first place, with changing attitudes and strategies on the part of control systems manufacturers. Since the beginning of the new century, several large OEMs have abandoned their closed product lines in favor of products that use standard communication technologies.¹¹ At the same time, to win the most innovative and complex projects – integrating different functions like lighting, access control, fire protection, HVAC – many large OEMs are trying to extend their expertise to system integration. Although to some extent, there has always been a conflicting model within OEMs that do both manufacturing and installation, this ‘double role’ has become more commonplace in the last 4–5 years.¹² The OEMs’ efforts to position themselves as

¹¹ An example is the agreement, signed in 2000, between Echelon and Honeywell (Echelon, 2000a), according to which the latter adopted LonWorks as standard in its own products, and it undertook to acquire them mainly from Echelon. In the same year, TAC and Echelon signed an agreement according to which TAC and its integrators became authorized retailers of Echelon software products (Echelon, 2000c).

¹² According to one informant, “In the last 4–5 years, the big controls companies have increasingly been moving into systems integration . . . these guys quite often play two roles, they are consultants

system integrators and providers of open systems modify role structures in ways that do not match Echelon's expectations or those of many of its small allies. If large OEMs acquire the necessary skills to perform complex integration projects, their strong brand visibility and their ability to drive down prices may undercut specialized LonWorks integrators. Among the small OEMs and integrators that specialize in LonWorks networks, some have moved towards even more complex networks, not necessarily in building automation, widening their technological competencies; others have successfully defended their market niches by maintaining consulting relationships with large OEMs, providing specialized LonWorks-related skills in the execution of technically advanced building automation projects.

Furthermore, while traditional OEMs are trying to grab a larger slice of the market for integrated projects, innovative system architectures and new competence networks are also appearing. As projects become more complex, involving the seamless integration of different subsystems, other integrators, specialized in industrial information systems and with highly sophisticated IT skills, entered the building automation industry, bringing with them technological architectures and 'cognitive schemes' typical of industrial automation.

Complex building automation projects increasingly are structured in ways that mirror the communication and control architectures now commonplace in industrial automation: while communication protocols such as LonWorks, BacNet and/or EIB are used for communication *within* subsystems, communication *among* subsystems takes place through widely used data transmission protocols, such as Ethernet, and thanks to generic supervision software. An example is the Atari building in Lyon, completed in 2001. The integration project was managed jointly by the system integrator Meta Productique and by the software developer Newron Systems. The latter is a small company with advanced LonWorks software competencies that often provides consulting services to other companies, OEMs and integrators, which wish to use LonWorks in complex projects for which they do not possess specific expertise. Meta Productique, with a background in industrial automation, realized a system with mixed hierarchical and distributed architectures, all based on open communication protocols: communication between the various devices within the offices (lights, air conditioning, heating) relied on LonWorks technology, while the general supervision of the general building control subsystems (which included automated reconfiguration of office spaces, access control, remote metering of energy consumption) was performed through an industrial supervision software installed on a server; communication between the subsystems and the server occurred via the OPC protocol on Ethernet-TCP/IP. Like in industrial automation, structures such as these may limit LonWorks' scope of application to sub-networks or 'islands' of distributed control within complex networks that have an overall hierarchical structure.

The movement toward new, more complex system architectures that include both hierarchical and distributed control and merge different technologies, is also

at the beginning and potentially they also provide the system integration. . . Over the last 5 years we have seen increasing development, in our controls clients, of business entities that will do the whole thing."

facilitated by a phenomenon that trade magazines call ‘convergence,’ that is, the tendency to integrate control systems – including building automation systems – to a higher level within the company’s information technology system. Many in the control industry express concern about this trend. The trade journal *Automated Buildings* (Gowan, 2002; Hartman, 2003) warns that it may bring about a reshuffling of roles in agent space, whereby those system integrators with an electric engineering background (as is currently the case of the majority of integrators in the building automation sector) may be confined to the role of simple installers, unless they quickly upgrade their skills in IT.

These processes have hampered the emergence of a single market system for distributed control. As LonWorks has found applications in different and often separate market systems, the LonWorks scaffolding structures, in particular the user groups that were very active in the second half of the 1990s, have lost importance and associates. LIA, now called LonMark International, has incorporated many of the European LonUser groups, but membership growth has nonetheless stalled in the last three/four years.

We can interpret LonWorks’ story throughout the 1990s and in the first years of the following decade as an example of the ‘rise and fall of a market system program.’ Echelon set out to build a market system for its technology by supporting a set of scaffolding structures that would sustain interactions and promote shared narratives centered on distributed control networks and open systems. However, the inertia of existing market systems structures, as well as defensive actions undertaken by the large players in order to protect their technologies, has proved instrumental in slowing down the growth of innovative applications and in hampering LonWorks’ diffusion. In the meantime, the existing market system structures have endeavoured to – and have eventually succeeded in – confining LonWorks to marginal roles that they could incorporate, instead of breaking down boundaries between them.

Echelon has successfully adapted to these changes: it has reduced its ambition to establish itself as the provider of general-purpose distributed control technology, and it has started to focus on system-level artifacts suited for specific applications – such as street lighting systems and, especially, remote metering systems.¹³ Echelon is well-positioned to benefit from the new global focus on energy conservation resulting from fears of global warming and energy crisis. Sales of LonWorks products have in fact picked up and have continued to grow over the last 3 years.

10.4 Innovation Through Networking: Tlon

Tlon is a small company actively constructing an environment favorable to the application of its exclusive competencies. It operates in market systems that are different but complementary and constitute a context where it can carve out a specialized

¹³ Some informants have identified 2003 as the year in which Echelon decided to invest heavily in remote metering applications and in powerline communication, one of Echelon’s technical strong points; the ENEL project, which we describe in Section 5, has been instrumental in bringing about this change of direction.

niche. This case exemplifies how LonWorks technology provides the opportunity for individual entrepreneurs, working together within generative relationships (Lane & Maxfield, 1996) and through scaffolding structures, to weave together projects that cut across existing market system boundaries. Generating these new ‘interstitial zones’ in agent-artifact space – which involve the integration of the focal technology with other technologies, some general purpose, some specific to the particular market systems that the projects link – may produce numerous difficulties but also potentially high rewards.

Founded in 1997, Tlon has headquarters in Schwäbisch Hall, in the Baden-Württemberg area. In many ways, it is a typical small European firm: the company coincides with the life and business objectives of its founder and of his closest collaborators, most investment is self-financed, return on research is modest, and return on investment takes place in the medium-long term. Innovation is mostly incremental, based on the application of existing technologies to new contexts; however, the company also takes part in sophisticated projects that may deliver relevant innovations.

Tlon is headed by an entrepreneur, VT, an electronic engineer with a long experience in the control industry. VT started to experiment with LonWorks technology in 1995, when he was managing the electronic controls department of a company that produced dyeing machines for the textile industry. After a couple of years’ experience he realized that his newly acquired competencies could be profitably applied to a wide range of control problems, and decided to start a ‘technology consulting’ company for customers interested in innovative electronic controls.

The project that boosted Tlon’s reputation was commissioned by a coffee machine manufacturer. This company was not satisfied with the controls produced in-house: the software was too complicated, it was unable to manage simultaneously all the functions it was supposed to perform, and it was hard to reconfigure. Tlon set up a system that constituted a textbook example of what it means to move from centralized to distributed control in manufacturing. The software was deconstructed into elementary components, which were then inserted into different Neuron chips and configured into a communication and control network. This system integrated different functions in flexible ways, allowing for easy and continuous upgrading. Between 1998 and 2001, Tlon introduced three successive incremental innovations in the same company, all based on LonWorks technology.¹⁴

In the first period of his entrepreneurial career (1998–2003), VT consolidated his relationship with Echelon. Participation in the relevant LonWorks scaffolds (LIA and the German LonUsers association, LNO) and the setting up of Infranet Partners (a consortium of firms interested in developing control network infrastructure, set up in 1999 by VT and a British colleague whom he had met through LonMark)

¹⁴ The first innovation was the development of a new system for the control of individual machines; the second concerned creation of a control network polling payment information from vending machines and transmitting it to the company’s accounting system; the third was the implementation of a remote control functionality for the network of vending machines, to reduce maintenance costs and enable the producer to monitor the machines’ state in real-time.

enabled Tlon to remain close to the circle of companies performing the most complex LonWorks-based projects. Tlon also started to cooperate with local agents that were interested in applying distributed control functionalities to other fields. Local technical schools contacted Tlon to carry out case studies and dissertations about problems that could be solved with LonWorks technology. Tlon gathered new business ideas and human resources from all these relationships.

When Echelon decided to concentrate on remote metering, focusing their attention on power-line communication and on solving large network problems, integrators working on complex projects had already begun to move toward new technological solutions, including wireless technologies such as ZigBee (for which, according to VT, Echelon has no specific expertise). As we noted in the previous section, system integrators and software developers still use LonWorks at the field level for communication among devices within control subsystems, but have now begun to use TCP/IP and software developed on an open source basis for integration at higher levels.

VT thinks that LonWorks is unlikely to become a popular choice for communication and control within simple systems, such as home automation or simple building automation projects: although the technology works well, it is expensive, and the market is already dominated by electricians and distributors working with large companies, like Siemens and Honeywell, which enjoy a strong reputation with final customers. Instead, Tlon has decided to concentrate on complex projects that only very skilled integrators can undertake. In particular, it focuses on complex networks, often comprising geographically dispersed devices, aimed at reducing energy consumption and improving device management. They have realized a network for the management of cooking devices in order to reduce energy consumption, a network for the remote management of HVAC systems in private homes, and a network for the management of photovoltaic panels in schools and government buildings, among other projects.

These projects, all developed locally, are conceived as components of a broader project, called 'Infranet Valley,' which VT hopes to realize in Schwäbisch Hall, and, maybe, to replicate elsewhere. The Infranet Valley project would extend the distributed control framework to networks that are geographically dispersed, to deliver functionalities to the individuals living and working in a certain territory. Initially (2000–2002), the plan was to construct a very large local network based on LonWorks technology, comprising around five million devices. The network would allow the local government to optimize energy consumption in local schools, swimming pools, and administrative buildings. Later (since 2003), the project has been reframed in terms of integration of specialized sub-networks, where communication is based on LonWorks and other protocols. Although VT continues to use LonWorks for device-level communication and wherever it provides the best communication solution, he has decided to acquire skills in the development of open-source software that he considers strategic for the realization of future projects. This change in perspective happened, and not by chance, when Echelon re-oriented its activities toward metering technologies and decided to pursue research on power-line rather than on wireless communication.

Infranet Valley is a long-term project, which builds on many small, local initiatives. VT continually attempts to build generative relationships with local actors (local officials, training and academic institutions, industrial associations), trying to create networks of operators interested in solving community problems, not specifically connected with distributed control. Tlon relies on numerous scaffolds: besides LIA, LNO, and Infranet Partners, it also takes part in associations and projects related to energy management, integrated management of public services, and initiatives promoted by local or national public bodies. VT thinks that networking at the local and international level can help his firm actively construct a market niche for complex local networks that is more promising than the market share or profit results that could be obtained by trying to position Tlon within the more established market system for simple automation projects. Participation in many relationship networks and scaffolds is expensive, but the positioning in the high end of the market guarantees the resources necessary to continue investments of this kind.

10.5 Maintaining Hierarchical Control: ENEL

In the late 1990s, a new application opened up for LonWorks technology that has proved instrumental in reorienting Echelon's activities: the remote control of utility meters. The first, and, so far the largest, application is ENEL's 'Telegestore' project, involving the installation of electronic meters in 27 million Italian households, which, in recent years, has constituted Echelon's main source of revenue.¹⁵ At its inception between the spring of 2000 and early 2001, the Telegestore project was extremely innovative. It was the first large-scale installation of this kind and has since become an important 'cultural' reference for many agents, opening the way for similar installations. In addition, at that time many in the industry anticipated that the electricity meter could be transformed into a 'residential gateway,' connecting, through power lines, the electricity provider upstream with any intelligent device downstream, including other meters, local plants, even home appliances. Some hoped that this innovation would leverage the launch of a radical change in the provision of user services through power-line communication, and, in particular, that it would finally launch the long anticipated and much coveted home automation market. As it turned out, though, the Telegestore project developed along more conservative lines, without spurring radical innovations in the provision of user services via the electricity grid. While it is interesting to investigate what prevented the process from going ahead according to the initial expectations, our main focus is on how LonWorks technology has been incorporated into existing structures in agent-artifact space, so that the realization of new functionalities has been carried out by existing competence networks under the control of a central agent, ENEL.

¹⁵ Between 2001 and 2005, ENEL has been by far Echelon's main customer; for example, in 2004 it has brought in \$64.1 million out of \$109.9 million of total revenue.

Founded in 1962, ENEL (Ente Nazionale per l'Energia Elettrica) is a large enterprise with public capital.¹⁶ For almost four decades, it was a *de iure* monopolist in electricity generation and distribution in Italy. Its position began to weaken during the 1990s with the introduction of provisions aimed at liberalizing the energy market, in particular the so called 'Bersani Act' (March 16, 1999) which established that, from January 1st 2003, no company could produce or import more than 50% of the total electric power produced or imported in Italy.

Therefore, in the early 2000's, ENEL was forced to give up an important share of its business. ENEL's management devised a new role for the company as a 'multi-utility,' diversifying into other networked services by reallocating the proceeds accrued from the mandatory sale of parts of its activity. ENEL tried to expand into water distribution, with Acquedotto Pugliese, but the initiative was not successful. It also acquired the gas distribution network Camuzzi SpA, the second largest gas supplier in Italy. Telecommunications were considered particularly important: between 1999 and 2001, ENEL became the second Italian telecommunications provider, through the acquisition of part of the mobile telephony operator Wind and the acquisition of the fixed telephony operator Infostrada. In this context of convergence between electricity distribution and telecommunications, a large remote customer management project began to take shape.

Many large European utilities had performed experiments in remote management in the 1990s. ENEL had carried out a pilot project called SITRED, experimenting with meters that recorded information about consumption and about the state of the network, sent it to a control center through power lines, and received instructions through the same means. The main technical difference between SITRED and the successive Telegestore project was the use of an electromechanical meter and the simultaneous data accumulation by an electronic support connected to each meter, which could be remotely controlled. The project involved several of ENEL's established suppliers: Siemens, Landis and Gyr, Schlumberger, Copeco, Ducati, ABB Elettroconduttore, Feme, and Bticino. The system was successfully tested on 70,000 households in and around Rome.

By the mid-1990s, the group of engineers in charge of this and other projects was convinced of the technical feasibility of a meter that could become a 'user interface.' This project moved to the forefront shortly after the appointment of Franco Tatò as ENEL's CEO in 1996. Tatò, a manager with international experience, had a background in communications, and it was under his direction that ENEL undertook a strategy of diversification into other sectors.

We can identify several reasons behind ENEL's decision to venture into a large project such as Telegestore. In 1997, an economic feasibility study highlighted the possibility of large savings for ENEL, both in meter reading (which could be performed automatically instead of manually) and in the control of energy flows,

¹⁶ ENEL was formed from the merger of over 1,200 private and public local companies, and, for a long time, was the second largest Italian group, after FIAT, in terms of revenue and employees.

reducing the revenue losses due to network failures and customer misreporting or tampering.

More generally, in the strategy of transformation into a multi-utility, the opportunity to maintain and reinforce the relationships with the final customers and to increase the range of available services (from the provision of customized energy distribution, to the remote management of customers' accounts, to mobile telephony) could provide a remarkable advantage *vis-à-vis* the competitors and could hinder the entry of other providers.

Finally, the importance of being the first utility worldwide to activate such a cutting-edge service did not escape ENEL's management. This might enable ENEL to sell the service to other providers, whether in Italy – where the local utilities were buying a large share of the electric distribution network – or abroad.

ENEL decided to disband the group that had worked on SITRED and to use internal capabilities more substantially. Unlike in SITRED, it was decided to use a fully electronic meter that was able to manage information directly. This decision seems to have involved a direct intervention on the part of Tatò, who wanted a radically innovative project, without the electromechanical/electronic technological compromise that could hamper further innovation.¹⁷

First, through an international tender, ENEL selected the British company AMPY to design an electronic meter with features meeting their current needs. To enable communication between the meter and the concentration and data processing center, ENEL could use its own electricity network and the mobile telephony network of its controlled company Wind. It was decided that communication between the concentration centers, each of which serves about 200 users, and the processing center would happen via GPRS, while communication between the meters and the concentration centers would be enabled by Echelon's LonWorks technology, over power-line. This choice seemed to have been motivated, according to sources internal to ENEL, by technical reasons. In fact, Echelon had developed a very efficient system for data transmission over power-line, and it would be able to resolve any communication problem rapidly using its own products.

Echelon and ENEL collaborated for about three years, from 2000 to 2003. In the first period, ENEL's internal team and two technicians from Echelon worked together in developing the devices. After this technical collaboration, kept behind closed doors, a preparatory agreement, dated May 10, 2000, was jointly presented by the two companies. The press release stated, 'Through this technology, numerous value-added services can be offered for the remote management of homes and offices, such as monitoring, use and reparation of home appliances and of security systems, control of air conditioning and lighting devices. The services offered can be accessed by customers through the Internet, fixed or mobile telephone, and will be managed by the ENEL group.'

¹⁷ It has been said that Tatò's radically innovative approach – implying the complete substitution of electromechanical technology with electronic technology, rather than their complementary use – came from his experience in Olivetti. Tatò, in fact, was Olivetti's manager during the complicated transition from typewriters to calculating machines.

Therefore, at least in the initial phases of the project, the LonWorks-enabled meter was interpreted by ENEL as a tool that could launch a radical innovation in the relationship with the final customer, which may eventually lead to a revolution in the utilities' market system, allowing the implementation of a wide range of services. The performances expected for the Telegestore system were not limited to the remote management of energy consumption, however; through the Neuron chip and two transceivers, it would be possible to connect, via power-line, the devices already present in the home (heating system, home appliances, access control) and transmit and receive data to and from a control center. This seemed the right opportunity to jumpstart – especially in Italy – the so far elusive home automation market.

In addition to ENEL and Echelon, another partner at the beginning of the project was Merloni Industries, the third largest European manufacturer of home appliances. The official agreement between ENEL and Merloni Elettrodomestici was dated October 19, 2000. The two companies agreed to cooperate to experiment with innovative forms of payment for home appliances connected via power-line to the intelligent meter. In the same period, Merloni consolidated its relationship with Echelon, with whom they had their first contact in 1993. The two companies issued statements announcing the deployment of Echelon's power-line transceivers within Merloni products (although the product under study, Leonardo, has never been commercialized). At the 2001 SMAU fair, Merloni presented the pay-per-use washing machine that communicated with the electronic meter: it would be rented and paid according to effective use. At LonWorld 2001, the most popular debate was The Residential Gateway – Gateway Challenges, where the first speaker was the manager in charge of the Telegestore project. The home automation revolution seemed about to happen.

In 2001, the Telegestore project began. Thirty million meters had to be produced and installed; the network of concentrators and the information network had to be set up. The meters were produced in different countries, mostly in Eastern Europe and in China, where relationships with the various sub-suppliers were mediated by Shenzhen Kaifa Technology Co.

The start of the Telegestore project in Italy was not painless. Many critics, within the Competition Authority, within the trade unions, and among the environmentalists, were sceptical that ENEL could manage such a complex innovation. These suspicions were fueled by the fact that the meters showed a worrying tendency to disconnect immediately customers that exceeded their contractually established wattage, unlike the old electromechanical meters that had elevated tolerances: this provoked harsh reaction from some consumers' associations.

Another source of controversy concerned Echelon's involvement as a supplier. The contract for the communication system – with a value of approximately 300 million Euro – had been assigned to Echelon without a competitive tender, and the agreement had been followed with the acquisition by ENEL of a share of about 9% in Echelon, through which ENEL gained a seat in Echelon's board. This close association with an American company was criticized by the new center-right

government formed in 2001, particularly by the right-wing nationalist *Alleanza Nazionale* party.¹⁸

Tatò resigned on May 24, 2002 – for numerous reasons, mainly political, probably not directly related to the Telegestore project – and Paolo Scaroni, another internationally renowned manager but with a background in manufacturing rather than communications, was appointed in his place. Under the new management, the project to transform ENEL into a multi-utility was abandoned, and in particular it was decided that ENEL would not engage in sectors different from energy distribution.

After Tatò's resignation, the relationships between ENEL and Echelon became difficult, as they had not consolidated any network of relationships below the CEO level, neither among the two groups' staffs nor in collaboration with external actors. There was not a clear, shared vision of the route that the project could follow, or a group that could build it. Moreover, evident contrasts between ENEL and Echelon appeared, concerning the promotion of the project to other utilities in Italy and abroad. Despite the continuation of the supply relationship, collaboration between ENEL and Echelon ceased in 2003: the controversy was brought before an international arbitration and was resolved in September 2005.¹⁹

ENEL decided to restructure the project's organization. In meter design, Ampy was replaced by Kaifa, which, together with ENEL, quickly perfected a replacement electronic meter; the communication protocol, appropriately modified, came from the SITRED project. STMicroelectronics adapted a chip kit capable of running the software. The announcement was made at the end of 2004, after the deal had been done. In order to sell Telegestore to local energy distribution firms, ENEL formed an external alliance with IBM, signed in March 2004.

In an article published in *Metering International* magazine (2004), Telegestore was presented as a service for managing electricity users and collecting and transmitting information, while the article was silent on the possibility of downstream connection.²⁰ In fact, the electronic meter, as installed in Italian homes, did not and could not connect the downstream appliances to the electricity provider's communication system. We can advance some hypotheses concerning the reasons why the project did not achieve all the functionalities that were initially expected. It is plausible that, at the start of the project, a communication technology that allowed connection with downstream devices or systems was still not available, or that it was

¹⁸ The parliamentary interrogation presented on 2/19/2002 by the Apulian MPs Ivano Leccisi and Ugo Lisi affirmed, 'how could it be that Echelon was chosen without public tender as supplier of ENEL, when the ENEL purchases will cover a share superior to half of its revenue,' and it was asked 'if the government thinks it admissible that an operation of this kind can be financed with public money; if the government does not think it appropriate, once possible responsibilities have been ascertained, to adopt immediate provisions vis-à-vis ENEL's current management.'

¹⁹ The International Chamber of Commerce largely swayed towards ENEL; however, it established much lower compensations than those that were demanded.

²⁰ There was only a brief, vague mention: 'The Telegestore system opens an outdoor-indoor communication channel. ENEL is technically evaluating the co-existence of the metering services with energy-related value added services.'

not sufficiently reliable or convenient. The decision probably was made to start with the technology that was available at the time, and, in the meantime, to continue to research a solution capable of providing more advanced communication performances. However, for various reasons (reshuffling of ENEL's top management and strategic priorities, problems in the ENEL/Echelon relationship) this process stopped, the 'new generation' Echelon/Ampy electronic meter was never realized, and the efficiency of rapid installation was prioritized.

Today the Telegestore project has come to a rapid completion. Sales to other Italian providers have been mediated by IBM: most local utilities in Italy have bought the Telegestore system, while another installation has been made in Cantabria, Spain.

Echelon independently launched its NES – Networked Energy Services – division in 2003, after carrying out experiments in Holland and New Zealand. At the end of 2004, it commercialized an electric meter capable of collecting information and receiving and sending commands.²¹ However, it was only at the end of 2005 that Echelon began to receive some reward for its efforts in terms of sales, thanks to a project with Swedish utility Vattenfall. The Vattenfall project also seems to be exclusively directed at meter reading and remote customer management, not differently from Telegestore. Since then, Echelon has received several very large orders; they have begun trials in various countries and have set up a growing network of NES Value-Added Resellers. For the first quarter of 2007, Echelon reported revenues of \$25M from the NES product line.

The Telegestore example highlights an altogether different way in which a technology can be incorporated into a potential market system, based on a specific new functionality that the technology helps make possible. Few relationships developed around the Telegestore project, and they were mostly mediated by ENEL: crosscutting relationships among the other participants did not develop or only developed minimally, and in no way did these relationships consolidate into a competence network. ENEL probably was responsible for this situation since it maintained tight control of collaborations, keeping the option to discontinue them at any time, and, this way, it fully controlled the process and the results. In carrying out this project, ENEL did not rely on external scaffolds. Its involvement in LonUsers was minimal (ENEL joined in the LonUsers Italia association in 2001, but it left a short time later, in 2003), and the relationship with Merloni did not produce any tangible results. With respect to LonWorks' market system, ENEL was not interested in distributed control *per se*. Rather, it used LonWorks technology as a competitive tool in order to attain better control of the existing market. Although an initial convergence had taken place between Echelon and ENEL's management's attributions – concerning the possible use of the electronic meter as a residential gateway enabling communication between devices inside and outside the home – this process was

²¹ The meters have been designed by a team of engineers that Echelon hired from a failed startup in the Silicon Valley area; they are manufactured by subcontractors in China and Hungary, and marketed by a number of NES Value-Added Resellers around the world.

interrupted, for reasons that we described above. The ENEL contract gave Echelon substantial revenues for several years, but it did not catalyze the launch of home automation nor did it boost the growth of a market system for distributed control. However, it convinced Echelon of the importance of orienting its activities towards remote metering applications and network-based services. That is, Echelon's attribution of its own role transformed from that of technology provider and market system coordinator to system – bundles of products and services – provider.

10.6 Lessons from LonWorks

Our analysis of several case studies confirms the interpretative power of a dynamic interactionist perspective, and it allows us to derive some useful implications for innovation theory.

Very often, when a new artifact is produced, it enters a pre-existing socioeconomic system characterized by established competence networks, physical scaffolding structures, consolidated role structures whose interactions with the new artifact strongly constrain the process of adoption, the ways in which the new artifact is used, and how its technical characteristics evolve. The production and commercialization of a new artifact does not simply trigger a cumulative diffusion process, representable with a logistic curve, but it catalyzes a variety of complex processes through which the artifact itself evolves and the set of other artifacts in conjunction with which it is used, installed, maintained and developed changes, sometimes spawning an entirely new artifact family. The cognitive attributions that agents make about the artifact's functionality change over time too, as well as the identities of the agents that are somehow involved in its supply and use.

The diffusion of LonWorks technology cannot be modelled as a simple 'epidemic diffusion' process, since the technology is changing over time and the artifact family is expanding. Even if we consider only the relatively unchanging building block, the Neuron Chip, the fact that the applications are so varied and that new ones are discovered over time implies that the size of the possible market for Neuron Chips keeps expanding, so that it becomes difficult to understand when the market may be 'saturated' by the new technology. Further, the population of adopters in different applications is heterogeneous and geographically dispersed, and it is hard to imagine that the diffusion process may happen with any regularity. Quantitatively, the diffusion of LonWorks artifacts seems to be characterized by major discontinuities, corresponding to important sale agreements being signed or called off, rather than a process of gradual technological contagion following the spread of information about the new technology.

Innovation processes are not simply driven by the technical characteristics of certain artifacts: the interpretations that different agents make of them are also crucial in driving their actions and hence in determining the overall shape of the process (Bijker, Hughes, T., & Pinch, 1987). From our examples, it is apparent that the technology had different meanings for different agents. While Echelon, at least initially,

saw its technology as a general purpose solution to control problems, large OEMs saw it first and foremost as a way to simplify communication within their systems and, at a later stage, as a tool that would allow them to access more sophisticated integration projects. For ENEL, LonWorks was an efficient solution to contingent communication problems; only few actors within ENEL shared Echelon's view of LonWorks as a technology enabling communication with the wider world of 'intelligent devices.' Tlon sees LonWorks as a technology that permits the realization of complex territorial networks: it does not only replace traditional products, it can also allow networking among artifacts that in the past were neither 'intelligent' nor integrated.

Issues of technological superiority and price advantage are negotiated in the social domain, and perceived or actual technological superiority does not immediately guarantee the success of a technology over its competitors. Rather than being known to everyone in the market system, even the technology's basic features are the object of negotiation and debate.

The process of constructing a market system for the new technology often requires creation of new role structures and new scaffolding structures, or the re-orientation of pre-existing ones. Often, the process develops across established industry boundaries and over a very long time span. For example, changes in role structures were necessary to sustain the creation of a market system for distributed control networks based on LonWorks technology: the creation of such networks would require the presence of independent systems integrators buying components 'off the shelf,' made available by OEMs who would agree to make their product conform to certain interoperability standards. In practice, the two roles of OEM and system integrator are not clearly separated, and this has given rise to different kinds of possible interaction schemes and different efforts on the part of the various agents to protect their own technologies and sources of revenues, thus, hampering the development of a market system along the lines originally envisioned by Echelon.

The consolidation of competence networks able to carry out the system's functionality is also very important to allow market systems to operate in practice. For example, it has taken a very long time for LonWorks competence networks to be established in building automation, and this has hampered the diffusion of the technology, so much so that it may have missed its 'window of opportunity' in favor of other technologies, imported from industrial automation via different competence networks. Another example is the lack of control competences on the part of electricians, which has certainly hindered the diffusion of home automation products.

The processes that construct a new market system result from a combination of innovation projects, which trigger subsequent cascades of changes but that are not lined up along a 'natural' trajectory (Nelson & Winter, 1982) determined by the artifact's intrinsic features. Each of the agents that we encountered developed its own course of action in conditions of ontological uncertainty (Lane & Maxfield, 2005) based on personal evaluations, attributions, and narrative structures, independent of those produced by Echelon: the sum of these individual actions cannot be considered as the predictable consequence of certain events. While the concept of 'technological trajectory' (Dosi, 1982; Nelson & Winter, 1982) may be useful to

describe, retrospectively and with the benefit of hindsight, the evolution of a broadly conceived technological system, it is not useful to interpret innovation processes in the making, since they are characterized by constant novelty, idiosyncrasy, and path dependency, and they are affected by the ‘hierarchically tangled’ (Lane, 2006) actions and interactions of agents located at different levels of social organization.

Acknowledgments The authors gratefully acknowledge financial support from the European Union research contract ISCOM-IST-2001-35505. We are grateful to the other researchers involved in the ISCOM project, who have provided helpful comments in numerous occasions, and particularly during seminars held in Paris, Reggio Emilia and Venezia in the period 2003–2006. We are particularly indebted to David A. Lane and Robert Maxfield, who have read several drafts of this article and have given us invaluable feedback and advice. Finally, we are very grateful to the many people we have interviewed in the course of this research, who have generously shared with us some of their time, expertise and views.

References

- Bijker, W., Hughes, T., & Pinch, T. (Eds.) (1987). *The social construction of technological systems*. Cambridge, MA: MIT Press.
- CSIA. (2002). *The market for control system integrators*. White Paper Prepared for The Control System Integrators Association.
- David, P. (1991). Computer and dynamo. The modern productivity paradox in a not too distant mirror. In *Technology and productivity: The challenge for economic policy* (pp. 315–345). Paris, France: OECD.
- Dosi, G. (1982). Technological paradigms and technological trajectories. *Research Policy*, 11, 147–208.
- Echelon. (1998a). Going open, wide open, White paper. http://www.echelon.com/support/documentation/papers/Going_Open.pdf. Accessed 13 November 2007.
- Echelon. (1998b). End to end solutions with LonWorks control technology, White paper. <http://www.echelon.com/solutions/opensystems/papers/end2end.pdf>. Accessed 13 November 2007.
- Echelon. (1999). Introduction to LonWorks systems, White paper. <http://www.echelon.com/support/documentation/manuals/general/078-0183-01A.pdf>. Accessed 13 November 2007.
- Echelon. (2000a). Honeywell enters into international purchase agreement with echelon, press release 07 February 2000. <http://www.echelon.com/company/press/2000/honeywellcommit.htm>. Accessed 13 November 2007.
- Echelon. (2000b). Implementing open, interoperable building control Systems, White paper. <http://www.echelon.com/solutions/building/papers/BacNetComp.pdf>. Accessed 13 November 2007.
- Echelon. (2000c). TAC AB signs international LonWorks® open systems reseller agreement companies to jointly promote LonWorks system to building controls market, Press Release 21 March 2000. <http://www.echelon.com/company/press/2000/tacab.htm>. Accessed 13 November 2007.
- Echelon. (2005). Open systems design guide, White paper. <http://www.echelon.com/support/documentation/papers/OpenSysDesignGuide.pdf>. Accessed 13 November 2007.
- Gowan, J. (2002). Industry focus shifts to web convergence. *Automated Buildings*, November, <http://www.automatedbuildings.com/news/nov02/articles/sin/sin.htm>. Accessed 13 November 2007.
- Hartman, T. (2003). Convergence: What is it, what will it mean, and when will it happen? *Automated Buildings*, April. <http://www.automatedbuildings.com/news/apr03/articles/thessup/thessup.htm>. Accessed 13 November 2007.
- Lane, D. A. (2006). Hierarchy, complexity, society. In D. Pumain (Ed.), *Hierarchies in natural and social systems* (pp. 81–120). Dordrecht, The Netherlands: Kluwer.

- Lane, D. A., & Maxfield, R. (1996). Strategy under complexity: Fostering generative relationships. *Long Range Planning*, 29, 215–231.
- Lane, D. A., & Maxfield, R. (1997). Foresight, complexity and strategy. In W.B. Arthur, S. N. Durlauf, & D. A. Lane (Eds.), *The economy as an evolving complex system II* (pp. 169–198), *SFI studies in the sciences of complexity* (Vol. 27). Boston, MA: Addison-Wessley.
- Lane, D.A., & Maxfield, R. (2005). Ontological uncertainty and innovation. *Journal of Evolutionary Economics*, 15(1), 3–50.
- Mollerman, E., & Broekhuis, M. (2001). Sociotechnical systems: towards an organizational learning approach. *Journal of Engineering and Technology Management*, 18(3), 271–294.
- Nelson, R., & Winter, S. (1982). *An evolutionary theory of economic change*. Cambridge, MA: Harvard University Press.
- Ropohl, G. (1999). Philosophy of socio-technical systems. *Techné: Journal of the Society for Philosophy and Technology*, 4(3) 59–71.

Chapter 11

Innovation Policy: Levels and Levers

Federica Rossi and Margherita Russo

11.1 Introduction

In this chapter, we focus on the policy implications of the theoretical approach to understanding innovation outlined in the previous contributions in this volume. How can this approach help policymakers to implement effective interventions, foster innovation processes and create structures to sustain them over time? To address this question, we first discuss the problematic relationship between innovation theory and policy objectives and implementation; we illustrate this issue, in Section 11.2, by discussing the example of current European innovation policy. In Section 11.3, we explore the rationale for innovation policy provided by our complexity-based perspective to innovation (Lane and Maxfield, 1997, 2005; Lane et al., 2006; Cane et al., this volume; Read et al., this volume; Russo, 2000), and we introduce the main features of this approach. We claim that policy analysis should mainly be concerned with the processes through which policy problems and solutions can be identified, rather than with the direct formulation of general policy recommendations. To substantiate this claim, in Section 11.4, we present a specific example of policy analysis based upon our approach and its associated tools. In Section 11.5, drawing from this exercise, we present some methodological remarks as well as an agenda for future research.

11.2 Theory and Implementation in Innovation Policy: The European Scenario

When designing, implementing and evaluating policies, awareness of the theoretical framework that inspires them is crucial to ensure consistency between policy measures and tools available for their monitoring and evaluation. Policy analysis

M. Russo (✉)

Dipartimento di Economia Politica, University of Modena and Reggio Emilia, Viale Berengario 51, 41100 Modena, Italy
e-mail: margherita.russo@unimore.it

should not only investigate the most effective policy instruments, it should also clarify their theoretical underpinnings, which may carry very different implications for policy. In addition, it must be remembered that policy measures are implemented within specific institutional and administrative contexts that guide them in practice, not necessarily in the same direction as the stated objectives.

In Europe, institutional decisions concerning innovation policy in the last ten years have been driven by a theoretical framework that has been made public through numerous 'guidance' documents issued by Community institutions, the European Commission, and the European Council in particular.¹ Support for innovation was first acknowledged as a public policy goal in the Green Paper on Innovation (1995), followed by the First Action Plan for Innovation in Europe (1996), which included numerous policy suggestions. Since then, and particularly after the Lisbon European Council (2000), innovation has gained increasing importance in the context of European development policies, which have the objective to improve and consolidate the competitiveness of the economic system, and which often pay special attention to small firms (European Commission, 2004; European Council, 2000, 2005). The connections among innovation, competition, and development and the role of innovation in the knowledge economy are the keystones of these policies, which are explicitly inspired by a 'market approach' (European Commission, 2000, 2003c). This term is used both to emphasize the need for private intervention alongside public incentives, and, quite often, to portray a view of innovation in terms of 'supply' and 'demand.' In this case, innovation is seen as a sequence of individual actions so that it is possible to identify certain actors that 'demand' and others that 'supply' innovations (the former are firms producing goods and services, while the latter are universities and research institutions).

The theoretical framework shaping European innovation policy has been influenced by the academic discourse on the topic of innovation. Over the last twenty years, economic and sociological theories of innovation have changed markedly, abandoning the traditional linear view of innovation in favor of systemic approaches. The linear view of innovation conceptualizes the innovation process as a sequence of well-defined, temporally and conceptually distinct stages. Although rarely codified in the economic literature, the linear model is widely shared, often implicitly, in academic discourse.² The model postulates that innovation starts with basic research activities performed by an individual inventor or, more often, by a research group, which lead to an invention – that is, a new idea or a new entity that is not yet ready for commercial exploitation. Subsequently, applied research and development activities, usually performed within industrial research laboratories, lead to embedding the invention into an artifact or process that can be commercially exploited. The

¹ This issue was explored in detail in a previous paper (Rossi, 2007), where we attempted to reconstruct the theoretical view of innovation that underpins Community recommendations, and we explored the innovation-supporting interventions that have been sponsored with EU funds.

² For a comprehensive reconstruction of the historical development of the linear model, see Godin (2006).

resulting innovations are sold in the marketplace, adopted by users, and imitated by other companies. Individual adoption choices lead to a process of diffusion of such innovations.

Recent approaches, instead, tend to conceptualize innovation in systemic terms, as a process that involves, at each moment, many actors (sometimes impossible to identify in advance), their relationships, and the social and economic context in which they are embedded. Kline and Rosenberg (1986) suggested that innovation should be represented by a chain-linked model, in which the various aspects of economic and scientific activity, internal and external to the firm, are linked together by multiple relationships of causality and feedback. Economic issues, technical issues, and the existence of a demand for innovation are all interdependent elements of the process of innovation. This model has opened the way to numerous systemic conceptions of the innovation process, seen as the result of dynamic interactions among heterogeneous elements.³

The most recent perspectives in economics see innovation as a process of creation of new, often tacit, knowledge. Increasing attention for the cognitive aspects of innovation has fostered a corresponding surge in interest for interactions among agents as sources of new knowledge: direct interactions among people are, in fact, the main modes of transmission and creation of tacit knowledge.⁴ Researchers have begun to study various forms of cooperation between firms that are directed at developing innovations (Freeman, 1991; Mowery and Teece, 1996), including user-producer interactions (Von Hippel, 1978; Lundvall, 1985; Russo, 2000). Sociologists and organization theorists have underlined the importance of the cognitive distance among agents in stimulating innovation (Lundvall, 1992; Nooteboom, 1999), while other scholars have claimed that it is instead geographical proximity among firms – which often implies cognitive proximity – that fosters innovation.

The influential literature on national systems of innovation – which emerged at the beginning of the 1990s with the path-breaking contributions by Freeman (1988), Lundvall (1988, 1992), and Nelson (1988, 1993) – has highlighted the interplay of a wide range of factors, organizations, and policies influencing the capabilities of a nation's firms to innovate (Nelson, 1993).

In the last ten years, such heterodox approaches to the analysis of innovation and technological change have influenced policymaking at European institutions.⁵ This theoretical redirection even has been acknowledged explicitly in some of the European Commission's documents.⁶

³ Among these, an approach that, in the late 1990s, found favor among academics and policymakers is the so-called 'triple helix' (Etzkowitz and Leydesdorff, 2000). Developed in the evolutionary economics framework, this approach suggests that cycles of production, innovation and policy-making mutually evolve as in a triple helix, in which dynamic selection takes place both within each helix and between helixes, fostering interactive and recursive relationships.

⁴ This issue was first raised in the literature by Hägerstrand (1965, 1970) and Polanyi (1969).

⁵ Mytelka and Smith (2002) reconstruct the role that some heterodox economic theories have had in influencing the policymakers' thinking within institutions like the European Commission and OECD, but not within others, such as the World Bank.

⁶ For example, in COM(2003)112 and COM(2003)27.

A first consequence of the adoption of a systemic approach to innovation has been the transition from framing innovation policy exclusively in the context of research and industrial policy to a more 'transversal' approach. In fact, some of the 'systemic' policy objectives outlined by the Commission – for instance, 'fostering an innovation culture,' 'establishing a framework conducive to innovation,' and 'better articulating research and innovation' as stated in the First Action Plan for Innovation in Europe (1996) – can be attained only through a mixture of interventions in several policy fields, involving education, social, industrial, enterprise, development, and research policies. The Commission recognizes that innovation policies must be implemented through interventions that involve not only the activities of basic scientific research, development, and commercialization of research outcomes – according to the linear model described above – but also small and medium firms and the social and institutional contexts in which they operate.⁷ In the same direction, another strand of policy analysis is linking innovation not only to the actions of isolated companies, but also to the activities of 'clusters' intended as aggregations of organizations.⁸

Another related consequence is that besides the usual 'top-down' interventions, 'bottom-up' interventions are emphasized, where the role of Community institutions mainly is to enable and coordinate policies rather than to dictate their contents (Triulzi, 1999).

However, despite the widespread attention paid to innovation issues by researchers and policymakers alike, and, despite the quantity of funds that are being channelled into innovation-supporting activities,⁹ the relationship between innovation theory and implementation of innovation policy is problematic. A large gap remains between the comprehensive approach to innovation advocated by the Commission and the range of interventions that are being funded in practice. The latter can be subsumed within a narrow list of topics: providing information services that facilitate interactions among different kinds of institutional actors; simplifying and extending access to patent protection; increasing firms' private research expenditure; supporting small innovative firms and start-ups through interventions aimed at simplifying

⁷ COM(2003)112 states, "The evolution of the innovation concept - from the linear model having R&D as the starting point to the systemic model in which innovation arises from complex interactions between individuals, organizations and their operating environment - demonstrates that innovation policies must extend their focus beyond the link with research. Since it is through enterprises that the economic benefit of the successful exploitation of novelty is captured, the enterprise is at the heart of the innovation process. Innovation policy must have its ultimate effect on enterprises: their behaviour, capabilities, and operating environment" (European Commission, 2003, p.4).

⁸ COM(2008)652 states that the EU "identified strengthening clusters in Europe as one of the nine strategic priorities for successfully promoting innovation"(European Commission, 2008, p.2)."

⁹ According to our estimates, expenditure on innovation-related interventions in the EU (broadly intended to include Framework Programme interventions as well as innovation-supporting measures sponsored by the Structural Funds) increased from approximately 6,052 million euro per year in the period 1994–1997 to approximately 7,404 million euro per year in the period 2002–2005 (figures computed from data presented in Rossi, 2007).

bureaucracy and granting access to innovation financing; and supporting innovation through public procurement. Although the Commission explicitly recognizes the systemic nature of innovation phenomena, actual interventions are generally not consistent with these premises.¹⁰ Therefore, even when policymakers engage in conscious efforts to develop a sophisticated theoretical framework on which to ground their innovation policies, implementation is far from easy or automatic.

This happens for several reasons. First, while the theoretical framework and the interventions performed are continually evolving, they are not perfectly synchronized. The relationship among these actions is mediated by numerous institutional levels and by processes that take place on different time and social scales so that the actors that are responsible for developing a broad theoretical framework to guide policy are generally different from those that devise concrete policy measures.

Secondly, as we have remarked previously, some objectives of innovation policy would require implementation of coordinated interventions involving multiple policy fields, from education to social inclusion, from research to entrepreneurship. In this respect, European innovation policy appears to suffer from several constraints: the current policy framework, characterized by vertically separated policy fields, the organizational structure of the Commission, and the funds' rules and scope for intervention all place limitations on the kinds of policies that can be designed to support a complex process like innovation. Functional separation among policy areas at the Commission level is often mirrored by similar administrative boundaries at the regional level, so that integrated interventions – at any territorial level – are rarely implemented. Although European documents increasingly stress this issue, present rules in the deployment and use of EU funds (administrative procedures, evaluation criteria, and monitoring tools) are still hampering the realization of coordinated interventions.

Third, policy programs, once established, tend to consolidate, continue, and expand over time, with the risk that interventions may overlap and that self-referential communities of actors accessing most of the funding may be created, further hampering the policies' effectiveness.

Finally, the problem with policy implementation is not simply procedural, but also conceptual. Broad attempts at theorizing innovation processes do not lend themselves to a quick translation into simple 'policy recipes,' precisely because conceptualizing innovation in systemic terms – or, as we argue, as a complex, multi-level process – means that it is not possible to devise context-independent ways to support it.

For these reasons, improved theoretical understanding of innovation processes should not aim to provide policymakers with simple encompassing solutions, but it should help them formulate and address questions that are appropriate to the particular context in which they operate. In this sense, innovation policies should have a 'local' dimension, that is, they always should be 'rooted in localities identified by

¹⁰ See European Commission (2003a) "Investing in research: an action plan for Europe."

sets of relations within specific communities of people, firms and institutions,' as Bellandi and Di Tommaso (2006) remarked with reference to industrial policy.

11.3 Rethinking Innovation

11.3.1 Changing Rationales for Innovation Policy: From Market Failure to Process Building

Mainstream economics still views the innovation process as fundamentally linear, and is mostly concerned with the problem of inducing economic agents to produce enough scientific and technological knowledge to feed into this process. Once the 'optimal' amount of knowledge is produced, firms just have to 'use' it to develop new products, and 'the market' – characterized by a pure or, more frequently, monopolistic competition framework – will clear supply and demand of the new goods through the usual price mechanism.

In this context, market failures associated with knowledge production provide a powerful rationale to justify public intervention. According to some well-known arguments, because the outcomes of basic and even applied research activities are characterized by uncertainty (with respect to both the timing and the quality of results achieved), low appropriability, and non-rivalry, market mechanisms would lead to an underproduction of the knowledge necessary for innovation (Arrow, 1962). Hence, public funding of science – either through direct state intervention (creation of public research institutions) or through public spending (funding of the university system and the research system) – is needed to ensure that the socially optimal amount of scientific knowledge is produced (Nelson, 1959). Other institutional mechanisms – in particular the patent system and legal devices such as trade secrecy – are designed to induce profit-seeking firms to fund private R&D expenditure by increasing knowledge appropriability and imposing limitations to its use by rivals. Aside from the mechanisms devised to foster knowledge production, knowledge use (and related markets for knowledge) is also subject to market failures due to information asymmetries (Akerlof, 1970; Arrow, 1971) that justify existence and public support of other appropriate institutional mechanisms designed to curb those failures: various contractual forms, the promotion of standards and public regulations, copyright protection, and so forth.

When innovation is viewed as a complex, rather than linear, process, the rationale for public intervention changes. If we think of policy as having a 'local dimension,' in the sense outlined above, policymakers can intervene not only to promote correction of market failures, but also to achieve specific strategic objectives or to reach meta-economic ones – for example, promoting access to knowledge, education, and health; fostering social or environmental sustainability; and attaining a specific distribution of wealth or a status of development (Belliandi and Di Tommaso, 2006). Innovation promotion is a key element when implementing strategies to attain local development, and this, in turn, provides grounds to justify public intervention.

11.3.2 A Complexity Perspective to Innovation

Several contributions in this volume allow us to gain an improved understanding of innovation from a complexity perspective. This perspective adopts an ‘organization thinking’ rather than a ‘population thinking’ approach (Lane, Maxfield, Read and van der Leeuw, this volume). Supra-individual social structures are not seen simply as aggregates of component entities, deprived of agency, useful only to ‘monitor changes in frequency distributions of their component’s properties’ (Lane, Maxfield, Read and van der Leeuw, this volume), as species are in biology. Rather, to understand human socio-cultural change, it must be taken into account that both individuals and organizations, belonging to ‘tangled hierarchies’ (Lane, 2006), are endowed with agency. They can generate changes in the structure of agent-artifact space, and, in turn, their structure and the functionalities they support can be modified by the actions of entities positioned at other levels in the social hierarchy. In this analytical framework, there is no *ex-ante* selection of a level of analysis to understand innovation processes and social processes in general. The focus shifts from proving causal relationships between variables to understanding how different structures of relationships carry different functionalities over time, and, therefore, support different kinds of processes.

With this approach, not only do we obtain a better grasp of the processes constraining innovation, from macro scaling relationships constraining growth (West et al., Pumain et al., this volume) to cognitive constraints limiting organizational size and dynamics (van der Leeuw et al., this volume), but we improve our understanding of the micro- and meso-level processes underpinning innovation, thanks to the development of a specific theory based on the agent-artifact space ontology (Lane and Maxfield, 1997). This ontology provides a language to describe innovation processes; it is a language in which syntax describes causal relationships between entities and processes, and, as such, it is theoretical, not simply phenomenological (Lane and Maxfield, 2005).

The theory pays great attention to the role of uncertainty in innovation processes, although the concept of uncertainty used here goes beyond its common characterization in terms of probabilistic knowledge. Instead, the theory claims that individuals involved in innovation processes act in situations characterized by ‘ontological uncertainty,’ that is, situations where economic agents do not know what the relevant entities are that inhabit their world, which kinds of interactions these entities have among themselves, and how entities and interaction modalities change as a result of previous interactions. The impossibility to assess what entities will affect the results of the agents’ actions prevents any evaluation of future outcomes, even in probabilistic terms (for a discussion of the concept of ontological uncertainty and its relationship with the concept of probabilizable uncertainty, see Lane and Maxfield, 2005). To explain how innovation processes take place in conditions of ontological uncertainty, a three-level theoretical framework is presented. At the level of the individual agents (micro-level), this approach describes how ontological uncertainty can be managed by agents in the short term through the adoption of a ‘narrative theory of action.’ At the level of agent interactions (meso-level), it claims

that innovation processes can result from particular kinds of relationships called ‘generative relationships,’ and, that a relationship’s ‘generative potential’ can be monitored by paying attention to some of its features. At the level of market systems (macro-level), it claims that agents take part in (formal or informal) organizations called ‘scaffolding structures’ to better manage ontological uncertainty and create ‘competence networks’ able to sustain and reproduce the functionalities needed for the market system to survive over time.

Therefore, the main building blocks in this theory of innovation are the concepts of generative relationships, competence networks, scaffolding structures, and the role of narrative in driving action in situations characterized by ontological uncertainty (Lane, Malerba, Maxfield, and Orsenigo, 1996; Lane and Maxfield, 1997, 2005, this volume; Russo, 2000, 2005).

11.4 Innovation Theory and Innovation Policy: Lessons from an Empirical Investigation

The complexity perspective adopted here enables us to identify analytical tools that can be applied to policy design, implementation, monitoring and evaluation activities. In this section, we illustrate these tools through our analysis of a specific policy experiment, the ‘Technological Innovation in Tuscany’ program (henceforth RPIA-ITT).

This program, implemented in the period 2001–2004¹¹ and funded in the context of the Regional Programme of Innovative Actions within the European Regional Development Fund, was intended to stimulate technological innovation processes in the Tuscan economy through the creation of cooperation networks among heterogeneous organizations – large and small firms, research centers, universities, local public institutions, business services providers, training agencies, and finance institutions – with the purpose of integrating competences and testing new methodologies for promoting innovation. The decision to fund cooperation networks was relatively unusual and complied with recent recommendations to promote systemic development in production structures composed of SMEs (Audretsch, 2002; European Commission, 2003; European Council, 2000). This intervention was not only directed to supporting innovation processes, it was in its own right an experiment in innovative policy design. The Innovative Actions Programme, although assigned a relatively small budget, provides a framework for experimenting with new ways of community structural intervention.¹²

¹¹ We were not involved in policy design, but in the latter stages of policy analysis and assessment. The methodology and results obtained from this analysis are presented in detail in Russo and Rossi (2009).

¹² Several European regions, using available EU funds such as those assigned to the Ris, Ritts and Ris+ programs, have promoted policies for supporting innovation in SME local production systems (surveyed in the papers by Nauwelaers and Wintjes, 2003; Bachtler and Brown,

The policy measure that we studied appears to be particularly close to the spirit of the theoretical approach proposed here. Through this program, Tuscany's regional administration was trying to support innovation processes performed by competence networks characterized by heterogeneity, which could potentially foster generative relationships among local actors, and, thus, trigger cascades of changes in agent-artifact space.

The main objectives of our analysis were: (1) to assess whether the program had succeeded in promoting the creation of well-functioning networks capable of integrating heterogeneous competences and of fostering systemic effects in the regional economy; (2) to understand the extent to which the program supported pre-existing networks of relationships or sparked the creation of new ones; and (3) to derive some suggestions that could be generalized to other innovation-supporting interventions.

In order to understand the structural characteristics of the networks of relationships underpinning the program and to explore some of the systemic effects that resulted from it, we relied on complementary use of social network analysis and qualitative interviews with the actors involved.

Network analysis was performed on two levels. First, we reconstructed the networks of relationships *within* each funded project, using the participants' joint involvement in the various work modules of the project as proxy for the existence of a relationship between them. Secondly, we explored the network of relationships underpinning the program as a whole. Here, we used the participation of the same organization in two project proposals as a proxy for the existence of a relationship between the projects. We used visualization techniques (Freeman, 2000), and we computed statistics relating to the network's cohesion and the nodes' centrality. The betweenness, closeness, and degree centrality indexes (Freeman, 1979; Wasserman and Faust, 1994; Degenne and Forsé, 1999) highlighted the organizations that were most and least actively involved in the program, and, therefore, helped us select the organizations to be interviewed. While the study of the structure of relationships underpinning each funded project allowed us to better understand which agents were able to facilitate the generative relationships that support innovation, the study of the general network's cohesion¹³ allowed us to assess if there were one or more cohesive subgroups of actors whose initiative was fundamental in recruiting a large number of organizations to the program.

We present some of the 'lessons' that we have drawn from this case study as an example of how our theoretical approach can provide the appropriate lenses through which innovation policies can be designed, analyzed and monitored.

Lesson 1: Developing innovative tools to monitor and assess the networks' generative potential. One of the implicit assumptions underlying the design of the RPIA-ITT program is that innovation processes can be fostered by exploiting existing relationships and by supporting and consolidating generative relationships among organizations that are not accustomed to interacting with each other. Surprisingly,

2004; Landabaso and Mouton, 2005; Rossi, 2007), but not many of them have explicitly focused on sustaining cooperation networks.

¹³ See Moody and White (2003) for a critical survey on this notion.

however, we found that the process of network creation and the networks' evolution over time were not carefully monitored by the policymaker, even in the context of a program explicitly designed for policy experimentation, where great attention should have been paid to unanticipated effects.

To assess the projects' achievements, the regional administration used indicators relating only to the products realized by each network (patents, prototypes, software, publications, workshops, training courses). In our view, however, it would have been more fitting to the program's aims to focus on the interaction processes that enabled such products to be obtained and to assess how changes in network composition, in terms of partners involved and their competences, affected a network's success; which organizations proved to be more successful in recruiting partners and obtaining funding; and what kinds of interactions were more conducive to successful innovation activities. We were able to answer these questions, at least in part, by integrating the information collected by the regional administration with the results emerging from our network analysis and from the qualitative interviews (Spradley, 1979; Agar, 1996, 2004; Russo and Rossi, 2008). Our approach allowed us to provide the regional administration with suggestions concerning the kind of monitoring tools (which information should be collected, in which form, how data should be organized and some suggestions for their interpretation) that they should have set up from the start to assess the success of this policy program (Russo and Rossi, 2007).

Lesson 2: Timing of policies, projects, and innovation processes. It is widely recognized that the detailed development and, particularly, the timing of innovation processes cannot be foreseen, even with respect to innovations that have already been acknowledged as commercially viable (Rosenberg, 1994; Lane and Maxfield, 1997, 2005, this volume; Rossi, Bertossi, Gurisatti, Sovieni, this volume). Exploitation of results itself is a process that cannot always be implemented, and it often is not even clearly identified, in the limited time available for policy intervention, which, in the RPIA-ITT, was a scant 13 months. In this respect, the RPIA-ITT program suffered from a problem that we had already observed while surveying EU innovation policies (Rossi, 2007): the lack of attention for what happens to products and services once they were brought to market to start competing with other products and services. Besides statements about the importance of understanding innovation as a system, even in the stage of the definition of general policy directions it appears that innovation continues to be conceived as a phenomenon that unfolds according to well-defined stages and for which it is possible to identify a clear beginning and an end. The effects that new products and services, once marketed, have on the socioeconomic system, remain out of sight just when they start producing (or not) those effects on growth in order to obtain which innovation policies are designed.

To avoid this shortcoming, the time span in which effects of the policy program are appraised should be reconsidered. By studying the effects that the program has produced over a longer time span, the policymaker should be able to assess the generative capacity of the relationships activated in the course of the program and of the new relationships emerging due to the activities performed, as well as the actors'

ability to initiate cascades of changes in agent-artifact space. In the RPIA-ITT case, the results described in the concluding reports should have been updated some time (at least 12–18 months) after the conclusion of the program. As it emerged from our interviews with people in charge of funded projects, changes that took place after the projects were formally concluded were of utmost importance both for the networks and agents involved and for the impact of that project on the regional system of innovation. The assessment exercises describing the effects of each project should also consider to what extent the projects led to further projects or benefited from the simultaneous implementation of other projects.

Lesson 3: Multivocal actors. Studying the structure of relationships underpinning the cooperation networks allowed us to better understand which agents are able to facilitate the generative relationships that support innovation. For example, we found that a number of business service providers were playing a special and important role in the RPIA-ITT networks. These service providers had different structural characteristics, different behaviors, and different objectives. However, most of them were active in the fields of training, certification, and technology transfer, a set of activities that allowed them to acquire a good knowledge of a wide range of companies' needs and potential (in terms of missing competences or idiosyncrasies) and to weave a close fabric of relationships with manufacturing firms and other local actors, such as trade associations and local governments. All this brought them close to many different contexts from which they learned several languages. Their 'multivocality' was important to identify local needs and sustain network creation. In particular, they were instrumental in bridging the world of applied research with small firms that had not previously been involved in collaborations with external organizations. The latter could be either 'follower' firms that were willing to participate in the projects once a core of participants had been established or very small manufacturing firms whose activity was entirely focused on production, which were unlikely to establish dialogue with academia or with industrial research centers. In many instances, the service providers proved able to set up new projects starting from their experience accrued through participation to previous EU-funded projects. These experiences also equipped them with the ability to monitor funding opportunities and to manage the relevant administrative-accounting procedures. These skills are crucial for organizing bottom-up policy interventions involving the participation of small firms.

Lesson 4: Central actors in generating networks. From the analysis of the RPIA-ITT program, we were able to study how the program tapped into a pre-existing network of relationships and was, in turn, able to influence that network. We found that some actors were central in presenting projects and implementing funded proposals. Although only 10% of the participants controlled almost half of the financial resources of the entire program, they were able, through multiple direct and indirect links, to involve a large number of other actors, many of whom had no previous experience of contact with research centers or universities. While the regional policymaker did not specifically target or monitor these actors, the key roles that they played in the program – from coordinating the project proposals to recruiting potential partners – became apparent both from our interviews and from our network

analysis exercises. We also collected, through the interviews, some evidence that certain relationships formed in the context of the RPIA-ITT were likely to give rise to other joint collaborations and new applications for public funding – a typical ‘emergent’ effect of the policy program, which is generally ignored in the evaluation exercises.

Starting from this experimental program, it may be possible to take advantage of this experience in the design of future programs, exploiting the knowledge of how generative relationships can be created within the local system, and, at the same time, reducing the difficulties inherent in fostering joint action among different organizations.

11.5 Implications and Applications of a Complexity Perspective

As we have previously claimed, the strength of our approach mainly is related to its ability to solicit the right questions for policy. Starting from an improved understanding and description of innovation processes, policymakers can design more effective interventions by exploring combinations between top-down and bottom-up measures, paying attention to meso-level structures, and reflecting on who should be involved in policy design and planning, who should be the recipients of innovation policies, what kinds of processes should these policies attempt to influence, and in what ways.

Performing the analysis of a specific case study through these theoretical lenses has not only allowed us to derive some recommendations specific to that program, but it has also suggested some methodological considerations that we present as concluding remarks.

We consider generative relationships as the privileged loci where shifts in attributions of identity and functionality take place. To foster innovation, policymakers should attempt to increase and monitor relationships with high ‘generative potential’ (Lane et al., 1996; Lane and Maxfield, 1997). To do so, policymakers would first need to explore what kinds of interactions – among which kinds of organizations and concerning which kinds of activities – support innovation processes; what are the most likely interaction loci that promote the emergence of generative relationships; and how can interactions with high generative potential be identified, monitored, and supported. Also, when interactions are not producing desirable results in terms of innovation, policymakers should explore innovative instruments to foster change. For example, they could seek to identify the main ‘narrative structures’ driving the actions of organizations in a certain area or industry and understand whether these are hampering or promoting innovation.

At a higher level than bilateral relationships, Lane and Maxfield (this volume) claim that innovation processes are sustained by specific cognitive and physical scaffolds and require the creation of competence networks. To sustain competence networks supporting production and innovation, policymakers should understand their structure and scope: they should identify whether local actors belong to local,

regional, national, or international competence networks and which structures, if any, coordinate the competences required at the local or industry level with the training needs of individuals and organizations. They should explore how such structures can be monitored and supported, whether there are any ‘missing links’ in the competence networks at the local or national levels, whether coordination with other policy fields (education, social, industrial) is required to design appropriate interventions, and, finally, whether it is possible to design policies that foster the emergence of new competence networks – promoting interactions between organizations that are involved in producing, using, installing the same technology or similar technologies – thereby encouraging the development of new applications.

In order to implement effective interventions, it is also crucial to identify the key agents and scaffolding structures that support local innovation processes so that policies can be designed to work with them in promoting innovation. For example, in the RPIA-ITT case study, we observed that the business service providers’ multi-locality is a key competence in order to select potential participants in a program requiring network activities, and that some scaffolding structures, such as a network of research centers in optoelectronic, were important drivers of innovation processes. Both business services providers and research centers acted as important scaffolding structures for the innovative projects supported by the RPIA-ITT program, although their roles were very different: the former played a crucial role in network formation and management, while the latter were instrumental in formulating the project proposals and in ensuring access to regional funds (Russo and Rossi, 2008). Policymakers should understand how such scaffolds can be identified and monitored, and, if necessary, how they can be supported. Based on our empirical research, we would also like to formulate a more general claim. We argue that a public agent that intends to promote innovation processes should rely on the ability to mobilize resources, to engage them in creative interaction, to accept and drive change, and to foster non-routine activities and processes. One example is a practice of the kind advocated by Hirschman (1967, 1995 rev.) with his ‘hiding hand principle,’ which he claimed was most effective in the implementation of development projects: a practice which intends to relieve actors from the difficulties connected with confronting unexpected events, that are endemic in innovation processes, by concentrating resources on improving their ability to creatively solve problems. Such abilities will then increase their attitude to initiate further changes in the future, even after public incentives have finished.

Policymakers must also have the necessary tools to understand the wider range of effects that their policies engender. Policy interventions have implications at the micro-, the meso-, and the macro-level. They directly influence micro-interactions among individuals and organizations, whether such policies consist in distributing funds, tax relief, or provision of information. They can also determine meso-level changes in the organization of relationships within a certain industry or territory. Finally, they can affect macro variables such as imports and exports and investment and overall expenditure both in public and private R&D.

In a complexity-based approach, innovation policies should be evaluated with respect to the systemic effects that they produce. Policies change the space of

interactions, whereby new actors are attracted to the innovation processes, recurrent patterns of interactions are consolidated, and new organizational structures are created. Rather than assigning more resources to standard ways of monitoring these policies, new ways of monitoring and evaluation should be sought: policymakers should devise new indicators, new evaluation procedures that take into account different time scales that different innovation processes require to produce tangible results, and new ways of assessing the processes of interaction among the involved actors in the course of the program. Some potentially effective tools have been identified through our empirical analysis, where we claimed that, apart from traditional firm-level indicators (number of patents, of new products, expenditure in R&D, expansion in the number of users or potential users of a technology, etc.), policy effects should be measured by new ‘potentially generative’ relationships activated, changes in the structure of competence networks, new scaffolds created, and changes in the patterns of use of artifacts, or systems of artifacts.

Complex system analysis provides other useful instruments to anticipate and assess potential impacts of policy programs, such as tools from network analysis (see, for example, the work on multi-nets described by White, this volume) and agent-based models (Lane, Serra, Villani, and Ansaloni, 2005, this volume), which policymakers recognize as potentially very important.¹⁴

We believe that simulation models constitute a fruitful direction for further research into policy issues, allowing policymakers to explore alternative scenarios. Simulation models, applied to the analysis of specific policy programs, can provide important insights with respect to the most effective ways in which public and private resources can be deployed to foster innovation. While these ‘artificial laboratories’ represent socioeconomic reality in a necessarily simplified way, they enable comparisons among the effects of different policy instruments, and in this way they provide an important theoretical support for policymakers in the choice of what specific set of measures to implement.

Once again, our experience with the analysis of RPIA-ITT provides a useful reference. This exercise has allowed us to collect information about the features of an experimental policy program directed at sponsoring innovation projects performed by networks of organizations. This information can prove useful to set up an artificial context which, although extremely simplified, would contain at least some elements of realism with respect to the relevant degree of agent heterogeneity, the number and types of organizations involved in the experiment, and the structure of the networks of relationships among them. After some of these realistic parameters

¹⁴ According to a recent report from a European Complex Systems Society workshop: “Complex systems science can help elucidate the systemic processes that embed decisions in case of crisis into a web of forces and interactions across biological, ecological, and societal dimensions. Relatively simple models inspired by complex systems science and its precursors in non-linear dynamical systems theory could have tremendous communication value in making clear what are otherwise non-intuitive ideas – for example, that we should expect small changes sometimes to lead to large consequences, or that perfect prediction is not always or even often an option”. (Dum and Buchanan, 2007).

have been entered in our artificial context, it would be possible to tweak them and visualize, thanks to the simulation, different scenarios in terms of the effects that the different policy parameters produce, and, in this way, compare the effectiveness of different interventions.

Acknowledgments This chapter collects and organizes the main results that have emerged from research carried out in the context of the following projects: The Information Society as a Complex System (ISCOM) funded by the European Commission DG Information Society Technology, 'Politiche per l'innovazione e sistemi di PMI', funded by Dipartimento delle Politiche di Sviluppo e Coesione, Ministero del Tesoro (May–December 2004), and 'Analisi delle azioni innovative e modellizzazione dei risultati', funded by Regione Toscana – Action 5 of RPIA (May–December 2004).

References

- Agar, M. (1996). *The professional stranger: An informal introduction to ethnography*. San Diego, CA: Academic Press.
- Agar, M. (2004). We have met the other and we're all nonlinear: Ethnography as a nonlinear dynamic system. *Complexity*, 10(2), 16–24.
- Akerlof, G. (1970). The market for lemons: Qualitative uncertainty and the market mechanism. *Quarterly Journal of Economics* 84, 488–500.
- Arrow, K. J. (1962). Economic welfare and the allocation of resources for invention. In R.R. Nelson (Ed.), *The rate and direction of inventive activity: Economic and social factors*. (pp. 609–626) NBER, Princeton, NJ: Princeton University Press.
- Arrow, K. J. (1971). *Essays in the theory of risk-bearing*. Chicago, IL: Markham Press.
- Audretsch, D. (2002). *Entrepreneurship: A survey of the literature*. Prepared for the European Commission, Enterprise Directorate General.
- Bachtler, J., and Brown, R. (2004). *Innovation and regional development: transition towards a knowledge-based economy*. Report prepared for the Knowledge based regional development and innovation Conference, Florence, Italy, 25 and 26 November.
- Bellandi, M., and Di Tommaso, M. (2006). The local dimensions of industrial policy. In P. Bianchi, and S. Labory (Eds.), *International handbook on industrial policy* (pp. 342–361). Cheltenham, UK: Edward Elgar.
- Degenne, A., and Forsé, M. (1999). *Introducing social networks*. London, UK: Sage.
- Dum, R., and Buchanan, M. (2006). Final report on the European Complex Systems Society Workshop "Science of complex systems: can science guide policies?" 19–20 December, Brussels, Belgium.
- Etzkowitz, H., and Leydesdorff, L. (2000). The dynamics of innovation: from national systems and "Mode 2" to a triple helix of university-industry-government relations. *Research Policy*, 29, 109–123.
- European Commission. (1995). Green Paper on Innovation, COM(1995)688.
- European Commission. (1996). First Action Plan for Innovation in Europe, COM(1996)589.
- European Commission. (2000). Towards a European Research Area. COM(2000)6.
- European Commission. (2003a). Investing in Research: an Action Plan for Europe. COM(2003)226.
- European Commission. (2003b). Green Paper on Entrepreneurship, COM(2003)27.
- European Commission. (2003c). Innovation Policy: Updating the Union's Approach in the Context of the Lisbon Strategy COM(2003)112.
- European Commission. (2004). Delivering Lisbon – Reforms for the Enlarged Union, COM(2004)29.

- European Commission. (2008). Towards world-class clusters in the European Union: Implementing the broad-based innovation strategy, COM(2008)652.
- European Council. (2000). Council Decision on a Multiannual Programme for Enterprise and Entrepreneurship, and in Particular for Small and Medium Enterprises.
- European Council. (2005). Presidency conclusions of the Brussels European Council 22 and 23 March.
- Freeman, C. (1988). Japan: A new national system of innovation? In G. Dosi, C. Freeman, R. R. Nelson, G. Silverberg, and L. Soete (Eds.), *Technical change and economic theory* (pp. 38–66). London, UK: Pinter Publishers.
- Freeman, C. (1991). Networks of innovators: A synthesis of research issues. *Research Policy*, 20, 499–515.
- Freeman, L. C. (1979). Centrality in social networks: I. conceptual clarification. *Social Networks*, 1, 215–239.
- Freeman, L. C. (2000). Visualizing social networks. *Journal of Social Structure*, 1(1).
- Godin, B. (2006). The linear model of innovation: The historical construction of an analytical framework. *Science Technology Human Values*, 31, 639.
- Hägerstrand, T. (1965). Quantitative techniques for analysis of the speed of information technology. In C. A. Anderson, M. J. Bowman (Eds.), *Education and economic development* (pp. 244–280). Chicago, IL: Aldine Publishing Company.
- Hägerstrand, T. (1970). What about people in regional science? *Papers of the Regional Science Association*, 24, 7–21.
- Hirschman, A.O. (1967, 1995 rev.). *Development Projects Observed*, Washington, DC: Brookings Institution.
- Kline, S., and Rosenberg, N. (1986). An overview of innovation. In R. Landau, and N. Rosenberg (Eds.), *The positive sum strategy* (pp. 275–305). Washington, DC: National Academy Press.
- Landabaso, M., and Mouton, B. (2005). Towards a different regional innovation policy: eight years of European experience through the European Regional Development Fund innovative actions. In M. van Geenhuizen, D. V. Gibson, and M. V. Heitor (Eds.), *Regional development and conditions for innovation in the network society* (pp. 209–240). Purdue, IL: Purdue University Press.
- Lane, D. A. (2006). Hierarchy, complexity, society. In D. Pumain (Ed.), *Hierarchies in natural and social systems* (pp. 81–120), Dordrecht, The Netherlands: Kluwer.
- Lane, D. A., and Maxfield, R. (1997). Foresight, complexity and strategy. In W. B. Arthur, S. N. Durlauf, and D. A. Lane (Eds.), *The economy as an evolving complex system II* (pp. 169–198). *SFI studies in the sciences of complexity* (Vol. 27). Boston, MA: Addison-Wessley.
- Lane, D. A., and Maxfield, R. (2005). Ontological uncertainty and innovation. *Journal of Evolutionary Economics*, 15(1), 3–50.
- Lane, D. A., Malerba, F., Maxfield, R., and Orsenigo, L. (1996). Choice and action. *Journal of Evolutionary Economics*, 6, 43–76.
- Lane, D. A., Serra, R., Villani, M., and Ansaloni, L. (2005). A theory-based dynamical model of innovation processes. *ComplexUs*, 2, 177–194.
- Lundvall, B.-A. (1985). *Product innovation and user-producer interaction*. Aalborg, Denmark: Aalborg University Press.
- Lundvall, B.-A. (1988). Innovation as an interactive process: from user-producer interaction to the national system of innovation. In G. Dosi, C. Freeman, R. Nelson, G. Silverberg, and L. Soete (Eds.), *Technical change and economic theory* (pp. 349–369). London, UK: Pinter Publishers.
- Lundvall, B.-A., (Ed.) (1992). *National systems of innovation: Towards a theory of innovation and interactive learning*. London, UK: Pinter Publishers.
- Moody, J., and White, D. R. (2003). Structural cohesion and embeddedness: A hierarchical concept of social groups. *American Sociological Review*, 68(1), 103–127.
- Mowery D., and Teece, D. (1996). Strategic alliances and industrial research. In R. Rosenbloom, and W. Spencer (Eds.), *Engines of innovation: U.S. industrial research at the end of an era* (pp. 111–129). Boston, MA: Harvard Business School Press.

- Mytelka, L., and Smith, K. (2002). Policy learning and innovation theory: An interactive and co-evolving process. *Research Policy*, 31, 1467–1479.
- Nauwelaers, C., and Wintjes, R. (2003). Towards a new paradigm for innovation policies? In B. Asheim, A. Isaksen, C. Nauwelaers, and F. Töttdling (Eds.), *Regional innovation policy for small-medium enterprises* (pp. 123–220). Cheltenham, UK: Edward Elgar.
- Nelson, R. R. (1959). The simple economics of basic scientific research. *The Journal of Political Economy*, 67(3), 297–306.
- Nelson, R. R. (1988). Institutions supporting technical change in the United States. In G. Dosi, C. Freeman, R. Nelson, G. Silverberg, and L. Soete (Eds.), *Technical change and economic theory* (pp. 312–329). London, UK: Pinter Publishers.
- Nelson, R. R., (Ed.) (1993). *National innovation systems. A comparative analysis*. Oxford, UK: Oxford University Press.
- Nooteboom, B. (1999). Innovation, learning and industrial organization, *Cambridge Journal of Economics*, 23, 127–150.
- Polanyi, M. (1969). *Knowing and Being*, Chicago: University of Chicago Press.
- Rosenberg, N. (1994). *Exploring the black box*. Cambridge, UK: Cambridge University Press.
- Rossi, F. (2007). *Innovation policy in the European Union: Instruments and objectives*. MPRA Paper, Munich, Germany: University Library of Munich.
- Russo, M. (2000). Complementary innovations and generative relationships: An ethnographic study. *Economics of Innovation and New Technology*, 9, 517–557.
- Russo, M. (2005). Innovation processes and competition: China challenges Sassuolo. How do innovation processes in industrial districts are affected by competition from new market systems. In Proceedings of the conference in honor of Professor Sebastiano Brusco, *Clusters, Industrial Districts and Firms: The Challenge of Globalization*, September 12–13. http://www.economia.unimore.it/convegni_seminari/CG_sept03/naviga.html. Accessed 13 November 2007.
- Russo M., and Rossi, F. (2007). Politiche per l'innovazione: dalla valutazione alla progettazione, Dipartimento di Economia Politica dell'Università di Modena e Reggio Emilia, Materiali di Discussione n.565/07.
- Russo M., and Rossi, F. (2009). Cooperation partnerships and innovation. A complex system perspective to the design, management and evaluation of a EU regional innovation policy programme. *Evaluation*, 15/1.
- Spradley, J. P. (1979). *The ethnographic interview*. New York, NY: Harcourt Brace Jovanovich College.
- Triulzi, U. (1999). *Dal Mercato Comune alla Moneta Unica*, Rome, Italy: SEAM.
- Von Hippel, E. (1978). Users as innovators. *Technology Review*, 80(3): 1131–1139.
- Wasserman, S., and Faust, K. (1994). *Social network analysis: Methods and applications*. Cambridge, UK: Cambridge University Press.

Part IV
Modeling Innovation and Social Change

Chapter 12

The Future of Urban Systems: Exploratory Models

Denise Pumain, Lena Sanders, Anne Bretagnolle, Benoît Glisse
and H el ene Mathian

12.1 Introduction

Urban systems are complex systems, mainly because of the non-linear growth processes that lead to very unequal concentration of population and activities in towns and cities over historical time. We have seen in Chapters 7 and 8 of this book that supralinear scaling relationships were a distinctive feature of the structure and dynamics of urban systems. At the end of Chapter 8, we have shown a few examples of trajectories of the weight of individual cities relative to the system they form. These trajectories show inflexions, or even reversals in trend, alternating periods of urban growth and prosperity in cities when an innovation cycle is located there and (at least relative) decline and impoverishment when a former specialization cannot be so successful or even maintained in some urban locations. Because their sustainability depends mainly on the result of their interactions with other places, cities are permanently submitted to the necessity of transforming themselves to improve their position in the system of cities. The competition process for the attraction of innovation among cities, which ends in their unequal development, is a very complex one, because of the multiple interlocked networks that connect all city activities, including linkages between their inhabitants and artifacts. The spatial inter-urban patterns that are generated by so many different kinds of interaction flows and their effects on differential city growth and societal evolution show an incredible variety in shape and magnitude and cannot be predicted from simple analytical models. Only simulation models can give tractable representations for such complex dynamics (Sanders, Pumain, Mathian, Pace-Gu erin, & Bura, 1997; Portugali, 2006).

In principle, the evolution of complex systems is unpredictable (Batty & Torrens, 2002). But simulation models, when correctly calibrated on past evolution, can help to explore issues among possible futures of these systems (Allen, 1997). This is especially possible in the case of urban systems, because their own structural dynamics is obviously slower than the succession of societal innovations that represent their main driving force, and, above all, because they exhibit a very strong *path*

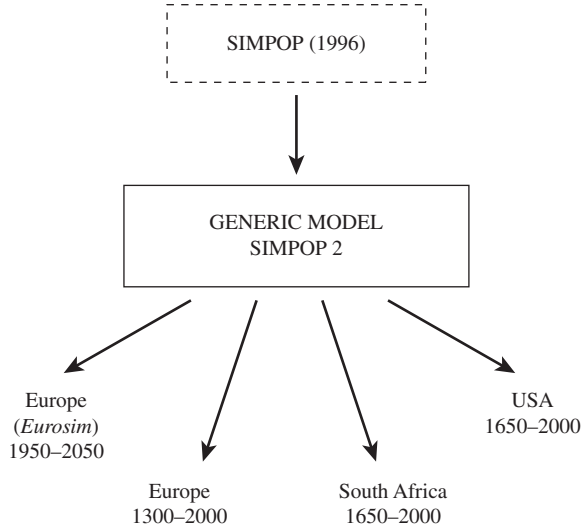
D. Pumain (✉)
Universit e Paris I Panth eon Sorbonne, UFR de G eographie, Paris, France

dependence in their evolution (Arthur, 1994). To succeed as an exploratory tool in predicting the future of a system of cities, a simulation model must satisfy at least two conditions. First, it has to be calibrated to properly represent the past dynamics of the system, since path dependence is a major property of the evolution of integrated urban systems; this raises the question of validation, which is a difficult and uncertain exercise in modeling by simulation, but much confidence can be gained when a generic model is able to reproduce, under reasonable parameterizations and initial conditions, either the effects of accidental exogenous events or the specific features of systems that are observed in different parts of the world. Second, when using the model as an exploratory tool for predicting the future, one has to include a number of probable changes in the demographic, political or technological context of the urban system, which may alter its further evolutionary path, and this will be done here in the case of the European urban system. Of course, this implies an abstract representation of the innovations, whose qualitative nature cannot in any case be predicted by any model.

To translate an urban evolutionary theory of urban systems (Pumain, 2000) into a simulation model, we designed a multi-agents system of interacting cities, in collaboration with the research group of Ferber (1995). The goal of this first prototype, called SIMPOP, was limited to the simulation of a few theoretical principles (Bura, Guérin-Pace, Mathian, Pumain, & Sanders, 1996). The model had a small number of cities (less than 400) and was calibrated roughly on the urban pattern of Southern France. The simulations demonstrated that one could reproduce the historical emergence of an urban hierarchy (over a period of two thousand years) from a set of rural villages (even from a uniform initial condition) *only* when interactions could occur between them (through a market and a competition for the acquisition of urban functions) and if new urban functions (i.e. innovations) were added more or less continuously during the process. The second version of this model, called SIMPOP2, which we present here, is adapted for a larger number of agents (about five thousand towns and cities in Europe). This new simulation tool is a more detailed and powerful multi-level multi-agents system. It is conceived as a generic model which can be applied to different levels of resolution in space and time and to a variety of geographical situations along with specific rules. We have derived four instantiations of this model, one called Eurosim¹ that aims at predicting the evolution of European large cities over the period 1950–2050, while the three others aim at exploring, by data-driven simulations, the past evolution of a variety of systems of cities in Europe (1300–2000), USA (1650–2000), and South Africa (1650–2000). One challenge is to understand which among the rules and parameters of the generic model have to be modified to represent the part of the urban dynamics that is particular to a specific time period or a region of the world (Fig. 12.1).

¹ The instantiation of this model was developed within the framework of another European research program called TiGrESS, (see <http://www.tigress.ac/reports/final/eurosim.pdf> and http://ec.europa.eu/research/environment/newsanddoc/article_2697_en.htm).

Fig. 12.1 Tree of SIMPOP generic and instantiated models



12.2 Urban Complexity and Multi-Agents Modeling

A variety of modeling techniques have been tried for simulating the dynamics of urban systems. Here, we recall only a few steps in this long history (Pumain, 1998). A precursor can be seen in the first Monte Carlo simulation of urban growth from rural and interurban migration flows (R. Morrill's method). A few attempts at using the formalism of catastrophe theory (Casti & Swain, 1975; Wilson, 1981) were not followed up. A first series of dynamic models were expressed as systems of non linear differential equations (Allen & Sanglier, 1979; White, 1978). These models described the evolution of state variables at a macro-level, the lower level interactions being summarized in mathematical relationships or in parameters. As interactions are non linear, the systems are not attracted towards a pre-determined equilibrium, a small change in the parameters of the model can modify the dynamic trajectory of the system and persist as a determinant of their further qualitative structure, according to a bifurcation. For instance, a small change in preference of consumers for large size and diversity of shops and a variation in the price of transportation can produce a spatial concentration of trade in a major urban center or its dispersion in a multitude of small centers. Even if some models made more explicit connections analytically between individual behavior and the resulting aggregated interactions (as for instance the synergetic model of interregional or interurban migrations first developed by Weidlich and Haag (1988) and applied to French cities evolution by Sanders (1992)), in practice there was very limited correspondence established with observations at a micro-level, since an "average" behavior was supposed to be representative of the individuals, and the applications were conducted with statistics on aggregated flows. Conversely, micro-simulation models integrated many details about the behavior and familial or professional career of individuals, but did not

pay so much attention to the evolution of the resulting structures at the macro level (Clarke, 1996; Holm & Sanders, 2007).

Compared to these earlier attempts to achieve self-organization in urban system models, the actual notion of emergent properties refers to a more explicit modeling of interactions, usually in agent based models or in multi-agent systems (Ferber, 1995). Multi-agents systems (MAS) are especially useful as simulation tools for modeling dynamics, when it is essentially explained by the heterogeneity of individual features and their interaction (Sanders, 2007). They enable the modeler to associate qualitative and quantitative rules and to integrate several levels of organization, diverse time scales, and dynamic relationships. They appear as a reasonably promising technique for simulating geographic worlds, mainly because of their ability to consider the environment of a system, their acceptance of a wide conceptual diversity of agents (allowing for multi-level analysis) and their flexibility regarding interaction rules, especially in spatial relationships. Multi-agents systems are much more flexible than differential equations for simulating spatial and evolving interactions, including quantitative and qualitative effects. Through the definition of rules at the individual level, they can reproduce the circulation of information between cognitive and decision making agents. They simulate, at the upper level, the emergence of collective or aggregated structures that can be tested statistically. The rules can be adapted for varying space and time scales of interaction under the course of history.

12.3 Ontology of the Generic Model

Before demonstrating how it can be used to explore a variety of urban evolutions, we describe briefly the general architecture and components of the SIMPOP model, that is, in the computer scientist vocabulary, its “ontology.” We have identified, in Chapter 6 of this book, a number of stylized facts that are common to all urban systems. They characterize, at the macro level of urban systems, major structural features and evolutionary processes (as size differentiation, functional diversity, and distributed growth) that emerge from the interurban competitive interactions. A model should simulate how the cities’ interactions produce these emergent properties. But drastic simplifications are required. Even if the computer’s capacities today allow simulations of the individual daily moves of inhabitants in a city with population of two million (for instance, Eubank et al., 2004), it would be impossible (and probably irrelevant) to represent, at the same detailed individual level, all the interactions between the 320 millions of Europe’s urban citizens distributed among 5,000 different urban agglomerations over several decades! Each model is an abstraction based on generalization. Since multi-agents modeling authorize a rather direct representation of a conceptual model, we decided to make three major abstractions for SIMPOP over the finest grained scale: first, we consider interactions between cities only, not between individuals, so that cities are considered as the “agents”; second, we select among all urban activities those that have a specific role in each city’s dynamics, i.e., their specialized functions, as main city attributes;

third, we retain only, among all real exchange flows that are generated by these functions, the relations that may create asymmetries, that is, “second order” interactions (Pumain, Bretagnolle, & Glisse, 2006).

12.3.1 The System of Cities and its Environment

We define urban systems as subsets of cities that can evolve through their interactions under a few external conditions. The systems we consider are never totally closed, nor can they evolve in a fully endogenous way. They correspond, in geographical theory, to subsets of cities that are submitted to the same general constraints (whatever they are, political or legal control, demographic and economic trends, cultural features, or constraints stemming from the use of the same limited amount of resources – geographers would call this coherent envelope a “territory”) and whose evolutions are interdependent, because of the many connections that link cities together. In the real world, a relevant frame for delineating systems of cities can be a national state (remaining roughly but widely valid during the last two centuries), but it may encompass a continent, or even the whole world in the case of certain “global” or specialized cities. The concrete systems of cities that we have chosen in our application belong to both types, continental or national (Europe, USA and South Africa).

The concept of environment in multi-agents modeling defines the medium of interactions between agents, which corresponds, in our case, to the location of cities (as they are immobile agents) and the societal conditions allowing them to communicate. This backcloth for inter-urban interactions is an evolutionary space, measured in space-time terms, because of progresses in communication techniques. We also describe a few more contextual exogenous components that cannot be generated through the interactions between cities but are necessary for the simulations. At first, there is a subset of elements that define the *initial conditions*, including a “map” locating the cities that are part of the system or that will be activated during its evolution, and their corresponding initial attribute values (especially a lognormal distribution of city sizes). This map can be a random pattern, generated by a stochastic process, or an observed set of geographical locations (Fig. 12.2). The first SIMPOP model focused on the emergence of a hierarchical system of cities and on the progressive structuring of the urban system according to size, functions and spatial proximities, starting from an isotropic spatial organization. The initial situation corresponded then to the almost uniform distribution of settlement pattern that is classically associated with a homogenous agrarian society. For the simulations with SIMPOP2, a larger diversity of initial conditions are necessary, on the one hand to start the simulation from any observed urban pattern, and, on the other hand, to test systematically the impact of different theoretical initial situations. As the simulations are made over longer historical periods, alternative theoretical patterns can also be used to control the possible effects of a geographical and political configuration on the general urban dynamics. Figure 12.2, for instance, represents

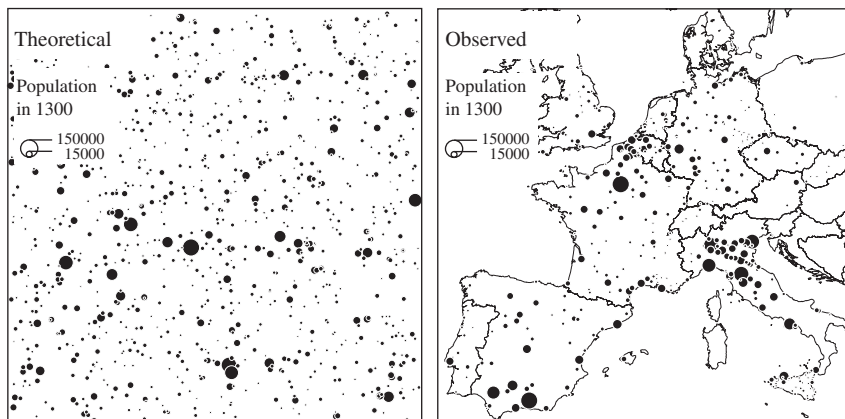


Fig. 12.2 Examples of initial urban patterns: theoretical and observed

two initial patterns for Europe in 1300: one corresponds to the observed situation; the other to a fictitious spatial distribution (a triangular grid) reproducing the numeric properties of the European urban system at that date. Population sizes are distributed spatially in a random way according to a lognormal distribution of same mean and standard deviation as observed. In the case of the Eurosim model, which represents the dynamics of the European system of cities for the period 1950–2050, the initial situation, urban Europe in 1950, corresponds to an observed and already well-shaped urban system, which has attained some maturity after a millennium of urban development. In that case, the inherited *form* of the urban system influences its future evolution, but new elements and configurations can emerge as a result of the way the cities catch (or not) new innovations, and the model is used to simulate the corresponding changes in cities' relative positions.

Another subset of exogenous data that complete the contextual “environment” of our system of cities are the number and type of the entities called “urban functions” (see below, Table 12.1) that represent innovation bundles (for production or services) with their historical date of appearance in the evolution of the system, and their associated attribute values. A third subset of external elements are parameters that describe the general demographic or economic growth of the society and the period under consideration. These variable parameters (in very limited number) are essential for calibrating the model; they have to be adapted to the historical and geographical context of each specific instantiation.

12.3.2 *Cities as Collective Agents*

In our epistemological framework, the city is not considered as a collection of individuals and enterprises whose simple aggregation would permit us to understand and reproduce the city's evolution. Rather, the city is seen as a complex entity that

Table 12.1 Urban functions and their dates of activation in SIMPOP2's instantiated models

		Eurosims:			
		Europe (1950–2000)	Europe (1300–2000)	USA (1650–2000)	South Africa (1650–2000)
Central functions (proximity principle)	Central 1		1300	1650	1650
	Central 2		1300	1800	1800
	Central 3	1950	1800	1850	1900
	Central 4	1950	1900	1900	1960
Territorial functions (political principle)	Territorial 1 (regional capital)	1950	1300	1800	1900
	Territorial 2 (Capital)	1950 (1990 for Berlin)	1500 15 cities	1800 1 city	1900 1 city
Specialized functions (Network principle)	Long distance trade		1300	1650 (east); 1860 (west)	1650
	Manufacturing 1 (Industrial Revolution)	1950	1800 Industrial Revolution	1830 Industrial Revolution	1860 Industry gold/diamond
	Manufacturing 2 (Electricity, automobile)	1950	1900 Electricity- automobile	1880 Oil, electricity- automobile	1930 Electricity- automobile
	Technopoles (NTIC)	1 st wave in 1950; 2 nd wave in 2000	1950	1940	1960
	Finance	1950			
	Tourism	1950			
	Hub (transport)	1950			
	NBIC	1990			

makes sense as a whole, characterized by its attributes and following some rules of evolution. This conception has led us to adopt a unique approach in the field of multi-agents systems (MAS) where most of the applications concerning social sciences are developed at the level of the individuals, with the idea of analyzing and understanding the structures that emerge at a higher level of observation from the interactions between their actions. In this respect, we agree with Openshaw (1997) that the lowest level of observation is not always the best one from a conceptual point of view. It is also important to state that there are not only two levels of interest, that of the individuals (micro-level) and that of the society (macro-level), but there are a whole set of intermediate levels of interest, including the cities.

While identifying a city as agent, our hypothesis is that the grounds for differentiated growth are better understood at that level than at the level of the households or of the individual political and economic actors. Of course the decision-making processes of each actor in the city has an impact on the city's development, but whatever this impact, it is *limited* compared to mechanisms of change, which obey meso-geographical regularities. Indeed, individuals' actions are of little weight with

respect to a city's trajectory in the long run, whose trend is not very sensitive to the diversity of individual intentions and decisions. These meso-geographical laws determine the context, in terms of possibilities and constraints, in which actors of different kinds and levels make their decisions. They determine the "bounds" of the possible future for a city, given its characteristics (size, accessibility, socio-economic profile, specialization. . .). These bounds can be interpreted as local attractors in the dynamic trajectory of a single city. The decisions of political or economic actors thus influence the direction of the trajectory towards one "bound" or the other. In other words, the meso-level dynamics give an "interval of plausibility," and the urban actors' decisions and actions determine *where* in this "interval" the change will occur.

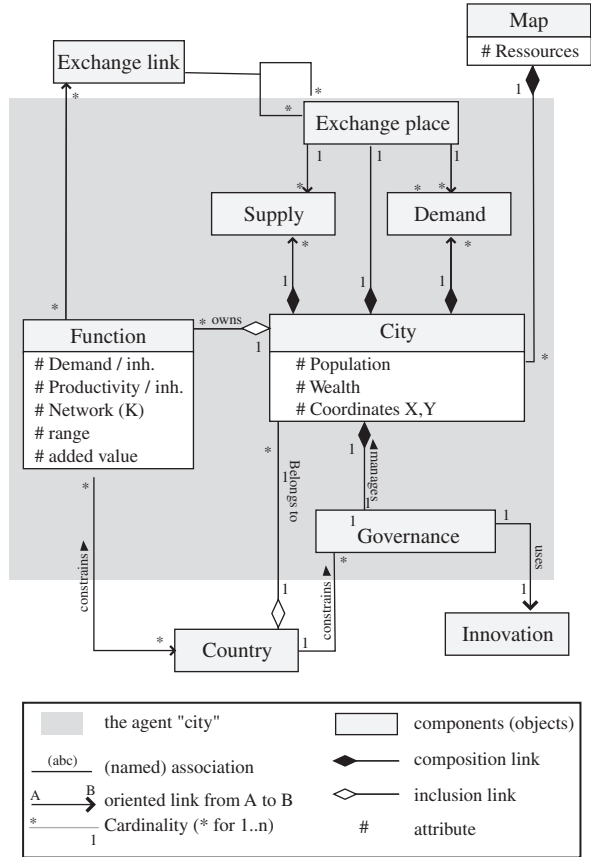
So the elementary entities of the model are cities, each represented by an "agent" in the terminology of multi-agent systems. This agent has a certain degree of autonomy as far as its decision making process is concerned; it handles information about itself, about the properties of the cities with which it is interacting, and the rules of evolution. It is able to communicate (through a set of interaction patterns called, in technical terms, a protocol of communication) with the agents representing the other cities. Through a collective entity called "governance" that represents the decisional capacity of urban actors, it acts not only as a reactive agent but can develop different types of strategies, including a more or less risky approach to the acquisition of new functions (investment). Urban functions are introduced exogenously during the simulation at given dates corresponding to the major innovation cycles (in an entity called "innovation," which is part of the context of the urban system). They can be attributed to cities in a passive way, according to a set of criteria that determine their allocation, or through the governance entity according to strategies of imitation (of neighboring or well-connected cities or cities that are similar in size or economic profile) or strategies oriented toward risk and innovation.

As our aim is to model the structure that emerges from the interactions between cities, our hypothesis is that the *interactions* between cities are the driving force in the evolution, and they determine the future of each city as well as the evolution of the macroscopic properties of the system of cities. These interactions stem from migration flows, commercial trade, information flows, knowledge exchange, etc. They determine how innovations will spread throughout the system. The hypothesis is that this is relevant for different types of territories and at different periods of time, for which towns and cities exist. That justifies the elaboration of a generic model, based on this conceptual framework and containing the common rules of evolution of towns and cities embedded in a system of cities (Fig. 12.3).

12.3.3 Main Attributes are Urban Functions

The simplest attribute of a city is a *location*, i.e. its geographical site. In our model, a map (see Fig. 12.3) is part of the initial condition, including many possible locations (given by two coordinates) for cities. During the simulation, new cities are activated

Fig. 12.3 Representation of SIMPOP2's ontology



at a rate corresponding to historical observations. The rules for this activation can vary according to the historical type of settlement in a country: in old urban systems, new cities arise randomly between the existing ones; while in countries of the new world, cities are generated along frontier lines. Their *relative location*, that allows cities to be more or less easily connected to others, may introduce a potential differentiation in their evolution. Accessible environmental resources, such as coast lines or natural corridors (allowing for the acquisition of trade functions) or mineral deposit zones that are also located on this map as exogenous information, are also part of the city attributes and can be useful for introducing the specific effects of resource-dependending innovations.

Each city is characterized on the one hand by its *population* and its *wealth*, which constitute the main state variables of the model; and, on the other hand, by its *functions* and the *distribution of its labor force* according to the associated economic activities. Two synthetic attributes, the *total population size* and the *accumulated wealth*, represent the strength of each city-agent in the system of cities. They summarize the past ability of the city to attract population and benefit from its

exchanges with other cities. These attributes depend, modulo stochastic factors, on another set of attributes, which is the portfolio of functions that define the capacity of action for each city in the system. An urban function represents the subset of urban activities which may generate asymmetries in the interactions between a city and other cities. The principle is taken from *economic base theory*: the potential of growth of a city, i.e. its ability to attract new population and new activities, depends on its ability to produce and to valorize “fundamental” productions or services that are export-oriented.² The core of the SIMPOP2 model concerns the formalization of the exchange market between cities. In short, the growth of each city will depend on the success of its exchanges. The type and level of function can evolve through time depending on the ability of the city to adopt innovations and to become attractive for some dynamic and leading activities.

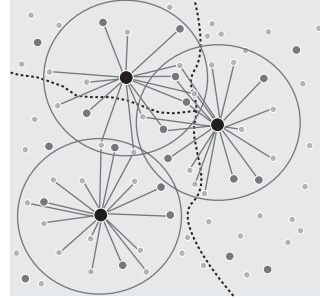
Three families of functions are distinguished, each one corresponding to a different principle of spatial interaction (Fig. 12.4 and Table 12.1):

1. *The central functions* generate interactions according to a *principle of proximity* (or *gravity principle*), following a spatial interaction model that is currently described by equations similar to Newton’s formula and included in Christaller’s central place theory (Christaller, 1933). They include the most classical urban activities – commercial activities as well as services and some manufacturing industries, whose production is intended for a regional market, i.e. neighboring towns and cities. The spatial interaction principle is the same whatever the spatio-temporal context, but the associated ranges of influence vary according to this context. More specifically, there are four possible levels of central functions (central 1–4) that emerge successively during historical times; they are hierarchized according to the complexity of the services they perform and to the spatial range of their influence. Figure 12.4A illustrates the spatial operating of this family of functions and how the competition between cities occurs through the overlapping of their zones of influence.
2. *The territorial functions* include the administrative activities that operate within the frame of political or administrative boundaries. They include two levels, the specific functions of the capital of a national territory or that of a regional capital (Fig. 12.4B). The administrative services that are produced supply the demand of all cities and towns of the corresponding region or country only: there is no competition across the boundaries.
3. *The network functions* consist of very specialized activities that were created by major economic cycles with a large range of trade; their development depends on the relative position of the city in a system of specialized trading relationships (Fig. 12.4C) rather than on spatial proximities. These functions are of different kinds, according to the main economic cycles that created major urban specialization, as described in Chapter 6, and their operation obeys different rules depending on the cycle to which they belong. These rules express the probability

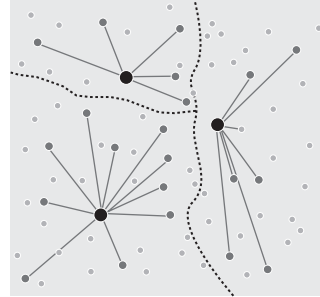
² Fundamental activities are classically opposed to induced activities, which are oriented towards the local urban or regional market.

Fig. 12.4 Three types of spatial interactions

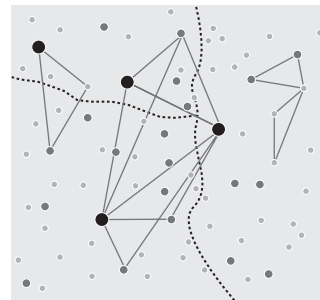
Gravity principle
for central functions



Territorial principle
for administrative functions



Network principle
for specialised functions



of different pairs of cities to be connected through trade, given their respective socio-economic profile (relations of complementarities for example, see below in Table 12.4).

The list of urban functions may vary according to the instantiated model. For instance, Eurosim³ establishes more categories among recent urban functions (Tourism, Hub, Finance and NBIC), while ignoring maritime trade, which remain essential for the simulation of European urban system over more ancient periods of time. It also has less distinct levels of central functions, and uses NBIC (converging nano, bio, information and cognition technologies) as the expected leading innovation cycle for exploring the next decades of urban specialization.

³ Knowledge about the functional evolution of European cities was used to control the hypotheses investigated (Cattan, Pumain, Rozenblat, & Saint-Julien, 1994; Cicille & Rozenblat, 2004).

12.3.4 Variables and Parameters

The aim of a simulation can be better understood by identifying essential variables of the model (Table 12.2). As seen above, the *state variables* of the model are the population, wealth of cities, and amount of labor force in each of their functions. They are computed at the city level and later aggregated to the level of the system of cities as a whole when the simulations are analyzed. But, other parameters have to be introduced to characterize the evolutionary context of urban dynamics. These *contextual variables* correspond to a higher level of interactions compared to those simulated in the model. They are given exogenously. They enable a calibration of the model on observed historical data. They define a mean growth rate for population and wealth (or economic product) and, for each type of function, their date of emergence, productivity level, demand level, and increase in value (profit gained from trade in addition to the actual amount exchanged). Productivity and demand could not always be retrieved from actual data and were sometimes first estimated and then tuned by calibration.

Intermediate variables are computed during the simulation; they characterize the dynamics of trade between cities (as unsold goods, or size of the customer network) – see Section 12.4, below.

Key parameters are decisive for calibrating the model. They keep the same value for every city but can evolve through time. Their value cannot be observed and has to be determined by trial and error, but has to remain in some domain of validity (for instance corresponding to a plausible historical succession of values or a logical cross-comparison of urban functions). This is a rather long and difficult process, since the effects of these parameters interact with one another. The simulations enable the modeler to check the sensitivity of the model to the variations

Table 12.2 Main variables and parameters in SIMPOP2

Category	Parameters
State variables	Population, Wealth, Labor force by urban function
Contextual variables (exogenously defined for each urban system)	Population and wealth: mean growth rates Date of emergence of each function Productivity, demand, added value, for each function
Intermediate variables (endogenous)	Unsold goods, Unsatisfied demand Size of the networks
Key parameters (calibrated)	Range of exchanges associated to the different functions Size of exchange networks for specialized cities Attraction level on labor force % of valuable customers Returns from the market on urban growth Barrier effects of boundaries

of these parameters and to choose the values that give the best fit to observed data. Six parameters are used, among which three (described below) modulate the city population growth according to the results of exchanges through urban functions (Table 12.2). The other three key parameters (range of exchanges, size of exchange networks and barrier effects) will be described below with the rules that use them.

1. *Share of growth not generated by the model.* This parameter helps modulate the mean growth rate that is introduced in the model as an exogenous demographic trend and that cannot be injected in the model as such, since interactions between cities also generate a significant share of urban growth. Remember that this mean growth rate that is allocated to each city represents an exogenous conjuncture, which reflects the historical trends in demographic evolution.
2. *Return from market on demographic growth.* This is a feedback effect from economy on demography, which is linked to the balance of exchanges of cities for each function. This balance leads to a wealth increase or decrease. This parameter varies according to the profit that can be generated by each urban function. It enters in the equation where the wealth of cities is computed at the end of a cycle of exchange and influences the city's demographic growth (Table 12.5, Equation 12.3).
3. *Attraction on labor force.* Agents can adapt their labor force to the economic context that is perceived by cities through their acquaintances, by increasing or diminishing the number of employed in a given function. This translates to population growth or decline. When demand exceeds supply, this labor force is increased; while, on the contrary, it is decreased if there are unsold goods and services.

12.4 Rules of the Model

The various instantiations of the SIMPOP2 model share a number of similarities that define the SIMPOP2 paradigm (see the generic model introduced in Sections 12.1 and 12.2). In addition, they specialize to address a specific problem or context. They each deal with a defined case study as well as precise objectives. We describe here the main rules that are given to Eurosim and SIMPOP2 Europe and relate them to the scopes of the two models as well as the scope of the overall SIMPOP2 paradigm.

12.4.1 Time and Interactions

Due to the computational nature of Multi-Agent Based Simulation, the evolution of state variables is, by nature, discrete. Though, Multi-Agent Systems rely on two scheduling techniques (Michel, Ferber, & Gutknecht, 2001; Ramat, 2007):⁴

⁴ The first technique is discrete time: the simulation is divided into iterations within each the agents that are activated. The scheduler is like a clock that can translate the virtual time of the model, the

SIMPOP2 uses first a discrete time approach and then an event based one to model the numerous round trips of the interactions occurring within a time step. The flow of time is divided according to the temporal resolution of the models. For Eurosim, which covers periods from 1950 to 2050, iteration is made each year. For SIMPOP2 Europe, which runs over a larger period, from 1300 to 2000, an iteration represents 10 years. This is a modeling choice that enables simple comparisons between the simulation results and observations from the real world. Also, the more the system of cities advances in time, the more the interactions become overwhelming. It is then reasonable to introduce more iterations, which can possibly capture quicker transformations of the system.

As explained in Section 12.1, the interactions are the driving force of the model. The model assumes that, during each iteration, the city-agents fulfill their trade exchanges to the end. That is, trade occurs till no more can be done (the round trip effect mentioned previously). This assumption is reasonable, considering the length of the corresponding period in real time. Within an iteration, the agents are then enabled to trade in a recursive fashion. This can be interpreted as an event driven process: an unsold supply tries to meet an unsatisfied demand and vice-versa.

12.4.2 *Detail of an Iteration*

Figure 12.5 gives a synthetic representation of an iteration in SIMPOP2 that we shall now detail. At first, the amounts produced by each city and the demand expressed by its population are computed. This is made by using contextual parameters that define the productivity and demand at individual levels for each type of goods and services provided by each urban function at a given time.

The *computation of the trade networks* is then made for each urban function, and may vary according to the specific models. Different constraints are imposed on the topology of the trade networks according to the three spatial types of interactions (Fig. 12.4). The transportation improvements are visible, for the SIMPOP2 Europe model, in the central functions that appear one after another in the system (Table 12.1). They are also represented in the values of a parameter that expresses the maximal possible range for trade from any city (evolving through time and according to city size). Moreover, we introduce two key parameters that help estimate the potential market of a city by fixing a maximum size to trade networks and by allocating a proportion of trade devoted to previous customers, the remaining part being reallocated randomly each time. The values taken by these parameters are important for regulating the stability of the model, which is highly sensitive to them.

Table 12.3 summarizes a few rules for generating possible trade relationships between cities located at different distances or having different demands according to the functions they have.

iterations, to real time. The second technique is event based: the scheduler executes events one after another given their chronological order (events possess a timestamp). The events and their order can be known from the beginning or they can be created by the agents during the simulation.

Fig. 12.5 Steps of an iteration

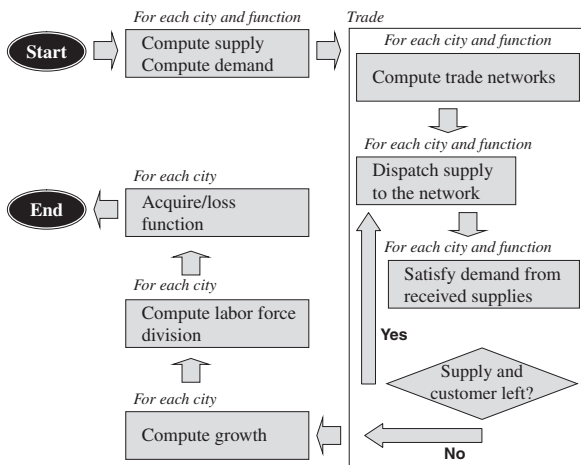
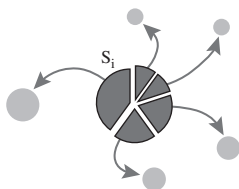


Table 12.3 Cities involved in trade networks by urban function

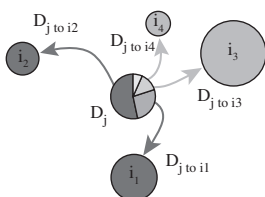
Urban Function	Type	SIMPOP2 Europe	Eurosim
Central1	Intern	–	–
Central2	Proximity	<70km	From 150 to 250 km
Central3	Proximity	<80km	n/a
Central4	Proximity	Before 1950, within 100km; 150 km after	n/a
Regional capital	Territory	–	n/a
Capital	Territory	–	–
Long distance trade	Network	Central 2 cities*	n/a
Manufacturing I	Network	Manufacturing I cities*	Manufacturing I or II
Manufacturing II	Network	Manufacturing cities*	Manufacturing cities
Tourism1	Network	n/a	All cities within the country and neighbouring countries
Tourism2	Network	n/a	Among largest and richest cities
Hub (transport)	Network	n/a	Largest cities
International Finance	Network	n/a	ManufacturingII, Tourism2, Hub, Technopole, NBIC or Capital cities
Technopole	Network	Central 4 or Manufacturing II cities*	Manufacturing II, techno, NBIC, capital or largest cities
NBIC	Network	n/a	NBIC and cities above one million of inhabitants

*there is a range constraint as well, but after 1800 almost all Europe is covered

Fig. 12.6 Rules of market exchanges

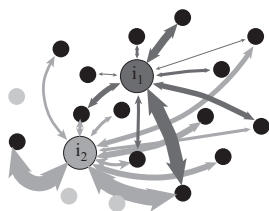


Supply :
each producing city (i) supplies its production (S_i) among its network of potential customers according to their demand or their distance .



Demand :
each city (j) shares its whole demand (D_j) among the suppliers i according to specific criteria (technological level, costs or previous links)

Fig. 12.7 Interactions resulting from market exchange. Line thickness indicates quantity



Effective exchange:
each city with demand D_j , buys from the city i the quantity :
$$\min (D_j \text{ to } i, S_{i \text{ to } j})$$

wealthier cities are served first. Simultaneously i and j exchange goods and wealth.

The resolution of trade networks involves a supply dispatching from the producers to its network of consumers and the actual “purchase.” Figure 12.6 shows how market exchanges are managed. The preference for Eurosim is based on low and high costs to be economically consistent. The wealth per inhabitant is used as a proxy for the level of wages and production costs. SIMPOP2 Europe considers distance and connectivity between seller and buyer. This is a trade off to consider both the impact of transportation in early ages and globalization at later stages.

The trades can lead to a complex topology. The networks can overlap and, thus, competition occurs. Figure 12.7 gives an example of two cities that try to sell the same production of a Eurosim urban function. Round trips are then necessary to dispatch the productions, because a perfect first match between supply and demand for a particular “supplier” or a particular “buyer” is unlikely.

At the end of the exchange process, the growth of the city-agent is computed by taking into account three elements, as detailed in Table 12.4. First, a positive or negative feedback from the trade exchanges is given, depending on their success (measured by unsold products or unsatisfied demand). This is translated into an *attraction on labor force* which increases or decreases the amount of employees within the corresponding urban function. Second, the demographic growth impulse given to the city from the general *demographic trend* of the country is computed. Third, a positive feedback is made from the wealth to the population growth through the *market return parameter*: urban functions generating high added values will bring a greater benefit to the city.

Table 12.4 Assessment of economical exchanges and computation of population variation at the end of an iteration

Total population variation at the end of an iteration

The total population at $t + 1$ is function of 3 components:

$$P_i^{t+1} = (P_i^t + \Delta^1 P_i^t)^*(1 + \partial_1 + \partial_2)$$

(1) Evaluation of labor force attractivity for each urban function k of city i :

First the variation of labor force between t and $t + 1$ for the sector k is evaluated by $P_{ik}^{t \rightarrow t+1} = s_{ik}^t \text{Pot} M_{ik}^t$, where $P_{ik}^{t \rightarrow t+1}$ designs the variation of workers in sector k , based on the potential of the trade network of this sector for the city i (PotMtik). The potential compares the demand of customers for the sector k , to the supply of the city at time t after the trading process. If there are unsold goods, the potential will be negative, conversely if there is unsatisfied demand, it will be positive.

s_{ik}^t is a parameter whose value follows a normal distribution $N(m_s, \sigma_s)$. For the short term simulations it may be interpreted as the “speed of adjustment.”

The variation of total active population due to the market adjustment is then given by:

$$\Delta^1 P_i^t = \sum_k P_{ik}^{t \rightarrow t+1} \quad (12.1)$$

(2) Demographic trend:

The second part of the evolution depends on the general demographic trend weighted by systemic effects:

$$\partial_1 = \alpha^{t*} G_h^t \quad (12.2)$$

where G_h^t is the global demographic trend observed at time t in the region h and α^t a parameter that evolves over time between 0 and 1. This global trend may vary over regions (according to differences in stages of demographic transition). α^t represents the share of growth that is not generated by the model.

(3) Market returns:

The third part of the evolution of the population depends on the city wealth increase:

$$\partial_2 = \beta^{t*} f(\Delta w_i^t) \quad (12.3)$$

w_i^t is the wealth of city i at time t , $f(\Delta w_i^t)$ estimates the balance of wealth of the city i between t and $t + 1$, β^t is the weight given to this third component of the population variation at time t . If there is no wealth increase, there is no effect of market return on city population growth.

When growth is computed, the city-agents update their labor force. The increase or decrease is made according to the intermediate variation determined in step 1 of the growth computation (Equation 12.1 in Table 12.4). A final value is set after respecting the constraint that the total active population represents about 45% of the total population (by normalizing). The city-agents can lose one urban function when their labor force reaches 0.

Regarding the rules for acquiring new urban function, they are given exogenously and may vary according to the model. In the Eurosim model, given the relatively short period of time, the allocation of functions is made *a priori* within

Table 12.5 Rules for the adoption of new urban functions in the SIMPOP2 Europe model

Urban function	Dates	Emerges among:
Central 1	1300	–
Central 2	1300	largest Central 1
Central 3	1800	largest Central 2
Central 4	1900	largest Central 3
Regional capital	1300	largest Central 2, minimum spacing
Capital	1500	largest Chieftown, minimum spacing
Long distance trade	1300–1800	largest Central 2, preferential locations (zones)
Manufacturing 1	1800	largest Central 2 or preferential locations
Manufacturing 2	1900	Central 3 or manufacturing 1
Technopole	1950	Central 4 or manufacturing 2

the initial situation in 1950 and the locations of the latest emerging function (corresponding to the NBIC specialization) are selected from this date, even if activated in 2000 only. The SIMPOP2 Europe model uses rules that are summarized in Table 12.5. The conditions that are requested reflect the most frequently observed transition in urban specializations as they were observed in the history of European cities.

12.5 A Multiscalar Method of Validation

In our simulation model, there is no optimizing constraint, and the evolution is open as in any exploratory simulation tool. However, we need some validation procedure to assess both its ability to reproduce past observations reasonably well or to evaluate the magnitude of deviations between a diversity of future scenarios.

To calibrate and validate the model, as well as to valorize the different results of the simulations, a multiscalar tool of “simulation data outputs mining” has been developed. The main objective is to test the coherence of the rules introduced in the model, its ability to produce trajectories that are realistic according to observations, and to get an insight in its sensitivity to initial conditions and parameter variations. But the potential outputs of the simulations represent a huge amount of data. For instance, in the case of Eurosim, the output for one simulation relates to 5,000 towns and cities, 100 time steps, 13 functions, and all interactions associated to communications and exchanges between the cities. Therefore a visualization and exploration tool for analyzing the outputs has been developed. This method investigates a data hypercube including three conceptual dimensions: time, state variables, and space, including interaction flows and multi-level organization. In addition, the calibration process includes methods for identifying which parameters influence the dynamics of the quantities and which influence the structures themselves (as urban hierarchies or configuration of exchange networks).

Thus, rather than computing only one objective function that would be a one-dimensional summary of the simulated values, we define a multidimensional framework for evaluating the simulation. This framework includes measurements on the

structural features of the urban systems: hierarchical, spatial, and functional, and this is evaluated at different scales. Different investigations are also introduced according to outputs: for instance, some, as population, wealth, labor force, are of interest from a thematic point of view while others are only used to check the coherence of the model and to help calibrating it.

The outputs are highlighted and summarized according to three entries corresponding to different geographical levels: the macro-level (European urban system for Eurosim or SIMPOP2-Europe) is the aggregate level; the local level concerns the cities themselves; between these two extremes, there are a series of intermediate results and outputs, corresponding to different kinds of geographical aggregation of the results from the city level:

- territorial: a regionalization can be defined, such as Eastern and Western regions for Europe, or national states;
- hierarchical: grouping cities by size class can help in detecting if there is a systematic size effect in the cities' evolution; and
- functional: grouping cities with the same specializations or with the same number of specializations to illustrate the effect of different functional levels on differential growth.

The number of variables to be represented, the use of different aggregations, and the use of different methodological filters produce a large amount of outputs. We use the complementarities between different ways of analyzing the outputs to get a complete overview of the different states of the urban system during the successive periods. This also enables a multi-dimensional comparison of the different simulations, which then facilitate the calibration of the model.

The outputs are analyzed through a series of methodological filters, as shown in Fig. 12.8. Thus, a standard report will include the description of the three geographical levels.

The macro level (global urban system) is described through:

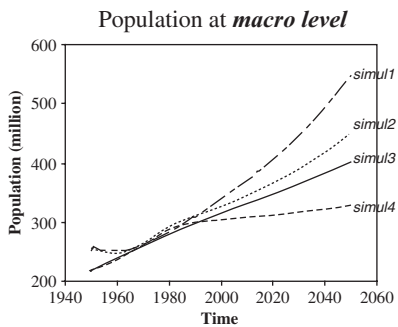
- global trajectories: population, wealth, variation rates of population and wealth, repartition of the labor force by activity sectors;
- hierarchical structure: rank size representation and modeling over time for population and wealth, primacy evolution; and
- spatial structure: global maps.

Intermediate levels (regions, subgroups of specialized cities) are analyzed with:

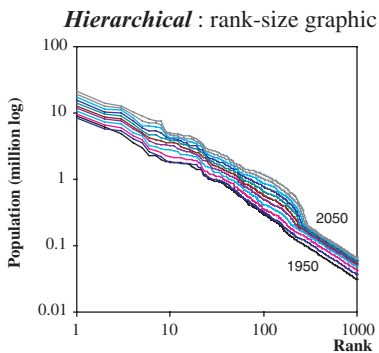
- demographical trends by regions, by family of specialized cities;
- evolution of the global exchanges by sector of activities; and
- decomposition of the different components of supply and demand by sector of activities.

And the local levels (cities) are illustrated through:

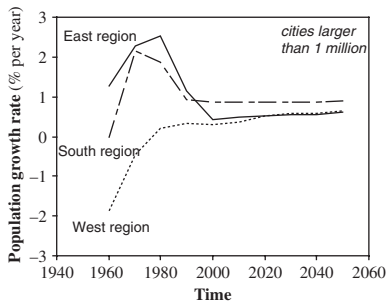
(1) Trajectories at different levels for state variables



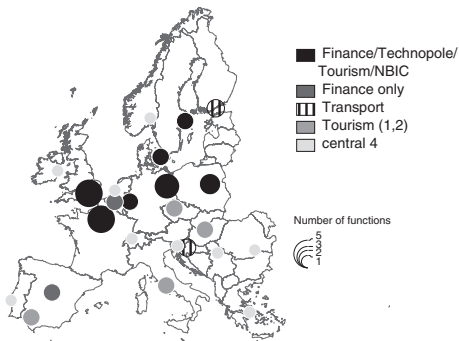
(2) Structures



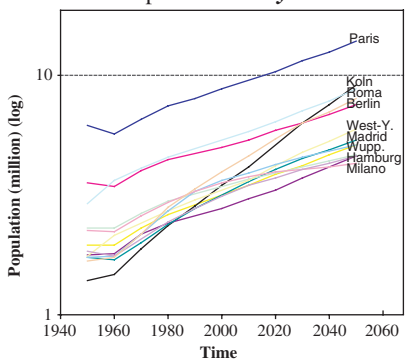
Growth rate at *meso level*



***Spatial* : location of urban functions**



Population at *city level*



***Interactions* : network of exchanges**



Fig. 12.8 An integrated set of outputs for validation

- local trajectories for all cities attributes: population, wealth, labor force, growth rates. . . and
- spatial structures of exchange: maps of the market networks, maps of exchanges flows, evolution of the size of the networks.

Depending on the steps of the calibration and validation phases, specialized reports may be edited for insights on specific dimensions. Multiple checks guaranty the coherence of the calibration, whatever the scale.

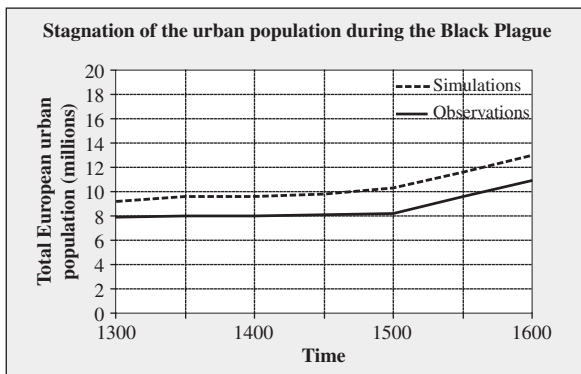
12.6 Results of Simulation

The assessment of SIMPOP2 abilities as an exploratory tool for the future of urban system is still in progress (Pumain et al., 2006), and further research is described on SIMPOP's website (<http://www.simpop.parisgeo.cnrs.fr/>). Meanwhile, we have selected a few results that seem of importance for validating our simulation approach. From the available experiments with the model, we present three examples illustrating its ability to represent urban dynamics in a consistent way and, then, discuss how it can be used in designing scenarios for the future of urban systems.

12.6.1 Simulating the Resilience of Urban Systems After External Perturbations

The model is flexible enough to reproduce huge variations in state variables at macro level as well as in many individual trajectories of cities that happened in the long history of urban systems. Some examples of such catastrophic political events that cannot be embedded in the “normal” evolutionary process of an urban system are the momentary recessions due to wars (observable during Napoleonic Wars at the beginning of 19th century, or the world wars 1914–1918 and 1939–1945). Another interesting “random” historical event is the Black Plague starting in 1348, which was followed by a period of urban population decline in huge but variable proportions from 20 to 50% according to the European regions. We tried to simulate such a catastrophic event with the model, at first by replacing the value of the demographic trend parameter that had been smoothed over the whole period by zero growth during the five last decades of the 14th century. The model reacted well, proving its sensitivity, but was not able to simulate the totality of the sharp decrease observed. It is only by also reducing the intensity of trade exchanges during that period (modifying the parameters of market return and attraction on labor force), that we were able to reproduce the observed population decrease during the 1350–1400 period and its rapid recovery during the following fifty years (Fig. 12.9). The time of reaction of the model to a change in key parameter values is not too long and the model can thus be used for analyzing the effects of societal events or changes in urban practices that occur on medium time scales (a few decades).

Fig. 12.9 Simulating system resilience after external choc



12.6.2 Global Cities Since the Middle Age?

While calibrating the SIMPOP2 model, we discovered that it was impossible to reproduce the size of a few cities whose observed populations were incommensurably large, compared to the results of the simulations. This happened at all periods in the evolution of the European urban system, since the Middle Age. While the population of large cities generally is simulated correctly below a given rank (for instance the second in 1500, the fifth in 1700, see Table 12.6), the population of the cities having a higher rank in the urban hierarchy is under-estimated in a systematic way. For example, in 1500, the size of the largest city (Paris) is 225,000, according to observation, whereas it is only 149,000 in the simulations, despite all our efforts to improve the fitness of the calibration by modifying all key parameters. Thus, the

Table 12.6 Observed and simulated sizes (in thousands) of the largest cities

Dates	Observations*	Simulations *	Rank of the next well fitted city
1500	Paris, 225	149	2
	London, 575	193	
1700	Paris, 575	189	5
	Naples, 500	182	
	Amsterdam, 200	162	
1800	London, 948	254	4
	Paris, 550	243	
	Naples, 430	239	
1950	London, 8900	2932	5
	Paris, 6200	2865	
	Ruhr, 4100	2733	
	Berlin, 3500	2439	
2000	Paris, 10500	6995	3
	London, 9200	6976	

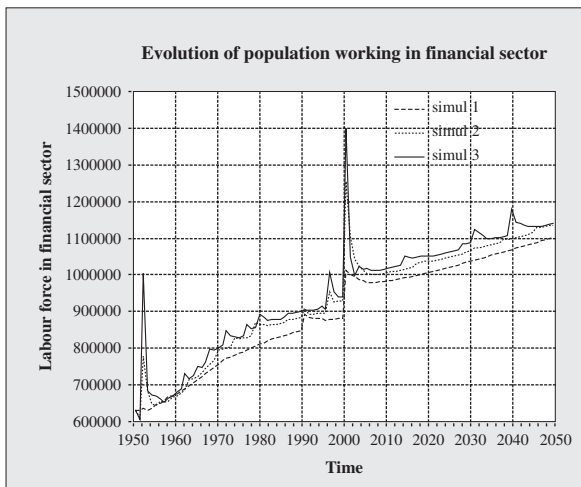
Source: Bairoch et al. 1987, Geopolis 1994, Géographie-cités 2000.

functions that are present in the model are not powerful enough for generating such extra large sizes, even when they are accumulated altogether in one city. According to the period, there were between one and five such “too large” cities. A hypothesis that is suggested by other historical observations is that these cities have in common a further emerging property, stemming both from their actual combination of functions as state capital, node in maritime trade, or focus of industrial activities, and from their exceptionally central position in enlarging exchange networks progressively out passing the frame of continental Europe as well. If such exceptional urban situations were identified by Fernand Braudel as centers of “world economy,” we could add another dimension by identifying centers of “world politics” (which could be a complementary explanation in the cases of Paris and Naples for instance) and summarize these two functions under the label of “world city function” (operating at a wider scale than the considered system). Its implementation within the model will help us to take into account the exceptional trajectories of cities like Paris at the head of the first large nation state since 16th century, the role of cities like London and Amsterdam in the Atlantic maritime trade in 17th century, the function of empire capitals of Paris and London during the colonialism period in 19th century and later on in industrial networks or financial activities, that are recognized today as symptomatic of “global cities”. Thus, the model suggests that the function “world city” is by no means an innovation of the last decades of 20th century! Further, the European system has to be considered in co-evolution with the rest of the world, earlier than expected, through these world cities that act as inter-systems gates (which are multiplying nowadays due to globalization).

12.6.3 Reaction of the Urban Systems to an Innovation of the 20th Century

Another example of resilience and adaptation of our model of urban systems was provided while testing the sensitivity of the Eurosim model to exogenous events. As mentioned in Table 12.1, the function named “Technopole” is acquired at different dates by specialized cities: a few large ones own it since 1950, and, according to the principle of hierarchical diffusion, a few medium sized cities (seven cities of 200,000 to 1 million inhabitants) acquire it in 2000. The simultaneous acquisition of this specialization by so many cities catching this new urban function, introduced a strong perturbation in the system, as illustrated in Fig. 12.10. The curves represent the evolution of the number of employees in a completely different urban function, the financial one, for different simulations corresponding to three slightly different configurations of parameters used during the test (Sanders, Favaro, Glisse, Mathian, & Pumain, 2006). The evolutions are globally similar, including a large peak in year 2000, with only small differences in the intensity and timing of shorter fluctuations. This result illustrates the effectiveness of interactions between the different functions: due to the requirement of new capital funds for investing in the technological innovations, the cities owing the “*technopole*” function are among those that have the strongest demand for the finance sector. Seven new cities expressing a demand

Fig. 12.10 Perturbation and adaptation after an innovation



result thus in an imbalance between supply and demand and the existence of an important potential of unsatisfied demand. The rule expressing the return from the market on the labor force quite naturally leads to an important increase in the finance sector employment for each supplying city. This increase is particularly important when the parameter measuring the speed of adjustment is high, which is the case for the two curves which register the highest peaks on Fig. 12.10. Moreover, this reaction shows the ability of the system to integrate sudden change, as the three curves representing the labor force in the finance sector recover their previous growth trend only two periods later. The shock is integrated, the effect of the newcomers is diluted, and the urban system has shown its fundamental resilience and adaptive capacity.

12.6.4 Predicting Future Trajectories for Individual Cities

The Eurosime model has also been used to test different scenarios on the evolution of the European cities during the 1950–2050 period (Sanders et al., 2007). To investigate what kind of consequences different contexts could generate in terms of urban structure as well as of individual cities' evolution, scenarios concerning possible future economic and demographic environments were imagined. Four extreme situations have been defined by combining hypotheses relative to the evolution of the demographic context on the one hand and to policies in matter of intra-European exchanges on the other hand:

- Two demographical alternatives are defined using IIASA's predictions: they correspond respectively to an hypothesis of *high demographic growth* (IIASA'S more optimistic predictions which means a very slightly positive growth rate); and of *low demographic growth* (IIASA'S more pessimistic previsions which means a clear negative trend for all Europe);

- Two political alternatives are introduced concerning the presence or absence of *barrier effects* between Eastern and Western Europe. The existence of barriers will reduce the possibilities of exchanges between cities located in the two geographical regions.

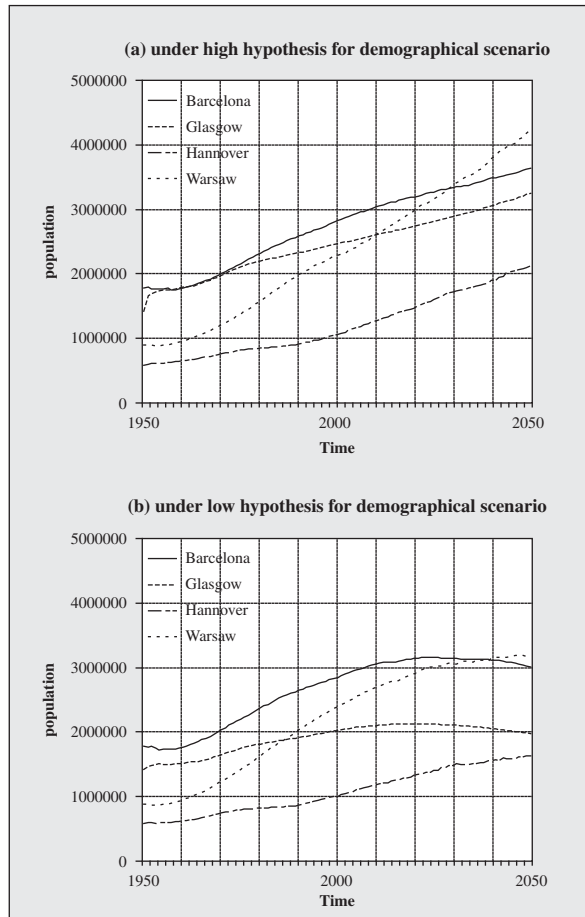
The model demonstrates that cities do not react the same way to such changes in demographic and political contexts. As an example, Fig. 12.11 represents the simulated evolutions of Barcelona's, Hamburg's, Warsaw's and Glasgow's populations according to the two extreme scenarios. These outputs also illustrate the ability of the model to produce qualitatively different city behaviors. Fig. 12.11 represents the evolutions of four cities with same specialization, "Technopole", for the period 1959–2000 according to these two extreme scenarios, no barrier and high demographic hypothesis on the one hand, barriers and low demographic hypothesis on the other. Quite naturally the trajectories corresponding to the first case predict higher growth (Fig. 12.11a) than for the second case (Fig. 12.11b). More interesting, qualitative differences appear between the two scenarios concerning the relative positions of the cities. The case of Glasgow for example is noticeable. This city suffers more than the others from the barrier effects. While the city almost rises to the level of Barcelona in the scenario without barriers, it remains far behind when barriers between blocks are introduced. In other respects, the growth of Warszawa seems to be more affected by the bad context of the second scenario than Barcelona, in the sense that the first catches up to the second more quickly in the first scenario.

There is no explicit rule in the model that would produce a more sustainable growth for economically diversified cities. The result expresses the combination of multiple interactions between couples of cities. As such, it is a consequence of self-organization processes. The model can then be used as an experimentation tool in order to explore the consequences of different constraints on interurban exchanges.

12.7 Conclusion: is the Future of Urban Systems Predictable?

Because of the many uncertainties about the future of cities, there is a need for exploratory models that could help determine the most plausible trends in their development. Of course, such models cannot be exactly predictive, since we know that predictions are often intractable for the underlying complex systems, especially in the long term. According to our method of data-driven simulation, the laws of urban dynamics presented in Chapters 6 and 8 are useful for validating an urban model, according to an acceptable representation of the past, but they have to be adapted and revised before using it as a predictive tool, as demonstrated in the Eurosim application. The SIMPOP2 model could also help in designing policies, by helping to estimate the relative cost of different choices. Is a polycentric development compatible with objectives of sustainable development? Can European cities keep their global competitiveness by sharing the investments dedicated to performing activities? Or should such investments remain concentrated in places offering the

Fig. 12.11 Future trajectories of a few technopoles according to Eurosim



highest returns? How strong are the links between the objectives of social cohesion and the spatial distribution of population and income?

Despite the accumulated knowledge from comparative studies on urban systems dynamics, many uncertainties remain about their possible evolution within a near future, i.e. the next hundred years. Two different kinds of events of the period may interfere with the existing dynamics that we have reported: those that come from inside the systems and those that arise from outside, from the societal environment of the system. In both cases, the conditions of interaction between cities are affected, and some effects of major shocks are already perceptible.

The first context that will introduce major changes in urban system dynamics is linked to the variations in the urban transition in different parts of the world. In developed countries, the question is how will the systems of cities evolve once they have “won” *all* the population in a given territory? Will they keep the same dynamical features as during their period of emergence and consolidation? What future can be expected when there is no longer migration from rural areas or local

demographic growth for sustaining the cities development? To what extent is a continuous population flow from outside (immigration from rural areas or foreign countries) or a minimum population growth necessary for maintaining their hierarchical organization? Some major turning points have already been observed in the evolution of urban systems. After a long period of spatial concentration, including an increase in urban population densities, the last four of five decades have been marked by local trends towards a de-concentration of resident population. Urbanized areas have expanded in surface much more rapidly than through demographic growth. This trend is sometimes interpreted as expressing the preference for rural places of residence; that would lead to a “counter-urbanization” (Berry, 1976), both at local and regional scales. On the other hand, trends toward concentration at a higher scale are observed. Will population and activities continue to concentrate in areas close to the largest metropolises? Are the small- and medium-sized isolated towns condemned to decline and disappear, as did so many villages in the past?

Both trends are suggested by our accumulated knowledge about past urban dynamics, but we should think of possible reversals that may happen because of completely different processes. Among the most frequently remarked potential changes are the demographic recession (population growth rates have been above 1% per year for two centuries, but they have recently become negative in some countries); the preoccupation for environmental quality and preservation of resources (which may hamper the further development of large cities); and new technologies for the circulation of information (which may change the relationships between the conception of cities, as places of work and residence). Thus, the “counter-urbanization” hypothesis could prevail and lead towards less inequalities in city sizes, a new population dispersal and a relative decay of large metropolises. The magnitude of the consequences of such processes can be implemented in the model by varying the values of some parameters. Meanwhile, measuring urban performance by population growth is a convenient, traditional way that facilitates comparisons in time, since there was a good parallelism between increase in the inhabitant number and the quantity of accumulated resources in systems that were not too heterogeneous. In the future and especially for comparisons at world scale, because the differences in standard of living are much higher, a more adapted measurement of the economic and social performance of cities like GDP or HDI would be needed (but is not yet provided by most statistical sources about cities all over the world, China excepted).

Reversals in dynamics also could come from “outside” the systems of cities, whatever their economic level: as the urban transition is continuing in developing countries, with unprecedented demographic growth rates, very large cities are becoming more and more the specificity of the urban systems in poor countries (Moriconi-Ebrard, 1993). In parallel, the globalization of economy and social information is developing new networks and increasing interdependencies between cities in the world (Taylor, Derruder, Saey, & Witlox, 2007). The disequilibrium between the hierarchy of city sizes according to their population and their gross product or income obviously is not sustainable over very long periods of time. Moreover, the evolution of national or continental urban systems according to their own evolutionary path cannot be prolonged independently of the overwhelming trend called

“globalization.” This trend may be seen as an “external shock” to many urban systems, because of its wide spatial extent and simultaneity in time, but, of course, it has been generated itself by the expansion of urban systems and the emulation of innovations that their dynamics is generating.

As it was the case for all previous innovations, the effects of that global integration on urban systems are predictable and measurable with the help of the SIMPOP2 model. At least in a first stage, the differences in accessibility to the newly created resources will be widening. As such, they may constrain the urban systems to keep a trend of concentration in the largest cities, because of the stronger competition between them. As demonstrated by Sassen (1991), very few centers in the world are concentrating the major parts of global finance, and it is not yet sure if the further developments of these activities can percolate in a larger number of “global cities” around the world. It is also uncertain if the traditional powerful urban centers have the capacity for maintaining their position in the emerging global city networks, or if they will be successfully challenged by new places of interest for investors (Hall, 1999). Meanwhile, new urban specialization may emerge in connection with the innovative economic sectors (Gaspar & Glaeser, 1996). The model can predict the general evolution of urban systems but, of course, not the exact location of these emerging new urban functions.

However, what we retain from our observations in Chapters 6 and 8 as well as from our experimentation with the model is that much regularity and universal rules can be expected in the evolution of any urban system. The major trend of historical path dependence illustrates best the capacity of resilience of such systems, their ability to absorb so many quantitative and qualitative changes in social organization without modifying their basic organization. Moreover, from that specific ability, *urban systems can be considered as “adaptors” of the spatial organization of societies subject to cultural, economic and technological changes. The SIMPOP2 model is a relevant, efficient and flexible tool for exploring their future evolution.*

References

- Allen, P. (1997). *Cities and regions as self-organizing systems: Models of complexity*. Amsterdam, The Netherlands: Gordon and Breach
- Allen, P., & Sanglier, M. (1979). Dynamic models of urban growth. *Journal of Social and Biological Structures*, 2, 269–298.
- Arthur, W. B. (1994). *Increasing returns and path dependence in the economy*. Ann Arbor, MI: University of Michigan Press.
- Batty M., & Torrens, P. (2002). Modeling complexity: the limits to prediction. *Cybergeo*, 201.
- Berry, B. J. L. (1976). *Urbanization and counter urbanization*. London, UK: Sage.
- Bura, S., Guérin-Pace, F., Mathian, H., Pumain, D., & Sanders, L. (1996). Multi-agent systems and the dynamics of a settlement system. *Geographical Analysis*, 2, 161–178.
- Casti, J., & Swain, H. (1975). Catastrophe theory and urban process. In *Lecture notes in computer science* (Vol. 40, pp. 388–406). London, UK: Springer-Verlag.
- Cattan, N., Pumain, D., Rozenblat, C., & Saint-Julien, T. (1994). *Le système des villes européennes*. Paris, France: Anthropos.
- Cicille P., & Rozenblat, C. (2004). *Les Villes Européennes*. Paris, France: DATAR, Documentation française.

- Clarke, G.P. (Ed.) (1996). *Microsimulation for urban and regional policy analysis*. London, UK: Pion.
- Eubank, S., Guclu, H., Kumar, V. S. A., Marathe, M. V., Srinivasan, A., et al. (2004). Modeling disease outbreaks in realistic urban social networks. *Nature*, 429, 180–184.
- Ferber, J. (1995). *Les systèmes multi-agents. Vers une intelligence collective*. Paris, France: Inter Editions.
- Gaspar, J., & Glaeser, E. L. (1996). *Information technology and the future of cities*. Harvard Institute of Economic Research Working Papers 1756, Cambridge, MA: Harvard University.
- Hall, P. (1999). The future of cities. *Computers, Environment and Urban Systems*, 23(3), 173–185.
- Holm, E., & Sanders, L. (2007). Spatial microsimulation models. In L. Sanders (Ed.), *Models in Spatial Analysis* (pp. 159–195). London, UK: ISTE.
- Michel, F., Ferber, J., & Gutknecht, O. (2001). Generic simulation tools based on MAS organization. In *Proceedings of Modeling Autonomous Agents in a Multi-Agent World*, 10th European Workshop on Multi-Agent Systems.
- Moriconi-Ebrard, F. (1993). *L'urbanisation du Monde*. Paris, France: Anthropos, Economica, Collection Villes.
- Openshaw, S. (1997). Developing GIS-relevant zone-based spatial analysis methods. In P. Longley, & M. Batty (Eds.), *Spatial analysis: Modeling in a GIS environment* (pp. 55–73). New-York: John Wiley.
- Portugali, J. (2006). *Complex artificial environments: Simulation, cognition and VR in the study and planning of cities*. New York: Springer.
- Pumain, D. (1998). Urban research and complexity. In C. S. Bertuglia, G. Bianchi, & A. Mela (Eds.), *The city and its sciences* (pp. 323–361). Heidelberg, Germany: Physica Verlag.
- Pumain, D. (2000). Settlement systems in the evolution. *Geografiska Annaler*, 82B(2), 73–87.
- Pumain, D., Bretagnolle, A., & Glisse, B. (2006). Modeling the future of cities. In *ECCS'06, Proceedings of the European Conference of Complex systems* (pp. 25–29). Oxford, UK: University of Oxford.
- Pumain, D., Bretagnolle, A., Glisse, B., Mathian, H., Vacchiani-Marcuzzo, C., et al. (2006). Final report of the European program “Innovation Society as a Complex System” (ISCOM, directed by D. Lane, www.iscom.unimo.it).
- Ramat, E. (2007). Introduction to discrete event modeling and simulation. In D. Phan, & F. Amblard (Eds.), *Agent-based modelling and simulation in the social and human sciences* (pp. 35–62). Oxford, UK: The Bardwell Press.
- Sanders, L. (1992). *Systèmes de villes et synergie*. Paris, France: Anthropos.
- Sanders, L. (2007). Agent models in urban geography. In D. Phan, & F. Amblard (Eds.), *Agent-based modelling and simulation in the social and human sciences* (pp. 147–168). Oxford, UK: The Bardwell Press.
- Sanders, L., Pumain, D., Mathian, H., Pace-Guérin, F., & Bura, S. (1997). SIMPOP: a multi-agent system for the study of urbanism. *Environment and Planning B*, 24, 287–305.
- Sanders, L., Favaro, J.-M., Glisse, B., Mathian, H., & Pumain, D. (2006). Dynamics of the European urban network. Final report of the European program *Time-Geographical approaches to Emergence and Sustainable Societies*, (<http://www.tigress.ac/reports/final.htm>).
- Sanders, L., Favaro, J. M., Glisse, B., Mathian, H., & Pumain, D. (2007). Intelligence artificielle et agents collectifs: le modèle EUROSIM. *Cybergeo*, 392.
- Sassen, S. (1991). *The global city: New York, London, Tokyo*. Princeton, NJ: Princeton University Press.
- Taylor, P. J., Derruder, B., Saey, P., & Witlox, F. (2007). *Cities in globalization*. London, UK: Routledge.
- Weidlich, W., & Haag, G. (1988). *Interregional migration – Dynamic theory and comparative analysis*. London, UK: Springer-Verlag.
- White, R. W. (1978). The simulation of central place dynamics: two sector systems and the rank size rule. *Geographical Analysis*, 10, 201–208.
- Wilson, A. (1981). *Catastrophe theory and bifurcation*. London, UK: Croom Helm.

Chapter 13

Modeling Innovation

Roberto Serra, Marco Villani and David Lane

13.1 Why Model Innovation, Which Models of Innovation?

The innovation theory (briefly, I_2T), which has been developed in the ISCOM project and which is presented in this book (Chapters 9 and 10), is based on the analysis of different case studies, spanning different time periods and different kinds of products, from the introduction of printing in the Renaissance, to key new technologies introduced in the 19th and 20th centuries, up to present-day ongoing innovation efforts.

This theory is qualitative (although fairly rigorous), and it does not claim to be able to provide either predictions or quantitative descriptions of the phenomena that it addresses. The theory makes statements concerning the core entities, relationships and processes that are necessary to understand a wide range of different phenomena and aims at uncovering some of their non-obvious features.

A natural question is whether modeling can be of any help in the development and refinement of such a theory. Indeed, quantitative models are often used to compare the theoretically expected behavior with the one observed, and, wherever quantitative predictions are possible, to forecast the behavior of the interesting variables.

But if the theory is inherently qualitative, what contribution can models provide? Of course, this question has been asked several times in the development of the social sciences, and it has received interesting, albeit partial, answers.

In the present context, models of particular interest are those that allow one to describe the behavior of a large collection of agents, thus bridging the gap between the level of a few human agents or organizations, which is addressed by I_2T , and the level of the behavior of a system comprising many more such agents. In order to be precise, let us stress that we concentrate here on dynamical models, which allow one to observe how the system changes in time.

This is by no means the only kind of model that might initiate a dialogue with the present theory. For example, a different kind of model could be also considered,

R. Serra (✉)

University of Modena and Reggio Emilia, Modena and Reggio Emilia, Italy

aimed at a sophisticated description of the interaction among two or few agents. However, in the development of our research, we attributed higher priority to exploring the behaviors of a collection of interacting agents, as suggested in, for example, Epstein and Axtell (1996), Gilbert and Terna (2000), and Axelrod and Tesfatsion (2006). The reason is that I_2T is endowed with several positive feedback mechanisms, and it is well known that in such cases, counterintuitive behaviors can be observed at a system level. Therefore, the behavior of a large group of agents cannot be easily inferred from that of a pair, or a few of them. Modeling can thus provide a major contribution to unfold the large-scale consequences of the theory.

The behavior of a collection of agents observed in the model can then provide information to the theorists about the consequences of their assumptions. Confirmation of expected behaviors is an interesting result. Even more interesting, there may be unanticipated or surprising results, which can lead to a re-thinking of the theory itself.

However, the interpretation of the results of the simulation is complicated by the fact that the model cannot simply be considered “the theory in action:” in order to make it manageable, it is indeed necessary to introduce simplifications and make choices that are not dictated by the theory.

Therefore, when there are surprising results at the aggregate level, it may be unclear whether these are the consequences of the simplifications introduced, or whether they represent a real outcome of the theory itself. It is therefore necessary to perform an in-depth analysis to disentangle the contribution of the (largely arbitrary) choices that are model-related from genuine unanticipated consequences of the theory.

While bridging the gap between the levels of a few agents to that of a collection of many of them appears to be the most important motive to resort to modeling, there are other important reasons that make modeling useful to develop a broad and rigorous perspective on innovation. They can be briefly summarized as follows:

- The development of a formal model requires that, in the restricted universe defined by the model itself, the claims of the theory be precisely stated. Therefore ambiguities and imprecisions can be more easily spotted and amended, thus sharpening the rigor of the theory
- The models developed provide a toy universe where one can easily experiment with the consequences of new suggestions, thereby leading to the possibility of exploring a wide set of alternatives
- The observations of the way in which the models evolve and of the different behaviors of the agents can provide a useful way to describe and communicate the theory itself and its major characteristics.

Let us now consider the type of models that are better suited for the purposes stated above. “Dynamical Model” is a general term that encompasses several alternative approaches, including differential equations and their time-discrete analogs, stochastic differential equations, cellular automata and many others. We have chosen to focus our work on so-called agent-based models (ABMs), which are

sometimes referred to as multi-agent systems, in particular in the Artificial Intelligence literature.

The main reasons for this focus are that in these models, it is possible:

- to deal with agents capable of sophisticated information processing,
- to endow agents with adaptive and learning capabilities that not only modify their behavior in time, but also change the rules which underlie this behavior,
- to introduce agents' goals and beliefs in a straightforward way, and
- to take fully into account the heterogeneity of agents, therefore escaping the constraints of the "representative agent" of traditional economic modeling.

ABMs are well suited, therefore, to relate hypotheses concerning the microscopic behavior of individual agents to emergent collective phenomena in a population composed of many interacting agents.

While models based on differential equations are too often quickly dismissed by devotees of ABMs as inadequate, we have recently argued that their expressive power is greater than what is usually assumed (Serra & Villani, 2006). However, in the present case we prefer to resort to ABMs because the reference theory requires that agents possess sophisticated ways to handle information about the world they inhabit, and because both intentionality and heterogeneity are thought to be important.

Attractive as agent-based models may be, they allow the modeler too much freedom with respect to the design of agents and their interactions. ABMs often have too many variables and parameters, but their worst shortcoming is that there is no theory of their behavior comparable to, say, that of differential equations. Therefore, the effects of the choice of the parameters cannot be foreseen or evaluated based on such a theory. Moreover, as observed above, these models usually include strong nonlinearities that are known to be able to give rise to unexpected outcomes.

Indeed, algorithmic models like ABMs may be highly arbitrary, so there is an *embarras de richesses*, and conceptual guidelines are needed to constrain the set of allowable models. Two major considerations may help to provide such constraints.

The first is that it is necessary to look for model behaviors that are robust with respect to different perturbations, which may involve changes in parameter values and, to some extent, in the choice of the functions that describe the agents and their behaviors. If simulation results match real world observations (which may be limited and imprecise, as so-called "stylized facts") only for a limited set of parameter values, e.g. for a very particular kind of decision rule, then suitable reasons should be found to justify these choices, before claiming that the model has anything to do with the real system it is supposed to simulate.

The second methodological consideration is that a model may be constrained strongly by its relationship with a theory of the social process that it is supposed to describe. In this way, the model should concentrate on those aspects that the theory identifies as the most relevant, dropping out or drastically simplifying a set of related, but (according to the theory) less important aspects. This approach is similar, in some sense, to that of classical hard science, and it is the one we take in this paper.

The theory provides the basic entities of the model and describes the kinds of relationships among them. The model represents a simplified universe inhabited by (some of) the theoretical entities, engaging in some of the kinds of relationships predicated by the theory. Simplification is necessary to deal with manageable systems: we do not look for an all-encompassing model, but, rather, we foresee a family of models that capture different features of the theory. The model described here is meant to represent the first member of this family.

The plan of the chapter is the following. In Section 13.2, we will briefly summarize those aspects of I_2T that are the more relevant to our purposes, and, in particular, we will focus on the indications and constraints that they provide for models consistent with the theory. In the following Section 13.3, we will outline a description of one such model, which we call I_2M (ISCOM Innovation Model). Although we made some efforts to keep the model as simple as possible, the need to include a strong relationship with the theory led to a model with a variety of features, and explaining them in detail would take up the whole chapter. (Readers interested in more details about the model may consult Lane, Serra, Villani, & Ansaloni, 2005).

A particular problem is the setting of the initial state of the model. Our approach, inspired by the notion of “exaptive bootstrapping” (see Chapters 1 and 14, this volume), is described in Section 13.4 below. In the same section, we also discuss the reference parameter values, chosen based on extensive simulations (which however cannot even approach a complete sampling of the large parameter space) and provide some details about the interactions between recipes and artifact space. Section 13.5 describes how the model responds to various kinds of perturbations, linking up with general interest in “avalanche size distributions” for complex systems (see for example Kauffman, 1993; Bak, 1996).

The penultimate two sections are dedicated to an exploration of some interesting behaviors of the model. In Section 13.6, we discuss the influence of the agents’ propensity to build long-lasting relationships on the probability of successfully realizing innovation projects. Section 13.7 investigates the influence of the initial condition (i.e. the state of the system at the beginning of the simulation) on formation of topological structures in artifact space. The final section sums up the results obtained so far with the model and provides some directions for further research.

13.2 Requirements from the Theory

We now examine the major constraints that the theory imposes on the model. We will not attempt here to summarize the I_2T , which is extensively discussed in Chapters 1, 9 and 10 of this volume, but we will limit ourselves to emphasizing the implications of the theory’s main assumptions on our modeling efforts.

The first claim of the theory is that innovation processes involve transformations of relationships among agents, among artifacts, and between agents and artifacts. This notion underlies the concepts of an “agent-artifact space” and, in particular, that of “reciprocity,” which essentially claims that artifacts mediate interactions

between agents, and vice versa, and, therefore, both agents and artifacts are necessary for a proper understanding of the behavior of market systems and of innovation. Despite first appearances, this is neither a minor nor an obvious assumption, since it is entirely conceivable to focus just on artifacts (“technological trajectories” have been given much credit in the past) or on agents alone, as is generally the case for neoclassical theories of innovation or epidemiology-based theories of innovation diffusion (see Chapter 11 for further discussion on this point).

One straightforward consequence of this claim is that it excludes the possibility simply to project agent-artifact space onto one of its two constitutive subspaces, ignoring the dynamics in the other subspace.

In I_2T , artifacts are given meanings by the agents that interact with them, and different agents take on different roles. The meaning of artifacts cannot be understood without taking into account the roles that different agents can play. Thus, artifacts may be given different meanings by different agents or by the same agent at different times.

Innovation, in the sense of I_2T , is not just novelty, but also a transformation in the structure of agent-artifact space, which unleashes a cascade of further changes. The innovation may involve the introduction of a new artifact, but also a change in relationships with other agents, or even a new interpretation of an existing artifact. In a sense, this theory can be seen as a theory of the interpretation of innovation, where “interpretation” actually means attribution of functionalities to artifacts and of roles to agents.

According to the theory, a new interpretation of artifact functionality typically arises in the context of so-called generative relationships. By interacting, a few agents come to invent and share a new interpretation, based on the discovery of different perspectives and uses of existing or expected artifacts. The generative potential of a relationship may be assessed in terms of the following criteria:

- heterogeneity: the agents are different from each other, they have different features and different goals, but the heterogeneity is not so intense as to prevent communication and interaction;
- aligned directedness: the agents are all interested in operating in the same region (or in neighboring regions) of agent-artifact space; and
- mutual directedness: the agents are interested in interacting with each other.

Moreover, agents must be allowed to interact and to take joint actions.

These are the aspects of the theory that appear to be most relevant for the modeling effort. This theory of innovation is highly sophisticated in describing the interactions between different players in innovation processes, and it cannot be entirely mapped onto a specific computer-based model. Therefore, the modeling activity aims at developing models that are based on abstractions of some key aspects of this qualitative theory.

The basic requirements for a model aiming at a dialogue with the theory, according to the principles expressed above, are then the following:

1. both agents and artifacts must be present,
2. the artifacts' meanings must be generated within the model itself: since I_2T claims that new meanings are generated through interactions among agents and artifacts, it would be inappropriate here to resort to an external oracle to decide *a priori* which meanings are better than others,
3. the roles of agents must also be generated within the model: indeed the theory claims that new roles are also generated through interactions among agents and artifacts,
4. agents must interact with artifacts and with other agents: interacting only with artifacts would prevent the possibility of describing agent-agent relationships,
5. an agent should be able to choose the other agents with whom to start a relationship; in general, an agent will be able to handle a finite number of relationships at a time, and it will choose a subset of the other agents as its partners. Agents must be allowed to dissolve a disappointing or unsatisfactory relationship to look for a better one, and
6. agents must have intentionality: they may be interested in certain types of artifacts or in entering a relationship with other agents.

13.3 An Outline of the I_2M Model

Let us now briefly introduce the main features of the model I_2M , which has been developed on the basis of the above principles. We will limit ourselves here to a very concise overview of the model, and we refer the interested reader to (Lane et al., 2005) for a more complete and detailed account.

In I_2M , agents can “produce” artifacts, which, in turn, can be used by other agents to build their own artifacts. An agent can produce several artifacts for the same agent, and it can produce one type of artifact for different customer agents.

Each agent has a set of recipes that allow it to build new artifacts from existing ones. Agents can try to expand the set of their recipes by applying genetic operators either to their own recipes or, by cooperating with another agent, to the joint set of the recipes of both. Moreover, each agent has a warehouse or store, where its products are put, and, from which, its customers can take these products.

The meaning of artifacts in this model is just what agents do with them, while the role of agents is defined by which artifacts they produce, with whom, and for whom. The role of agents is also partly defined by the social network they are embedded into, as explained in the following paragraph.

If we look at the network of agents, there is a directed link of the “strong tie” type between A and B if there is a customer-supplier relationship between the two. There are also weak ties between two agents (“acquaintances”), which refer to the fact that agent A knows something about agent B (e.g., its products). Note that other types of networks can also be considered; for example, the artifact network (where two artifacts are linked if there is an agent which uses one of them as an input to produce the other) and also a heterogeneous network, where both agents and artifacts are

explicitly represented. In this latter case, most relationships are between an artifact and an agent: there is no direct link between two artifacts and no strong tie between two agents, since each of these ties is mediated by a direct tie with an artifact, but there may be links between two agents of the weak-tie kind.

The value that an agent (say, A) gives to its relationship with another agent, B, is summarized in a single numerical variable (the “vote”), which is composed of the sum of two terms. The first term is increased or decreased based on the history of supplier/customer interactions between A and B, while the second term takes into account the results of previous co-operations in developing new projects, if any. A parameter, which can be changed in different simulations, determines the relative weight of these two terms.

A key point is the structure of artifact space. It is required that the space has an algebraic structure, and that suitable constructors can be defined to build new artifacts by combining existing ones. For reasons discussed elsewhere, we have adopted a numerical representation for artifacts and the use of mathematical operators instead of e.g., binary strings, λ -calculus or predicate calculus. Therefore, the agents are “producers” of numbers by applying mathematical operators to other numbers, and recipes are defined by a sequence of operators and the numerical arguments on which these operators act.

Each agent is also endowed with a numerical variable (called its “strength”) that measures how successful it has been so far. Strength increases proportionally to the number of artifacts sold and decreases proportionally to the number of active recipes. Note, however, that strength, in the present version of the model, cannot be interpreted as “money” since it is not conserved in the interactions between two agents (if A and B interact, and ΔS_A and ΔS_B represent the change in strength of the two agents due to this interaction, it may well happen that $\Delta S_A \neq -\Delta S_B$).

As far as innovation is concerned, an agent can invent new recipes or eliminate old ones. In the present version of the model, no new agents are generated, while agents can die because of lack of inputs or of customers.

The model is asynchronous: at each time-step, an agent is selected for update, and it tries to produce what its recipes allow it to do. Therefore, for each recipe, it looks for the input artifacts, and, if they are present in the stocks of its suppliers, it produces the output artifact and puts it into its stock (the supplier stocks are, of course, reduced). Production is assumed to be fast, i.e. it does not take multiple steps.

Besides performing the usual “buy-and-sell” dynamics, when its turn comes, an agent can also decide to innovate. Innovation is a two step-process. In the first step, the agent defines a goal, i.e. an artifact that it may add to the list of its products, while the second step concerns the attempt to generate a recipe that yields the goal as output.

In the goal-setting phase, an agent chooses one of the existing types of artifacts (which, recall, is a number M) and then either tries to imitate it (i.e. its goal is equal to M) or it modifies it by a jump (i.e. by multiplying M by a random number in a given range). It has been verified that imitation alone can lead to a sustainable structure for agent-artifact space, in which, however, innovation eventually halts, which need not be the case when (some) agents innovate via jumps (Lane et al., 2005).

After setting its goal, an agent tries to reach it either by using the operators of one of its recipes with different inputs or by combining two recipes to generate a new one with genetic algorithms. In this phase, the agent can decide to cooperate with another agent, sharing the recipes of both, to reach a common goal. The propensity of an agent to cooperate is ruled by a parameter, while for most simulations we discuss, choice of the partner is made with a probability distribution biased in favor of those agents to whom the choosing agent has assigned a high vote.

13.4 Typical Model Behaviors

13.4.1 Initialization

Setting up initial conditions “by hand” is time demanding, so we developed an automatic “initial condition builder.” First, raw materials are introduced; afterwards, agents are added one at a time. The randomly generated recipes of the newly added agent can use only already existing artifacts as inputs, producing in this way (possibly) new artifacts. The networks that result from such a process have the property that older artifacts are likely to be connected to more agents than new ones.

Networks generated according to the above procedure have the common feature that the unused artifacts (in particular, those produced by the most recently added agent, many of which are new to the system) cannot provide any source of activation to their owners, and, therefore, the agent that produces them cannot survive, causing a “domino effect” that generally destroys the entire network. To avoid this collapse, we can either provide an external reward region or introduce a modification of the dynamics (e.g. the possibility of changing a provider) that leads to self-sustaining rings, as shown in Fig. 13.1.

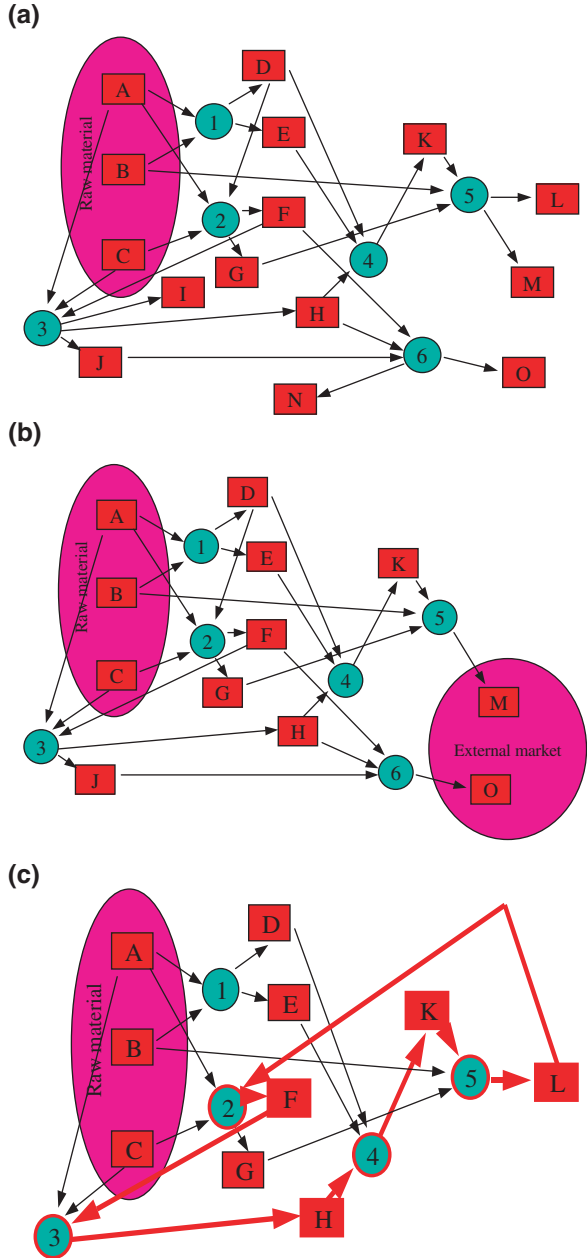
13.4.2 Time Behavior of Some Key Variables

We have performed many experiments with this model, in order to test the model and to understand some of its behaviors. In particular, we find that

1. Systems without innovation collapse (Fig. 13.1a) unless: there is an external market (Fig. 13.1b) or a self-sustaining loop is already present within the system (Fig. 13.1c);
2. Imitation alone is unable to introduce a significant number of novelties; and
3. The simultaneous presence of imitation and jump actions allows a strong increase of diversity in the resulting artifact space.

To provide some hints about the temporal behavior of the model, we do not present statistical results, but rather show typical model behaviors. During the simulation shown below, agents perform innovations that are either “social” (by linking up with new suppliers) or “technological” (by generating new goals and recipes).

Fig. 13.1 (a) The initial unstable situation, and (b), (c) two possible stable configurations. In (b), the agents' survival is guaranteed by the presence of an external market (artifacts N and L disappear), whereas in (c), it is guaranteed by the presence of a self-sustaining cycle (highlighted in red – artifacts M, O and J disappear)



Each agent knows only the artifacts produced by those agents with whom it has a tie, weak or strong.

We have examined different innovation procedures. In particular, when innovation in recipes is performed by applying genetic operators to the set of existing

recipes, we have considered both the case of “lonely” agents, who modify only their own recipes, and “collaborative” agents, who share their recipes in genetic search of new ones.

A typical example is shown in Fig. 13.2, where we can observe the temporal behavior of some variables. The parameters describing the characteristics of the run are collected in Table 13.1.

Figure 13.2a shows the number of artifacts that are present in the system, in particular, the total number of artifacts, the number of artifacts produced, and the number of artifacts used. The difference between the total number of artifacts and that of artifacts used is the number of still unused innovations.

For these counts, the same artifact type (i.e., a number), if produced by different agents, counts as different artifacts. In contrast, Fig. 13.2b shows the number of different artifact *types* present in the system. Clearly, many artifacts types are produced by more than one agent.

Figure 13.2c, d show the temporal behavior of some aspects of the artifact space, that is, the diameter (the difference between the highest and the lowest number

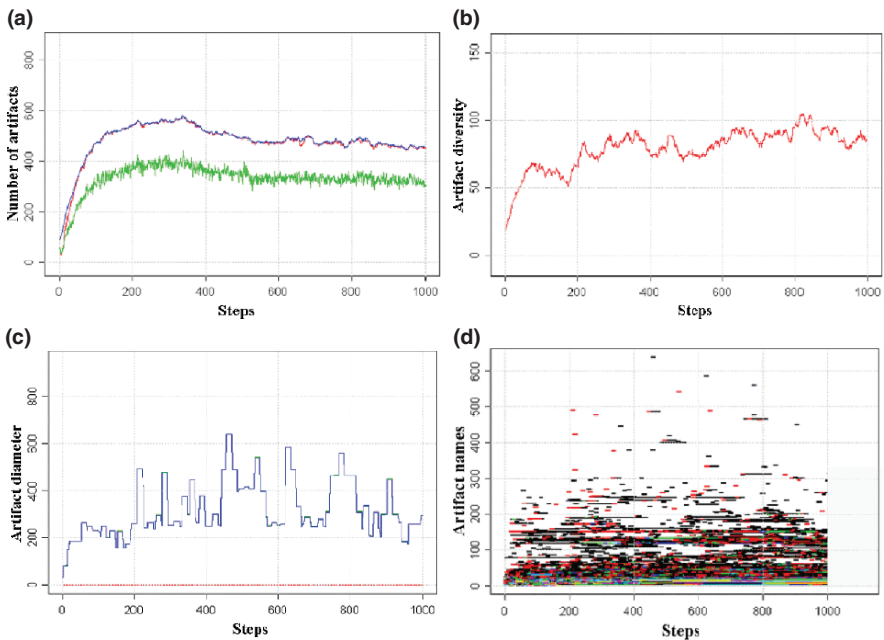


Fig. 13.2 The temporal behavior of some aspects of the artifact space. **(a)** The total number of artifacts present inside the system (*blue line*: total number of artifacts; *red line*: total number of produced artifacts; *green line*: total number of used artifacts). **(b)** The number of different “names” present inside the system (an index of the diversity present inside the system). **(c)** The diameter of the artifact space (the difference between the highest and the lowest “names” presents on the system). **(d)** a representation of the evolution of the artifact space (for each step there is a dot if the corresponding value is realized by at least an artifact, the different colors representing the different numbers name’s realizations)

Table 13.1 The parameters describing the characteristics of the run described in Section 13.4.2

Parameter	Value and/or description
Initial number of agents	40
Initial number of recipes for each agent	2
Number of allowed inputs for each recipe	2,3
Number of raw materials	2
Innovation probability during each step	0.2
Jump probability (once an agent initiate the innovation phase)	0.7
Jump range	[0.7, 6]
Innovation strategy	Each agent tries by itself; in case of failure, it tries again collaborating with another agent
Agents' knowledge	Each agent "knows" the artifacts: – of its providers – of its acquaintances
Number of random acquaintances for each agent	5
Vote strategy	Based upon the past successful joint projects

present on the system) and the "history" of the artifacts (where for each step there is a dot if the corresponding value is realized by at least one artifact, the different colors representing the different numbers of artifacts that realize the "name").

By observing Fig. 13.2, it is possible to monitor the expansion of artifact space, which is followed by a more stable phase. The peak in the total number of artifacts does not correspond to periods of maximum diversity, and the space occupation is not uniform (there are deserted areas close to zones with high density of artifacts).

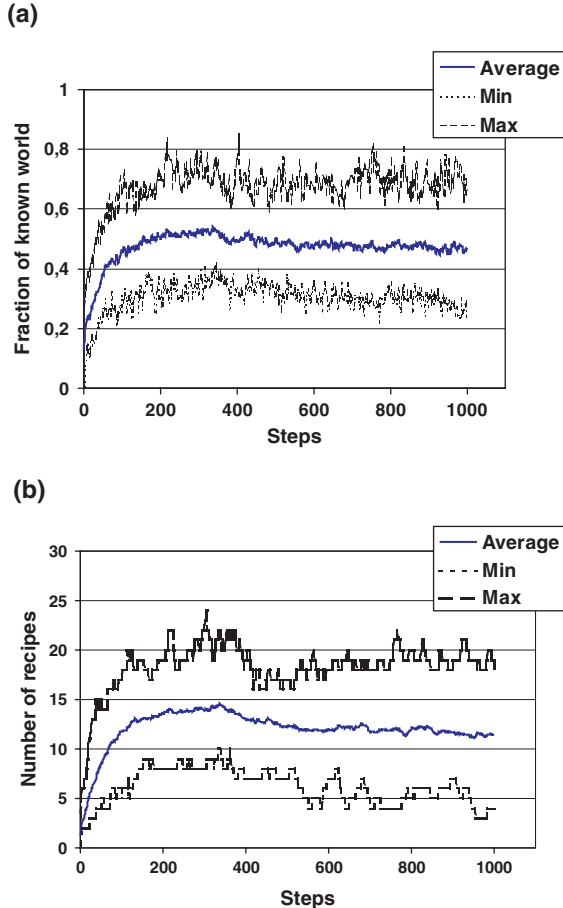
Other variables monitor the agents' behavior: for example, Fig. 13.3 shows the fraction of the system known and the number of recipes owned by each agent. It is possible to observe the emergence of different paths: despite the fact that, in this experiment, the agents share the same parameter values, in the end there are agents that know a very small fraction of the system and agents that are aware of more than 60% of the artifacts, and there are agents that own less than five recipes while some have more than fifteen. This high heterogeneity implies that not all the agents are allowed to develop to the same level and suggests the presence of particular structures inside the system.

13.4.3 Goals

In I_2M , agents have a "goal" in artifact space, i.e. a new artifact they aim at producing, either by themselves or in cooperation with another agent. The goal (G) itself is determined by the following method:

- First, one of the existing artifacts is chosen at random (heuristically, this is a way to sample the set of existing artifacts), let us call it the intermediate goal (IG).

Fig. 13.3 (a) The fraction of the system known by each agent; (b) the number of recipes owned by each agent. The plots report for each run step the average, the minimum and the maximum of the interested variable

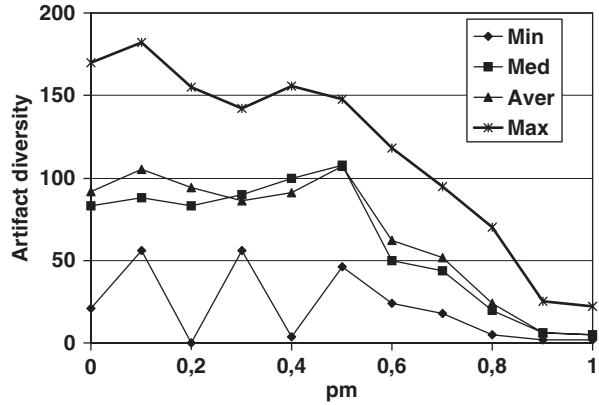


- Second, the IG is modified by a “jump.” Recall that artifacts are numbers: the number corresponding to the IG is multiplied by a random number belonging to a given range R (briefly: $G = J(IG)$).

Resorting to randomness simulates the definition of a goal in a system where it is hard to predict whether a particular artifact will be successful or not, and where all the properties of an artifact “in use” cannot be established *a priori*, solely on the basis of engineering design considerations.

An important parameter is the goal persistence, which measures how long an agent will stick to its goal if it has been so far impossible to reach it. We can assume that each agent has the probability p_m of maintaining its own goal, ranging from 0 (changing the goal each time the agent has to innovate) to 1 (always keeping the same goal until it is successfully produced). Qualitatively, we can regard p_m as a measure of agent flexibility, i.e. its propensity to change its objectives (of course, it is actually a decreasing function of “flexibility”).

Fig. 13.4 Diversity of artifacts present in the system (average, median, minimum, and maximum out of 10 runs) as a function of the agent probability p_m of maintaining the goal. If p_m is equal to 1, the agents try to realize the same goal until they succeed in reaching it



The effects of the agents' different innovation propensities are evident by observing Fig. 13.4, which shows the diversity of artifacts in the system after 2000 steps of simulation.

If the probability of changing objective is sufficiently high, the system is able to maintain a sustained growth of diversity (Fig. 13.5); otherwise, diversity levels off and remains always lower than in the previous case. A threshold seems to appear around $p_m \approx 0.5$. This effect is likely due to the fact that the perseverance in attempting to create an artifact that is hard or impossible to achieve with the available operators constrains some agents to continue in unsuccessful attempts, and, more importantly, prevents them from introducing other innovations.

13.4.4 Recipe Structure

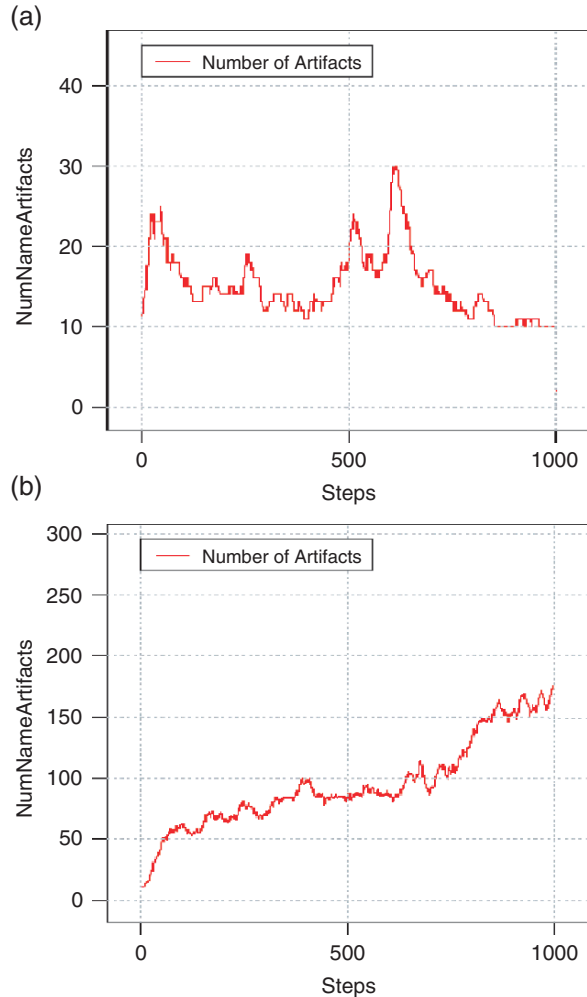
In the experiments described so far, the recipes can use only two kinds of operator, that is, addition and subtraction, and have two or three input artifacts. The recipe structure is very simple, but is it enough to produce the increase of diversity of Fig. 13.5b, where the environments at step 50 and at step 1000 are very different.

To better appreciate this phenomenon, we made some experiments where the recipes have initially the same proportion of 2, 3, 4, 5, 6, or 7 inputs – in order to simplify the exposition; in the following, we indicate the number of inputs of each recipe as the recipe “size.”

Not all kinds of recipes have the same survival probability: in fact, recipes unable to find all the desired inputs at the same time will, after a while, disappear, and obtaining seven different inputs at the same time is more difficult than finding only two or three.

Without any particular constraint, after a while we should observe an inhomogeneous distribution where short recipes overwhelm longer recipes: in effect, in absence of goals, this distribution is the final one (Fig. 13.6b). However, as just reported, we have a particular constraint: the agents do not try to generate new

Fig. 13.5 The artifact diversity for a system with high p_m (a) and low p_m (b) versus time. Note the change of scale on the graphs



recipes at random but rather they aim at specific points in artifact space. It is not necessarily the case that simple recipes can reach the agents' desired goals!

In fact, the runs where the agents have goals reveal a more complex history. At the starting point, more or less all the kinds of recipes have the same number of realizations (Fig. 13.6a), whereas at step 2000, a distribution where the short recipes overwhelm the longer ones begins to emerge. But as the artifact space becomes more fragmented by increasing the deserted areas interposed with zones having high density of artifacts, simple recipes are not able to fulfill all the requirements – for example, to be able to build a new artifact belonging to an area separated from all the available inputs areas by one or more deserted zones. In this case, recipes with more than two inputs can be useful: at step 5000, the agents more frequently use recipes with four or five inputs (Fig. 13.6c).

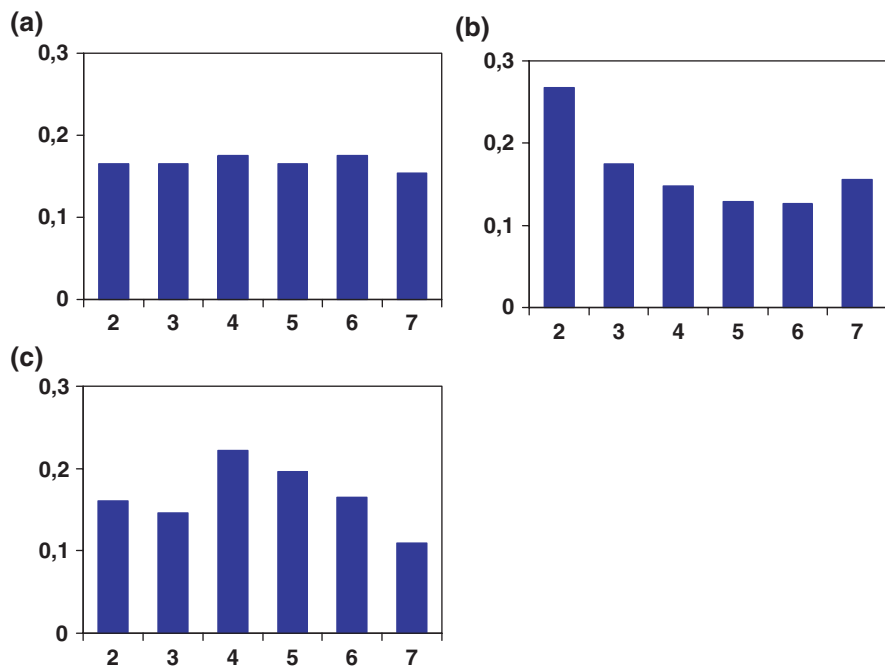


Fig. 13.6 The distribution of recipe size during a typical run where the allowed number of inputs spans from two to seven. (a) Initial distribution; (b) the final the distribution of recipe size in systems without goals; (c) the final distribution of recipe size in systems with goals

13.5 Distribution of Avalanches

13.5.1 Introduction

The model simulations show systems that change rapidly in time: artifacts continuously come to life and die, new functionalities are added to the already existing ones, agents can augment their own recipes or rapidly decline and eventually expire. Some questions naturally arise: how robust is the system? Are there bursts? What is the distribution of the size of avalanches of change?

Two different kinds of perturbations can be considered:

- internal perturbations, which can be observed by looking at the time plot of some relevant variable, e.g., the total number of artifacts. These are due to events that occur in the time evolution of the system. For example, it may happen that a recipe is eliminated that was used by other agents, and these latter are no longer able to produce some of their products. If the situation does not recover rapidly, this may lead to the loss of other recipes and perhaps to the death of one or more agents; and

- external perturbations, which can be generated from outside, for example by removing an agent or a recipe at a given time step. In this case we compare the evolution in the perturbed vs. unperturbed case.

To observe the system behavior, we can consider several variables. In the following discussion, we consider mainly the total number of artifacts as indicator of the “health state” of the system.

13.5.2 Internal Perturbations

The internal perturbations are the “natural fluctuations” of our system. By observing several runs we can often see that the system approaches a quasi steady state level of the number of artifacts, but we can observe also that there are continuous oscillations around this level, which have very different dimensions. Very often, the deviations are small, but sometimes they are quite large. The suspicion arises that we are observing a sort of “punctuated equilibrium,” where perturbations of all scales coexist. Big avalanches can have their cause inside the system, without the necessity of external interference. The succession of big avalanches, interposed by more numerous avalanches of minor magnitude, could give birth to a power law distribution for the fluctuations’ magnitude, as it seems to have been observed in several natural phenomena ranging from earthquakes to species extinctions (Bak, 1996).

To observe the magnitude of internal perturbations, we consider 20,000 steps of the evolution of systems with 40 agents. During the simulations, we observe the temporal distance among the peaks (see Fig. 13.7a) and the crises’ magnitude (that is, the distance between adjacent maxima and minima – see Fig. 13.7b). The series do not span many orders of magnitude, and, therefore, one cannot claim that the tails surely follow power law behaviors, but the shapes do not contradict such a hypothesis. For the variables shown in the figures the exponents of the best approximating power law are -1.5 and -2.8 , respectively.

13.5.3 External Perturbations

Another interesting feature of these systems is their robustness when they are subjected to external perturbations. Again, we consider mainly the total number of artifacts as indicator of the “health state” of the system.

Each experiment is composed of a pair of model runs. Each pair starts with the same initial condition and the same seed of the random number generator: therefore, the two runs are identical as long as they proceed without perturbations. At a predetermined step, a perturbation occurs: starting from that step, the two trajectories diverge. We tested several time points to perturb the system; here we present experiments where the perturbation is introduced at step 800, after the peak that characterizes the greater part of our runs.

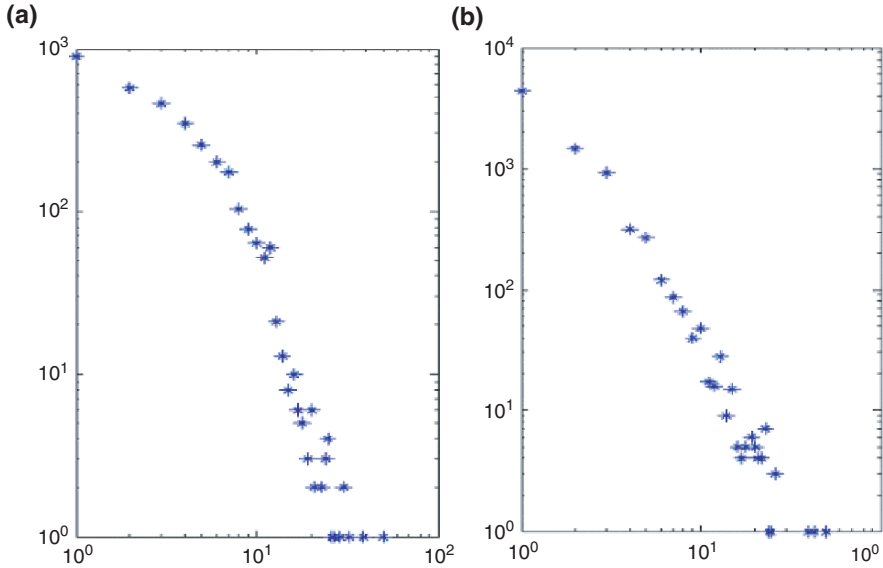


Fig. 13.7 (a) The distribution of temporal distances among the peaks of the total number of artifacts of a 40 agents system, during 20,000 steps (log-log scale). Note the power law tail, which has an exponent near -1.5 . (b) The distribution of crisis' magnitude of the total number of artifacts of a 40 agents system, during a 20,000 step evolution (log-log scale). Note the power law tail, which has an exponent near -2.8

It is possible to choose various measures of the difference between the two systems. In the following, we use the final area interposed between the perturbed and the unperturbed trajectories (see Fig. 13.8a). Three main kinds of perturbation are considered:

- the deletion of a randomly chosen agent,
- the deletion of a randomly chosen recipe, and
- a different sequence of random numbers (which functions as a control for the magnitude of the differences resulting from the first two cases).

We did 1,000 experiments (each experiment being a pair of runs) for each kind of perturbation. The following phenomena are observed above the level of system noise:

- a. the deletion of a randomly chosen agent is an event the system remembers; nevertheless after a while the perturbed system is able to reach a new equilibrium,
 - but the deletion of the first agent introduced into the system is a dramatic event: the perturbed system is not able to reach a new equilibrium and slowly moves to an increasing distance from the unperturbed system;

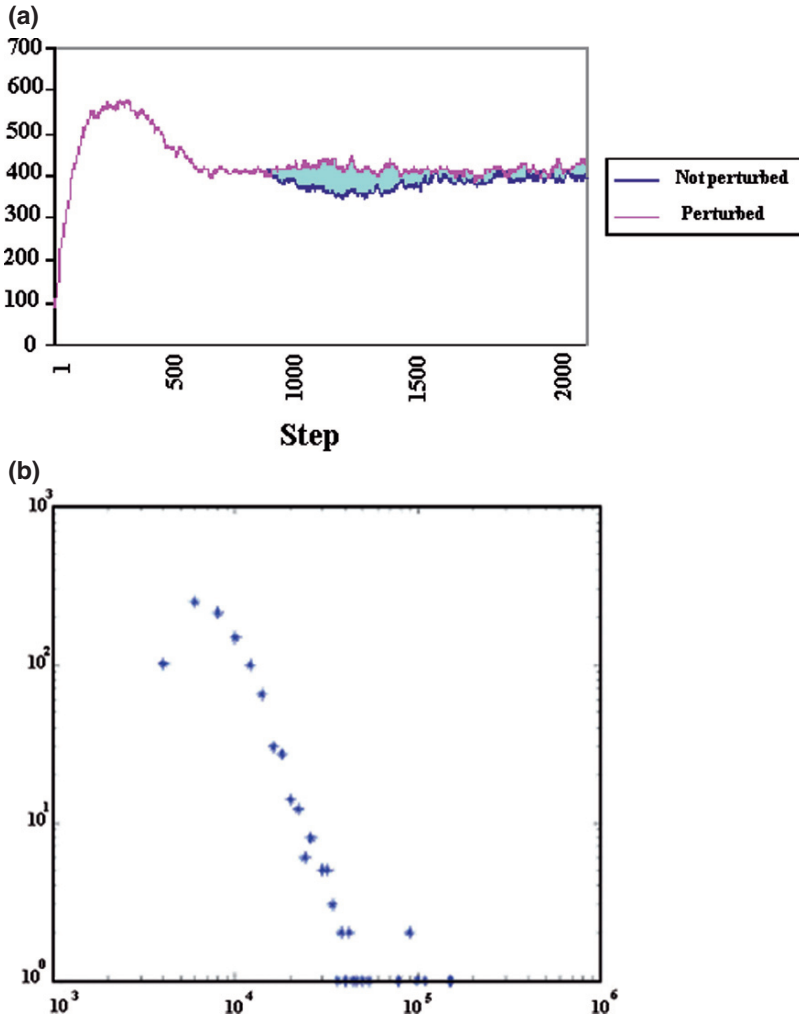


Fig. 13.8 (a) The final area interposed between a pair of perturbed and unperturbed trajectories; and (b) the final distribution of the final areas interposed between 1,000 pairs of perturbed and unperturbed trajectories (log-log plot)

- b. the recipe deletion is an event that does not effectively damage the perturbed system, which rapidly recovers the previous level; or
- c. a mere change of the random sequence does not effectively damage the perturbed system, which rapidly returns to the previous level (the results are similar to case (b), thus confirming that the recipe deletion has a small effect)

The fact that removal of an agent has bigger consequences than the removal of a recipe could be explained by remembering that one agent holds several recipes

(in fact, the big avalanches are caused by the elimination of the agents that have a number of recipes higher than the average), and that the deletion of an agent cannot be recovered. Even more interesting, it is clear that some agents are more important than other ones (the observation of point a, above): this aspect will be analyzed in depth in Section 13.6. On the other hand, the comparison between the deletion of a recipe and the change in the seed of random generator indicates that the deletion of a recipe is an almost negligible event (points b and c, above): a recipe can be easily rebuilt by the system.

Finally, we can observe that the final distribution of the measures (Fig. 13.8b) shows again a long tail, roughly linear in the log-log scale.

13.6 Stable Relationships

13.6.1 Introduction

A key issue of the set of theories upon which I₂M is built is the agents' ability to create (stable) relationships that foster innovations. In conditions of ontological uncertainty (Lane & Maxfield, 2005), the theory states that successful relationships should be based on the generative potential of the partners. That is, an agent should create preferential relations not with arbitrary agents, but only with the subset of agents that, at that moment, show high potential for generativity in relationships with the focal agent.

In I₂M the agents can maintain two kinds of relationships:

- an agent can “know” another agent (i.e. the first agent knows the existence and the outputs of the second one); this knowledge, if not subsequently enforced by means of an artifact exchange, has the duration of only one step, and
- an agent can have artifact exchanges with another agent (being one of its providers).

The second kind of relation is supported by the agent's recipes, the presence of which guarantees the temporal stability of the relationship. Therefore, exchange of artifacts involves stronger relations than those due to “acquaintances.”

These acquaintances allow the agent to locate its own goal and to choose the partner to realize it; therefore, it is useful that advantageous acquaintances be maintained while the less convenient ones be discarded. But how should an agent recognize advantageous acquaintances?

13.6.2 The Importance of Past History

We considered different situations, where the choice of the partners is either random or based upon the agent's history. This latter alternative is implemented as a vote, which each agent gives to each other agent it knows; the vote is a function of the

history of the relationships between the two. In general, it is the sum of two terms, one related to the history of buy-and-sell relationships, the other dependent on the history of attempts to develop a joint project (i.e. to share the inputs and recipes in order to reach a common goal).

Let us consider the model results when the vote is only influenced by the projects the two agents did together in the past. In this case, the vote dynamics is very simple:

$$V(t+1) = V(t) + \Delta_t - \lambda V(t), \quad (13.1)$$

where

$$\Delta_t = \begin{cases} +2 & \text{if there is a successful joint project at time } t \\ 0 & \text{otherwise} \end{cases}. \quad (13.2)$$

Unknown agents are given $V(t) = 0$, and occasional acquaintances (agents just known by chance) are given an initial $V(t) = \varepsilon$, with $\varepsilon < 1$. The third term of the equation is a forgetting term, which lowers the vote of those agents that no longer engage in partnerships. When an agent has to choose a partner to collaborate with the intention of realizing its goal, the choice is probabilistic, based on the votes.

The effects of this vote attribution mechanism is shown in Fig. 13.9, where we can see the vote given by a certain agent to the other agents and the relative table of presence/absence of stable collaborations (a stable collaboration being a relationship in which the vote is higher than Δ_t for hundreds of steps). Note that stable collaborations can arise, be interrupted, and later start again. The vote and the partnership mechanisms are able to establish what we can interpret as reciprocal trust (as manifested in reciprocal high votes).

13.6.3 First Results

Now we can compare the results obtained by the voting system with the results obtained by randomly selecting partners. The system has several variables, not all clearly implicated in this change of strategy: the major effects are evident on the variables more strongly involved in relationship processes. For example, it is possible to compute the number of projects (i.e. the innovations that reach the predefined goal by combining the recipes of two different agents) realized by each agent and, subsequently, to plot and compare the resulting distribution. The major effects of the introduction of the voting system are visible in Fig. 13.10a (in which the distributions are calculated by cumulating the results from 30 runs).

For this figure, note that:

1. The median and the average of the two distributions are similar; and
2. The distribution corresponding to the voting system has a more dispersed shape, indicating the existence of a subset of agents able to create more projects than the other one, at the price of blocking the creative action of another, more numerous subset.

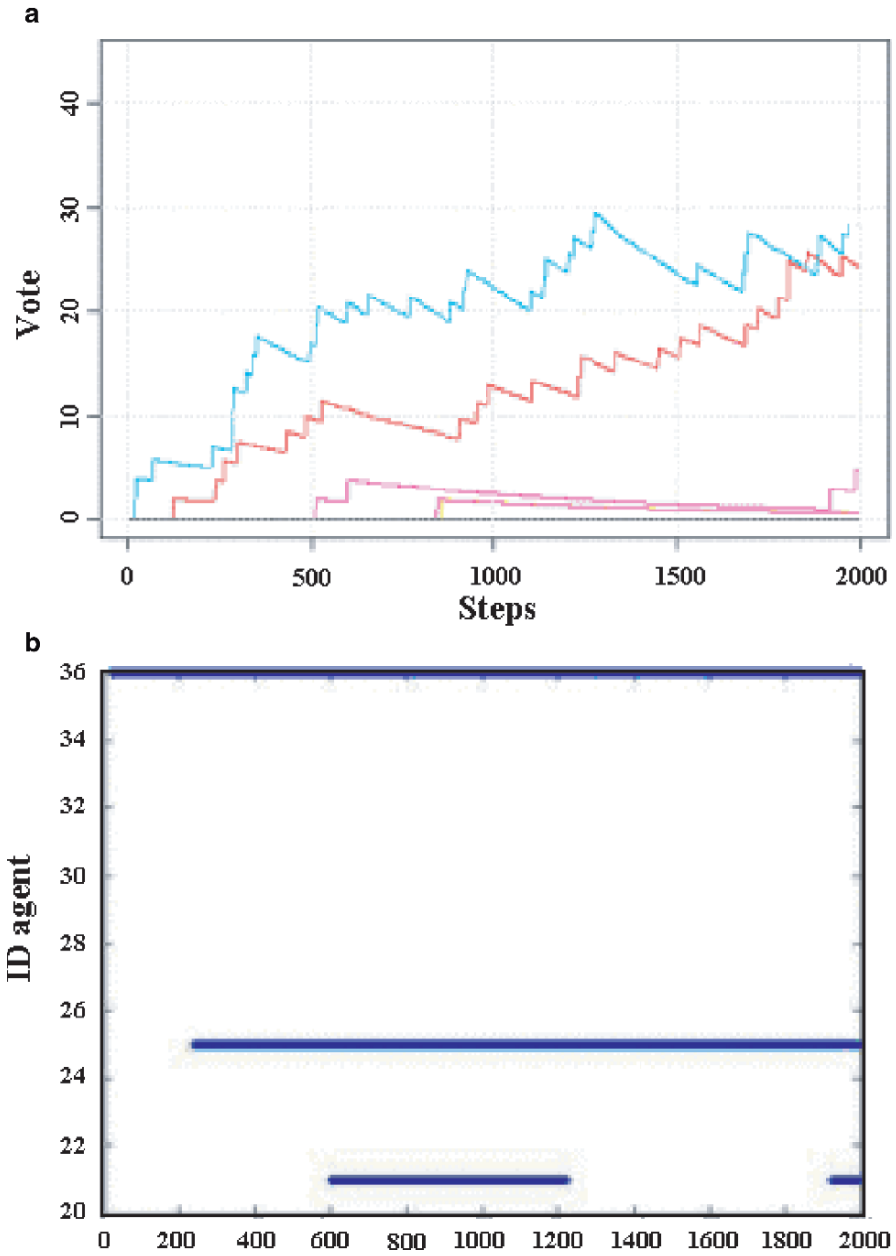
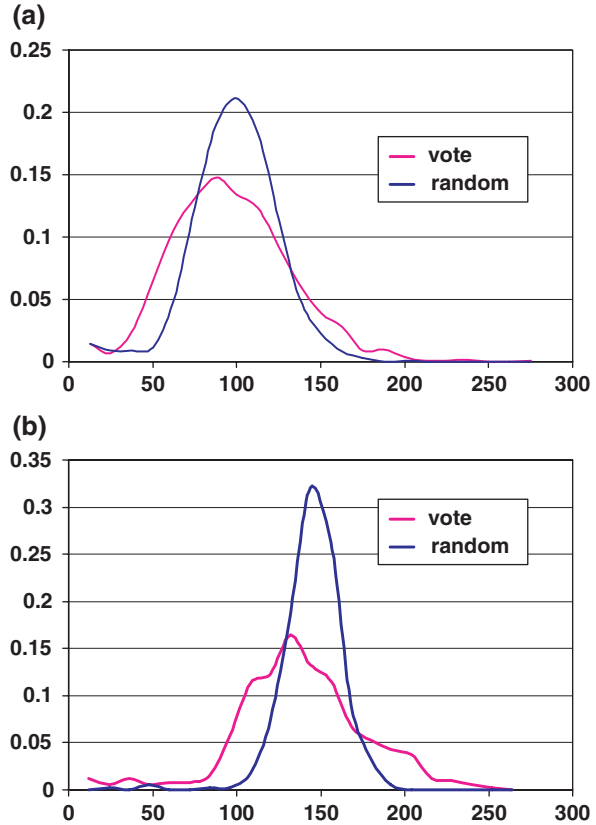


Fig. 13.9 (a) The vote given by agent 11 to the other agents during a 2,000 step run; and (b) the relative table of presence/absence of stable collaborations (the blue lines corresponding to the number of the collaborating agents)

Fig. 13.10 Distribution of the number of projects (i.e. the innovations that reach the predefined goal by combining the recipes of two different agents) realized by each agent, with (*red line*) or without (*blue line*) a voting system. **(a)** Summation of 30 model runs with standard parameters; and **(b)** summation of 30 model runs with rationally limited agents



This second characteristic is strengthened significantly by conditions of ontological uncertainty. The agents of Fig. 13.10b suffer from great reasoning limitations (the performance of their “genetic engines” have seriously deteriorated, passing from a population of 80 individuals and 40 generations to a scheme with 4 individuals and 2 generations). It is possible to observe that:

- To realize their goals, the agents have to make more frequent use of partnerships (the averages of the two distributions increase greatly);
- The distribution of agents randomly choosing their partners is higher and proportionately narrower than in the previous situation;
- The more intense recourse to joint innovation processes enhances the advantages as well as the disadvantages of utilizing the characteristics of the voting system. In particular, the fraction measuring the area between the right tails of the distributions doubles, passing from 7% to 13% (at the price of a similar increase of the area between the distributions’ left tails, which passes from 13% to 25%). In other words, the number of agents able to exploit the characteristics of the voting system increases at the price of an increased number of agents badly connected and debilitated.

It is worth observing that, in accordance with the statements of the I₂T, the systems where relationships matter and ontological uncertainty plays an important role foster the existence of particularly connected agents, able to exploit the generative nature of their relationships.

13.7 Structures in Artifact Space

In I₂M, different entities are present, and among these entities, several kinds of interactions take place. Generally, if we have to deal with entities and binary interactions, it is possible to schematize the corresponding systems by means of graphs (networks); several I₂M interactions have a binary shape, and, therefore, in our model, it is possible to define various interesting networks.

In this section, we focus our attention on a particular network that is able to reveal some details regarding the structure of the artifact space: the artifact type network. The artifact type network is a graph where two different numbers are linked if there is at least one recipe within the system that uses an artifact corresponding to the first number in order to build an artifact corresponding to the second number. In this case, the link is directed from the first number to the second one. By means of network analysis, it is possible to address questions regarding the global structure of the systems, by abstracting from the details of each single interaction.

In particular, we are interested in observing in our system the existence of so-called “technological waves,” that is, the sequential emergence of strongly inter-related sets of artifacts. A real example of these waves is the succession of the technologies related to “maritime trade” (the epoch of great maritime discoveries), “manufacturing coal and steel” (first industrial revolution), “manufacturing electricity and automobiles” (second industrial revolution), “informational and communication activities” (the contemporary period) (Chapter 8, this volume). The presence of these waves is important because it reveals the possibility of a perpetual novelty, where interacting systems of new kinds of artifacts can substitute already existing ones. With this aim, it is possible to distinguish several interesting questions, as for example:

- Are there technological waves present in our system?
- If so, under what conditions can they appear?

To investigate these questions, we have to identify particular subsystems within the whole artifact network. There are several network analysis techniques able to identify, at least approximately or partially, such kind of organizations: for example, strong component analysis (Scott, 2000), k-core analysis (Scott, 2000), and the percolation of k-cliques (Palla, Derényi, Farkas, & Vicsek, 2005). In the experiments we have discussed so far in this chapter, all these analyses identify only one organization, which emerges during the early simulation steps and is never replaced by other systems (see Fig. 13.11a for an example). This is the typical outcome when the starting materials are very near each other: a very stable world, in which newcomers

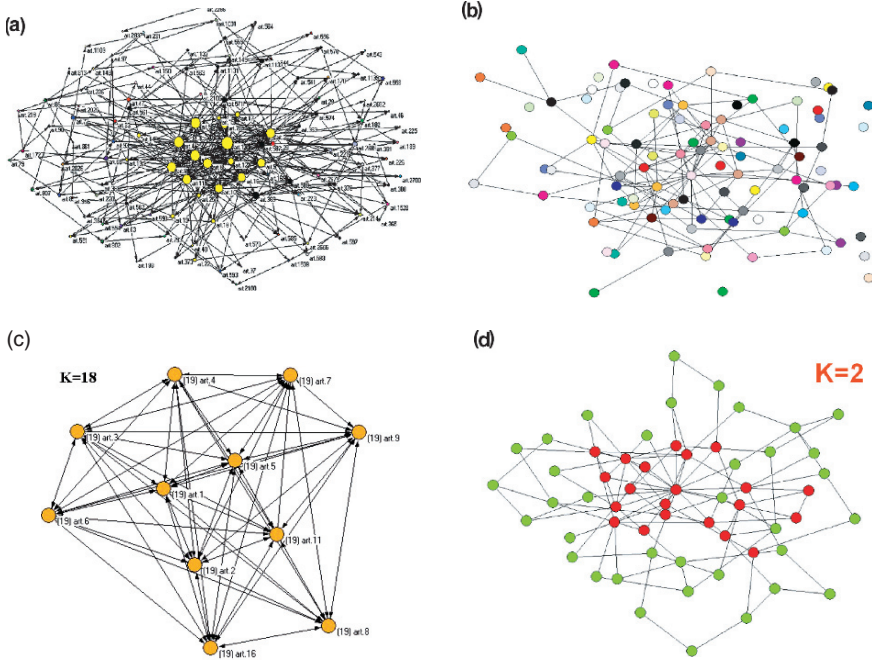


Fig. 13.11 The outcome of several analyses performed on the “names” network. **(a)** Strong component analysis: a huge strong component predominates, embracing almost all the nodes, whereas all other components are formed only by one or two nodes. The vertex sizes are proportional to the number of paths that go through the nodes. **(b)** The same as **(a)**, but in presence of two clusters of raw materials: there are no strong components. Each color indicates a different strong component (additionally, the visualization tool has a limited number of colors that are repeated for different components: in effect each strong component on the picture is formed by only one node) **(c)** K-core analysis; there is only one core for each value of k , up to the very high value of 18. The identified nodes are the same nodes that prevail in **(a)**. **(d)** The same as **(c)**, but in presence of two cluster of raw materials: we find a core only for values of k lower than 3 (in **(b)** the maximum values of k is 18)

do not replace, in a wholesale way, the basic structure constructed by already existing artifacts.

However, preliminary simulations of I_2M highlight the possibility that in addressing these issues, it is important to consider the diversity of the starting materials (the raw materials of our system). When these artifacts are not agglomerated in only one cluster, but, in contrast, show more than one preferential position, the final situation is quite different (see Fig. 13.11b, c).

Even more interesting is the evolution of the outcome of the communities identified by means of the percolation of k -cliques. The communities identified by this kind of analysis are the organizations closest to the definition of “technological waves,” and it is possible to observe the subsequent replacement of several communities in time. Last, but not least, we can observe the existence of more than two

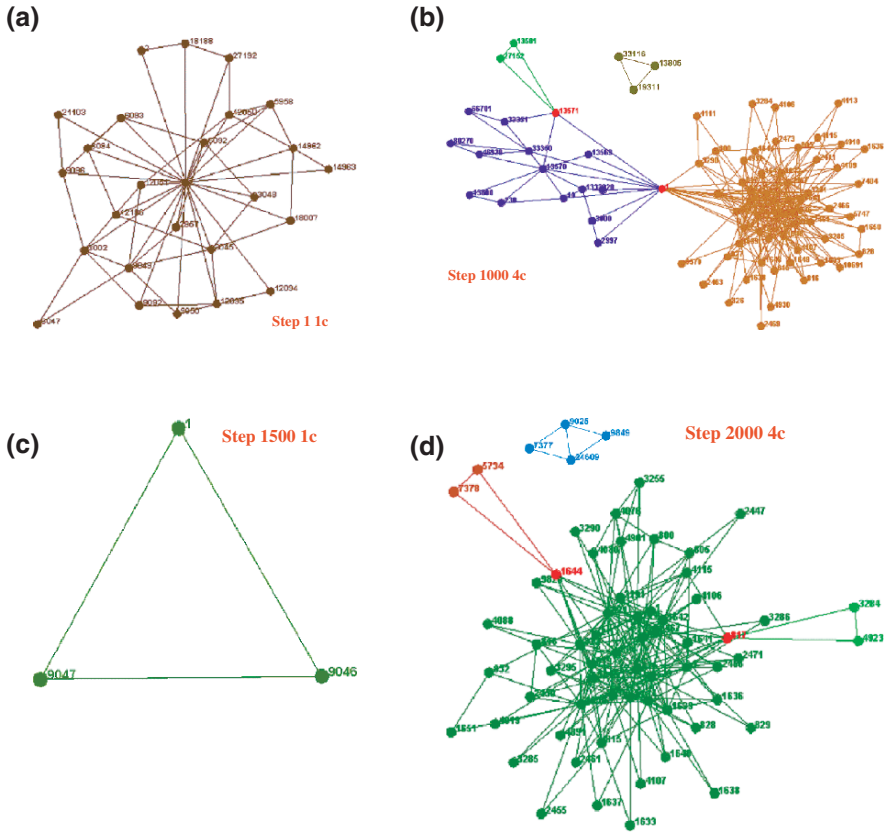


Fig. 13.12 The evolution of the outcome of the communities identified by means of the percolation of 3-cliques. It is possible to observe the subsequent replacement of several communities in time, including the presence of a crisis (c) and a subsequent recovery period (d). Last but not least, we can observe on occasion the coexistence of more than two communities, also if the initial clusters of row materials are only two: it is a clue toward the existence of non linear effects

communities, even when there are only two initial clusters of raw materials (see Fig. 13.12), indicating the possible existence of strongly non-linear effects.

13.8 Conclusions

After having presented some of its main features, we now consider what kind of model I_2M is. There are, of course, very many meanings of the word model (we learned from our colleagues in biology that for them even a rat or a mouse may be a “model” of a human organism!), but we will confine ourselves to mathematical or computer-based models, to which I_2M belongs.

We can compare it to some selected models, highlighting common features as well as differences, although by no means will we provide an exhaustive review or classification of these various models.

Let us first consider the Alchemy system (Fontana & Buss, 1994), which is loosely inspired by ideas concerning the origin of life and the role of self-maintaining collections of molecules. Alchemy is a kind of “abstract chemistry,” where there are functions instead of molecules. The formalism of lambda calculus captures a kind of structure-function relationship that parallels that of molecules.

In the course of time, interesting self-organizing phenomena can take place in the Alchemy reaction dynamics. In particular, cycles are formed, which are composed of functions that catalyze each other’s production. When one observes the interactions that take place among these high level entities, one finds algebraic structures whose behavior can be described without making reference to the low level ones that give rise to them.

Spontaneous cycle formation also is observed in I_2M , as described in the previous sections. Apart from technical differences, both models take a very abstract view of the phenomena that stimulated their conception. In a sense, Alchemy’s viewpoint is more abstract, as it gets rid of a number of the fundamental features of the biological world that inspired it. On the other hand, it is also more focused on a specific (abstract) kind of question, i.e. those that relate to the formation and interaction of high-level entities. Neither of the models aims at providing meaningful comparisons with actual quantitative data, but, rather, they provide proof-of-principle that interesting emergent phenomena can take place when many microentities interact.

A completely different model is that of a chain of macroscopic harmonic oscillators, where one of them has a mass, M , much larger than that of the others (which have all the same mass, m). In this case, it is possible, assuming that the law of classical mechanics hold, to write a set of exact equations of motion. If we suppose that the initial conditions are not precisely known, the deterministic description is substituted by a law (Liouville equation), which governs the time evolution of the probability distribution of the positions and velocities of all the particles. Let us now suppose that we are only interested in the behavior of the heavy particle, whose motion is slower than that of the others. We can then project the Liouville equation onto the subspace of the variables that describe the heavy particle. This projection takes the form of a perturbative expansion (Serra, Zanarini, Andretta, & Compiani, 1986), in which the leading term (Fokker-Planck equation) describes the heavy oscillator in the same way as a Brownian particle. This approximation is valid on a time scale that is shorter than that of the slow oscillations and fast with respect to that of the light oscillators.

The model provides an approximation to the behavior of the slow oscillator, which is based on an established theory. Moreover, estimates of the accuracy of the approximation, as well as correction terms, are also provided.

Both models (the oscillator chain and I_2M) are based on a theory, but the former is quantitative, while the innovation theory is qualitative and verbal. It is worth noticing that, in both cases, unanticipated consequences of the theory can be observed: this is obvious in the I_2M case, but the emergence of Brownian motion

in the oscillator chain also can be considered an unanticipated behavior. We observe that, in both cases, comparison with experimental data is not particularly relevant, although for different reasons. In particular, for the oscillator chain the reliability of the underlying theory provides a firm foundation for the model-based claims.

As a final example of a different kind of model, consider random Boolean models of genetic networks (Kauffman, 1993). These are based upon a description of the phenomenon (the regulation of the expression of various genes in the cell), but, in order to get a manageable model, many simplifications are introduced. Some of them represent approximations to the observed variables (like the use of Boolean functions for gene expression levels or the neglecting of protein dynamics while resorting to a “gene only” model), while other simplifications represent lack of knowledge by introducing randomness (e.g., the choice of the connections among genes is done at random, and also the Boolean functions are chosen at random).

In the case of gene regulation network models, it is possible to draw some general conclusions, like those that concern the relationship between number of different cell types and genome size (Kauffman, 1993), or the distribution of the size of perturbations in expression levels induced by knock-out (Serra, Villani, & Semeria, 2004; Serra, Villani, Graudenzi, & Kauffman, 2007). By comparing the model results with actual observed data, it is thus possible to justify *a posteriori* the simplifications that have been introduced.

Even here, the basic description of the mechanisms of gene regulation may be considered a partial theory of the phenomenon, and the model is rooted in this theory. Similar to I_2M , the model’s purpose is that of unfolding the systemic consequences of assumptions concerning the individual gene behavior. The main difference between the two cases is that experimental data are available for random Boolean networks. However, it is interesting to notice that, for about 30 years, the only data of this type have been that concerning the relationship of different properties (cell cycle length, number of cell types) with genome size, and that this is a highly questionable measure, given the uncertainties about the number of actual genes. In spite of this limitation, the model has been widely used as an approximation to the functioning of real cells and as a tool to explore some candidate generic features (e.g., dynamical criticality of real cells).

Briefly, we can observe that the three models considered above, as well as the innovation model, are based on a set of hypotheses concerning the microscopic interactions and are used to explore the global properties that emerge from these interactions. In a sense, every agent-based model of social or economic systems shares this property, and, therefore, is based on an implicit theory of the phenomena at hand. However, most models of this kind do not present a clear distinction between what the theory claims and what the model is and shows. A major part of our effort in the I_2M case is that of providing an explicit discussion of the relationship between model and theory.

The usefulness of this model will ultimately rest on its capability to activate a dialogue with the theory, in order to improve the latter. From this perspective, a major difficulty comes from the existence of several specific mechanisms, which are

necessary to make the model work, but which can affect its outcomes. Therefore, one of the major aims of future research will be to disentangle those contributions from others that come from basic theoretical assumptions.

It is also clear that the model still lacks some of the features of the theory: notably, the choice of the partner and the setting of the goal, which might be made more sophisticated than the naïve mechanisms so far introduced. In addition, the drive to innovation is questionable, since our agents are natural born innovators, and are not led to look for novelties due to specific reasons rooted in their particular circumstances. Therefore, the model could be modified by incorporating these characteristics. However, a higher priority should be attributed to a simplification of the model itself, making it free from perhaps unnecessary complications in handling the customer-supplier relationships, to make it more comprehensible without hiding its key innovation mechanisms.

References

- Axelrod, R., & Tesfatsion, L. (2006). A guide for newcomers to agent-based modeling in the social sciences. In L. Kenneth, K. L. Judd, & L. Tesfatsion (Eds.), *Handbook of computational economics, volume.2: Agent-based computational economics* (pp. 1647–1658). Amsterdam, The Netherlands: North-Holland.
- Bak, P. (1996). *How nature works*. New York, NY: Springer.
- Epstein, J. M., & Axtell, R. (1996). *Growing artificial societies: Social science from the bottom up*. Cambridge, MA: Massachusetts Institute of Technology Press.
- Fontana, W., & Buss, L. W. (1994). What would be conserved if ‘the tape were played twice’? *Proceedings of the National Academy of Sciences*, *91*, 757–761.
- Gilbert, N., & Terna, P. (2000). How to build and use agent-based models in social science. *Mind and Society*, *1*, 57–72.
- Kauffman, S. A. (1993). *The origins of order*. Oxford, UK: Oxford University Press.
- Lane, D., & Maxfield, R. (2005). Ontological uncertainty and innovation. *Journal of Evolutionary Economics*, *15*, 3–50.
- Lane, D.A., Serra, R., Villani, M., & Ansaloni, L. (2005). A theory-based dynamical model of innovation processes. *ComplexUs*, *2*, 177–194.
- Palla, G., Derényi, I., Farkas, I., & Vicsek, T. (2005). Uncovering the overlapping community structure of complex networks in nature and society. *Nature* *435*, 814.
- Scott, J. P. (2000). *Social network analysis: A handbook*. (2nd ed.). London, UK: Sage Publications Ltd.
- Serra, R., & Villani, M. (2006). Agents, equations and all that: on the role of agents in understanding complex systems. In M. Schaerf, & M. O. Stock (Eds.). *Reasoning, action and interaction in AI systems and theories. Springer Lecture Notes in Computer Science 4155*, 159–175.
- Serra, R., Villani, M., Graudenzi, A., & Kauffman, S. A. (2007). Why a simple model of genetic regulatory networks describes the distribution of avalanches in gene expression data. *Journal of Theoretical Biology*, *249*, 449–460.
- Serra, R., Villani, M., & Semeria, A. (2004). Genetic network models and statistical properties of gene expression data in knock-out experiments. *Journal of Theoretical Biology*, *227*, 149–157.
- Serra, R., Zanarini, G., Andretta, M., & Compiani, M. (1986). *Introduction to the physics of complex systems*. Oxford, UK: Pergamon Press.

Chapter 14

An Agent-Based Model of Information Flows in Social Dynamics

Davide Ferrari, Dwight Read and Sander van der Leeuw

14.1 Introduction

Elsewhere, one of us recently argued for the need to develop a wide range of modeling approaches in the social sciences (van der Leeuw, 2005), including archaeology and anthropology. This call reflects a growing interest, on the part of archaeologists and anthropologists, in constructing and using dynamic models, and a growing interest on the part of modelers in the domains these social sciences cover (e.g., Gilbert, 1991; van der Leeuw & McGlade, 1997; Ballot & Weisbuch, 2000; Janssen & Jager, 2000; Kohler & Gumerman, 2000; Janssen, 2002; Anderies, 2006; Kohler & van der Leeuw, 2007; etc.). In this paper, we will attempt to model a specific phenomenon that has occurred, and occurs, in all societies: the creation of novelty (combining invention and innovation).

The model is based, in part, on a case study by van der Leeuw and others among potters in several villages near the town of Pátzcuaro, State of Michoacán, Mexico, around 1990 (cf. van der Leeuw & Papousek, 1992; van der Leeuw Papousek, & Coudart, 1992). In this region, the repertoire of products varies considerably from one village to the next. In one, most potters make almost exclusively globular ‘ollas’ to keep water or other liquids. In another, the potters specialize in ‘cazuelas,’ deep, open dishes (with lids) in which one cooks. Yet other villages produce a much wider range of products, such as the famous village of Tzintzuntzan (Foster, 1948), where one finds plates, cups, saucers, dishes, jars, etc., as well as a wide range of purely decorative pots.

Our study is aimed at identifying and studying how artifacts mediate in the dynamic relationship between the ‘mental models’ of the makers and those of the users of the local ceramics. By studying the traces that the manufacturing techniques and procedures left on the pots themselves (following an approach developed in archaeology, see van der Leeuw, 1976), we succeeded in describing, systematically, the making of each of the products of the potters. That gave us the insight in the logic behind the manufacturing techniques to enable us to outline the mental model of

D. Ferrari (✉)

School of Statistics, University of Minnesota, 313 Ford Hall, 224 Church Street,
Minneapolis, MN 55455, USA

pottery making which the Michoacán potters refer to in their work, and which is the basis for their categorization of their products.

That model first distinguishes between 'open' shapes (plates, dishes, bowls, cups, etc.) and 'closed' shapes (ollas, jars, water storage vessels, but also many decorative vessels). The distinction is closely linked to the method of shaping these vessels. The 'open' ones are shaped over a mushroom-shaped mould that consists of one piece only, whereas the 'closed' shapes are molded in two or more vertical sections. Within each of these two major categories, other distinctions refer to the number of pieces of clay of which a vessel consists and to size and shape of the molds used. One way to summarize the range of shapes from the makers' perspective is by referring to 'rotational symmetry breaking,' i.e., to the ways in which the objects deviate from a perfect globe (represented by the ollas), either in the horizontal plane or the vertical one, or both. Rotational symmetry breaking in the horizontal plane can deform a circle into an ellipse, an oval, a square, a rectangle, a trapezoid, or any shape in between. Symmetry breaking in the vertical plane can truncate the circle, or deform it into a cone, a cylinder, or any number of other, more complex shapes. A description that combines the two kinds of rotational symmetry breaking can capture the whole range of basic shapes made. Secondary aspects of the makers' categorizations refer to additions, such as handles, spouts, rims, etc., but also to the decoration on the pots, or the way they are fired, etc.

Interestingly, the potters' mental model is very different from that of the users, who primarily refer to the function of the pot (carrying or storing water, cooking, frying, serving, pouring, drinking, etc.), to the occasion in which the pot is traditionally used (such as the piñatas), or to its size (whether it serves two, four, six, eight people, etc.). The difference raises the question how makers and users of the pottery come to a mutual understanding about what are considered 'good pots.' That understanding is, of course, one of the important factors determining what kinds of pots are found in the villages at any one time and what kinds of innovations are likely to occur in their repertoire. From the perspective of innovation, we need to draw the reader's attention to other differences as well. First, there is an important difference in the frequency of innovation between the manufacture of the pots, on the one hand, and their decoration on the other. The reason for the difference resides in the fact that in the case of manufacture, an error is usually irreversible (i.e., the object breaks or cracks and, therefore, is non-functional), whereas in the case of decoration, (a) errors can be corrected and (b) even if the decoration is not up to par, the pot can still be used for many functions. The rules of manufacture are thus prescribed in detail (and these prescriptions are 'hard-wired' in the shape of the molds that are used), while in the domain of decoration there is no detailed prescription, but rather a set of general rules combined with the proscription of certain things that are not deemed acceptable. As a result, potters are much freer to express themselves in the decoration, and in certain villages there is much more variety in that domain. This variety points to the fact that the nature of the raw materials, the tools, and the rules for using both play an important part in determining the potential for innovation.

However, the situation also differs from village to village. In a village (Capula) that is adjacent to the old Spanish city of Michoacán, Morelia, the traditional

manufacturing techniques have been adapted to a completely different range of shapes, while in the villages further away from the city, the potters make only the traditional shapes of the region. In part, that seems linked to the fact that, in the more remote villages, interaction primarily is among the inhabitants of the village, and ideas do not change unless they reach a consensus. There is, thus, an important degree of social control, and the potter who innovates risks social opprobrium. Near the city, on the other hand, many interactions occur among individual potters in the village and people in the city. Consensus among the village inhabitants plays a much smaller role, and innovation is, thus, (socially) easier. One is struck by the extent to which communication and consensus formation (alignment) determine the extent to which innovation can occur.

Although this brief and highly schematic description does not do justice to the complexity of the dynamics of categorization, communication, and innovation among potters in Michoacán, we do not have the space here to delve more deeply into the matter. We hope to have made the point that there are many different feedbacks between the social, the technical, and the functional aspects of pottery manufacture and use. The model we present in this paper is an attempt to begin looking at these in some detail, by providing an experimental tool that enables us to simulate the dynamics of interaction between people and, in particular, between makers and users of artifacts, and the role that these artifacts play in that interaction. Although the ideas that underpin it have been developed, in part, in the context of the Michoacán research, we claim that they are much more generally valid. They seem to be applicable to other (past and present) instances of pottery-making (e.g., van der Leeuw, 1976, 1993), but also to the manufacture of prehistoric stone tools in many parts of the world (e.g., Böeda, 1986, 1990; Pigeot, 1992; van der Leeuw, 2000), as well as to contemporary industrial manufacture, for example in the Italian Business Districts (Russo and Rossi, this volume). To optimize its usefulness, the model is therefore formulated in very general terms.

14.2 Categorization, Communication and Material Culture

In the most general of terms we can state that often, interactions between particular (groups of) agents¹ take the form of more or less recurring patterns that persist over time. Making (new) artifacts, for example, or generating (new) forms of organization in a society is a multifaceted process based on recurrent sets of actions, in which a number of agents with different roles are involved. In this process, new concepts and objects are discovered and invented by simultaneously surveying the opportunities offered by the social and material environment and adapting the existing categorizations of the agents. When looking for inventions (new strategies, new products,

¹ The term 'agent' is here used to refer to a social entity that engages in interactions with, and collective actions of, a larger collection of similar entities. Thus, an economic agent can be a person, a small group, or a larger, more structured set of individuals acting as an organization.

or organizational innovations) the different agents participating in this process each explore a variety of opportunities from their own perspective (Papousek, 1989). But for an innovation to emerge, there must also occur a degree of convergence among their perspectives so that the group, as a whole, aligns itself around the creation and use of the invention. The members of a society are engaged in pursuing ways to produce novelties (and, thus, new functionalities), and a number of convergent processes and transactions therefore take place among them.

Social dynamics can be viewed profitably in terms of exchanges of matter, energy, and information. In those exchanges, these commodities play different roles because they have different properties. The former two commodities are subject to the laws of conservation and can be transferred from one individual to another, but cannot be shared simultaneously. Information, however, can be, and is, shared. It follows that, from this perspective, information exchanges hold groups of individuals or societies together and may lead to convergent behaviors among specific groups of individuals. In other words, societies are held together and innovate because their members share a common system of beliefs, which enables them to interact and otherwise express themselves in ways that they mutually recognize (e.g., do things in compatible ways, identify and use similar resources, share technology, have similar norms and similar institutions). When attempting a representation of the social dynamics involved in the exchange of information, the scientist must face two specific questions. First, how can we conceive a notion of information that is suitable for describing the interaction among social agents? And second, is there a reasonable paradigm to describe the ways information is circulated? Neither has been answered in a satisfactory manner so far. In fact, there is significant uncertainty about which notion of information we should use when modeling such processes. In the last few decades, several anthropological studies have shown that individuals process information with the aid of categories (e.g., Selby and El Guindi, 1976); in addition, a large body of literature in social psychology has emphasized the central role of categorical thinking in social relationships (e.g., Markman & Gentner, 2001; Macrae & Bodenhausen, 2000). However, in other social sciences (such as economics and sociology) the potential of the representation of information in terms of categories has yet to be realized.

The process of categorization has generally been modeled from the perspective of the cognitive capacities underlying category formation, both from a developmental viewpoint that considers how a newborn child learns to form concepts, upon which categorization is based, and by asking how categories might be represented cognitively and new categories derived from existing ones. This work generally has focused on how categorization of external phenomena may take place through identification of properties or characteristics that are deemed salient for category definition. In this framework, the external world presents itself to the individual, and one of the primary goals of categorization is to discern patterning in external phenomena, which can be embedded in category definitions. Definitions that then enable effective interaction with the external world through properties or characteristics attributed to category members by virtue of entities being identified with a category. The categories that we construct in this manner can be quite broad and may make

differentiation among the phenomena included in a category, such as a category of hot objects, or can be quite specific and narrowly defined, depending on the context and degree of specificity required for the task at hand. Functionality can be associated with a category by considering what actions or behaviors are appropriate either in terms of how we interact with objects associated with the category (e.g., ‘don’t touch hot objects’) or how we might make use of objects categorized in this manner (e.g., a hot flat piece of metal with a handle can be used to iron clothes).

This approach to categorization of external phenomena has not addressed the complexities that arise with what anthropologists refer to as material culture, namely those material objects that are produced or manufactured with the intention that they will be used by other persons. Whereas categorization of material objects in the external world begins with already existing entities and is based on properties of these entities, the objects of material culture do not exist except through raw material being acted on by the producer, who is guided by her/his categorizations of raw material, manufacturing techniques, and so on. A producer does not simply act on raw material, but is guided in her/his actions by the goal of producing objects that, from the producer’s viewpoint, will be instances of a category of artifacts believed to be of utility to users. In brief, the producer is not just making things, but kinds of things. The actions of the producer in making an artifact are guided by what the producer intends to make and her/his knowledge of the actions that need to be taken in order to produce an instance of a particular kind of object. What is produced also relates to categorizations held by the user. The user, for her/his part, does not merely categorize what has been produced by the producer, but is guided by categorizations that relate to the intended use or activity in which the artifact might be employed. These categorizations are not independent of the producer and the producer’s categorizations, but are formed in part through interaction with the producer.

In theory at least, the categorization of already existing artifacts proceeds from the observation of the artifacts to the categories in accordance with the patterning perceived, but it is independent of pre-existing categories. But the making of artifacts is part of four modalities for interaction:

1. The interaction between the producer and the raw materials used places constraints on what can be produced.
2. The transmission of artifacts from producer to user (which, in most cases, is not accompanied by direct interaction) proposes the producer’s conception of the products to the user.
3. The execution by the user of the tasks for which the artifacts are designed provides a context for the evaluation of the suitability of the artifacts for the task.
4. The user’s communication to the producer about the suitability of the artifacts (in the form of a decision whether or not to use more of the same artifacts or in the form of an actual information exchange between the user and the producer) closes the loop.

Except in the case where the user and the producer are one and the same person, the categorizations employed by the producer as part of working with raw material

and forming an artifact of a particular kind, and the categorizations employed by the user with regard to the activities or tasks for which the artifact may be employed, will only partially overlap in terms of their respective category definitions. The overlap is not coincidental, nor is it simply due to the fact that both producer and user are categorizing the same artifacts, but it is an integral part of the triad made up of producer, user, and artifact, in which the producer and the user are engaged in communication mediated by artifacts. This communication entails searching through their respective categorizations and possible modifications of their categorizations until convergence occurs between the categorizations of the producer in terms of artifact production and the categorizations of the user in terms of the artifacts' suitability for the tasks and activities in which they will be involved. Stasis occurs when the respective categorizations are satisfactory from the viewpoint of the producer and the user, but does not require that the respective categorizations be identical.

Neither the producer nor the user determines the outcome of this process, nor is the outcome completely predictable by reference to the external conditions relevant to the use of the artifacts, as only a few of the possible dimensions for characterizing artifacts are sufficiently constrained by external considerations to dictate what values they will assume. A knife must be sharp for the task of cutting, for example, hence, the angle of the knife's edge is bounded by the task of cutting and the material from which the knife is made, but what material will be used, the size of the knife, the form it has, whether it might be decorated, etc., are not determined by the knowledge that the knife needs to have an edge that cuts effectively. Instead, what constitutes these other dimensions and the numerical or qualitative values they assume for any artifact relates to the respective categorizations of the producer and user, which is negotiated when the producer conveys to the user his current categorizations in the form of the artifacts, and the user conveys to the producer an evaluation of those artifacts based on their utility in doing the tasks at hand. The process will end when there is sufficient alignment between the two sets of categorizations so that the constraints of the producer in what s(he) can produce and the constraints of the user with regard to the task or activities at hand are both satisfied. Alignment needs not, and almost surely will not, imply identity between the sets of categorizations as there is no reason for the user to be aware of the full range of categorizations relevant to the producer, nor for the producer to be aware of all the categorizations of the users of the artifacts.

This process of convergence cannot be understood by focusing on either the producers or the users alone, as convergence takes place through imperfect communication of incomplete information between producer and user, in which the artifacts play a mediating role. We aim to model, through simulation, this process of convergence in the respective categorizations through exchange of information. The remainder of this paper proceeds as follows. The next section will present the basic features of the model and its main assumptions. Section 14.3 outlines the entities of the model and their properties. In Section 14.4, we describe the model dynamics at the macro, meso and micro levels. In Section 14.5 we present a first version of a computer implementation of the model, along with some results of our simulations. We conclude with a discussion of the main results.

14.3 The Model

14.3.1 Objectives of the Model

The aim of the model is to represent the dynamics involving invention (the generation of new categories) and communication (transmission of categories instantiated as real objects), which has the properties just outlined. In particular, our model links phenomena at three different levels:

The micro-level: At the level of the individual, we propose a simplified description of how individuals collect information about the external world, categorize it, and combine existing categories in order to create new ones. As part of this process, it is assumed that information is encoded at different levels of complexity (so that categories encompassing a smaller number of cognitive features represent ‘simpler’ real objects).

The meso-level: The second level concerns exchange of information among the individuals. In our model, individuals act in accordance with their social role (as either producer or user of artifacts), and the information is exchanged in terms of ‘packages of categories’ instantiated in the artifacts that are exchanged between individuals in the society. Interactions among individuals may generate new, more complex information, and different communication patterns will produce differences in the information encoded by individual members. This is a key element that ensures the generation of new ideas at the individual level.

The macro-level: The last level defines some of the characteristics of the society as a whole, such as any limits to the kinds of actions that its members are willing to undertake. These affect the categorization process and play a role in the interactions between individuals and groups. The macro-level can be thought of as defining the shared system of beliefs and the common physical and technological resources. The third level encompasses the previous two, and its dynamics play out on longer time-scales.

We insist on the fact that this model does not pretend to represent any of the detail of the cognitive processes underlying the social interactions among individuals. Instead, we provide a simplified description of such processes, enabling us to approximate phenomena involved at the longer time scales commonly used in archeology.

14.3.2 Model Description

In the model, the cognitive space is represented by a set of dimensions and the cognitive space of a particular agent is the set of cognitive categories defined by using a subset of these dimensions. For the sake of simplicity, a category will be characterized in the model by a binary vector of zeros and ones in which each entry in the vector indicates whether the dimension in the cognitive space associated with that entry is ‘relevant’ (entry is 1) or not (entry is 0) in determining whether an object

is a member of a category. For instance, consider a cognitive space with only three dimensions representing height, width, and depth. Then the binary vector (1,1,0) defines a category whose members will be objects for which height and width are relevant dimensions but depth is not (e.g., a piece of paper suitable for the purpose of writing would be a member of the category as it is the height and width of the paper and not its depth that makes it suitable for writing).

The agents in the model can exchange and store information provided by other agents in the form of categories represented by binary vectors. The central idea of our model is that, starting from a limited number of categories, the agents can ‘build’ (or ‘invent’) new categories through social interactions in which information about categories is conveyed from one agent to another agent. In the model, such new categories are represented by new binary vectors, constructed from binary vectors already available to the agents.

We assume that new categories (= binary vectors) are not generated without purpose and that, generally speaking, information is not exchanged just for the sake of doing so. Lane, Malerba, Maxfield, and Orsenigo (1996) have provided a robust argument to the effect that the agents that participate in interactions do so precisely because they expect such interactions to generate emergent constructions from which they will benefit. Here, we assume that the agents go through interactions with other agents because they are interested in producing ‘useful’ categories, which can be defined as those better satisfying their objectives than current ones.

In the model, each agent is characterized by a *social role*. In the current version of the model, only two roles are considered: producer or user². A producer produces physical objects – the artifacts – whereas a user makes use of these artifacts and evaluates their functionality. As the producers develop new cognitive representations in the form of binary vectors, they have two goals in mind: (1) instantiating representations of their categories in the form of artifacts that satisfy the users and (2) fulfilling their production constraints (or objectives).

The goal of the users is to find ways to employ the artifacts satisfactorily in relation to some set of tasks or activities. The artifacts are used in tasks or activities according to their categorization by the user, and their effectiveness in the performance of a task or activity associated with that categorization is evaluated. The objective of the users can thus be restated as ‘finding the categorization of artifacts that entails the highest possible evaluation of a given artifact when it is used in a task or activity befitting its categorization.’ The highest evaluation associated with a categorization and task will be called the ‘functionality score’ for that artifact. In this way, the objectives of the agents are explicitly included in the model along with the subjective nature of the goals.

One of the key points of the model is that the communication process between users and producers involves *incomplete information* transmission. Users cannot read producers’ minds and, therefore, cannot share directly their cognitive represen-

² We foresee implementing a richer variety of social roles, according to specific contexts, in later versions of the model.

tations with them, or vice-versa. Instead, the communication between the two parties must be based on a proper subset of cognitive features associated with categories that both users and producers are able to ‘read.’ This restriction on information transfer is motivated by the assumption that when producers’ and users’ cognitive categories are instantiated in artifacts and in functionalities, respectively, a loss of information occurs. Not all of the dimensions relevant to the producer are expressed in artifacts in a manner that can be ‘read’ by users, and functionalities expressed to the producer need not refer to all of the dimensions relevant to the user. Thus, some of the dimensions of the cognitive representations are not transmitted, and both users and producers will obtain only partial knowledge of each other’s categories.

Finally, the system of beliefs and the physical resources available to all agents play an important role in determining agents’ goals and objectives. These macro-level features influence the relevance that the agents assign to the different cognitive features that define each of the categories. Producers tend to be affected by factors that are mostly related to the implementation of new concepts in artifacts, whereas users are more concerned with the gains that can be obtained by different ways of exploiting the artifacts. Moreover, we assume that the shared system of beliefs affects agents on a slower time scale than ‘quick’ events such as interaction processes with other agents and creation of new categories.

14.4 The Main Parts of the Model

14.4.1 Categories

A category belonging to an agent is represented by a d -dimensional binary vector as $c_i \in \{0, 1\}^d$ ($i \in \{1, \dots, n\}$), where $n\psi$ is the number of categories owned by that agent and $\{0, 1\}^d\psi$ is a d -dimensional binary space. Each element in the vector denotes a different cognitive feature,³ and a specific category is identified by a string of zeros and ones, where each 0 and 1 represents absence or presence of a certain cognitive feature, or dimension.⁴ In our representation, the number of features of a given category characterizes the length of its description (and hence its information content). Categories that have many features (and higher information content) correspond to more complex artifacts, and those are usually the artifacts that require more energy and matter to exist and function in the real world. Whereas a prehistoric stone tool consists of a stone from which certain flakes have been

³ We assume that a category can be potentially composed by a large number of features. Thus, usually d is a very large number compared to the number of features n appearing in the categories owned by the agents.

⁴ The choice of binary vectors is consistent with the representation of the categories argued by Bruner (see, e.g., Bruner, Goodnow, & Austin, 1956). However, we anticipate the possibility of developing a more suitable representation of the cognitive space, for example by allowing each feature to be integer or real-valued variable instead of binary. Some interesting work on the representation of continuous cognitive spaces has been done by Gäardenfors (2000).

removed, in a particular sequence and with a specific technique, a Chinese Ming vase has undergone a much more complex process of manufacture, including the preparation of the clay, the shaping and drying of the object, the decoration and the glazing, and ultimately the firing. That process takes many more of the dimensions of the raw materials and of the technology into account. To bring us up to the modern world, the complexity of manufacture of such a vase is extremely simple compared to that of the manufacture of a car or a computer, which involves many different materials, many different components fulfilling very specific functions, etc. Hence, from the producer's perspective, the category 'Paleolithic stone tool' is precisely defined by the simple description of the material and the flaking technique and sequence; the category Ming vase refers to a description not only of the material and the shaping technique, but also the decoration and the glazing, as well as the processes the object undergoes in the kiln, etc. In the case of the car or the computer the producers' definition of the category refers to the very complex ways in which the many different components have been produced, and have been fitted together in order to create the working object 'car' or 'computer'. A similar argument can be made for the users' descriptions of these objects, which primarily refer to their functions.

14.4.2 Agents

An agent is represented as A_j , ($j \in \{1, \dots, K\}$) and is characterized by its role: it can be either a producer or a user. We denote the role by using the superscripts (p) and (u) (e.g., if the agent A_j is a producer, we write $A_j^{(p)}$). Both types of agents are described in terms of: (i) their cognitive structure (set of categories owned), and (ii) their objectives (target functions). Each producer $A_j^{(p)}$ ($j \in \{1, \dots, k\}$) has a cognitive structure represented by: the number and composition of categories owned, and a real number $s_{ij} \in [0, 1]$, which denotes the functionality score associated with the category c_{ij} . The functionality score is a record of the result of a previous interaction entertained with a user about a given category. It represents a feedback signal that the producer receives from the user concerning user's 'satisfaction' with a certain artifact in relation to the current objectives (more precisely, about the category instantiated by that artifact). If no interaction took place, the score is zero.

The producers' objectives are expressed by a target function $f_j^{(p)}$ that assigns weights to the cognitive features of a given category. The target function takes values on the d -dimensional binary domain and maps onto the interval of real numbers between zero and one ($f_j^{(p)} : \{0, 1\}^d \rightarrow [0, 1]$). In particular, the target function used here is simply defined as the inner product $f_j^{(p)}(c) = \langle c, w_j^{(p)} \rangle$, where $w_j^{(p)}$ is a vector of weights whose elements add to one. Note that if each dimension of the cognitive space is equally important (all the elements of $w_j^{(p)}$ are equal), the features have the same weight. Another extreme case occurs when only one feature is seen as relevant (one of the elements of $w_j^{(p)}$ is one and the rest are zeros).

Similarly, each user $A_j^{(u)}$, ($j \in \{k + 1, \dots, K\}$) is described by the number and composition of the categories c_{ij} owned and by a target function $f_j^{(u)}$ that expresses the importance that users assign to different cognitive features that appear in these categories. Furthermore, each agent is characterized by a set of parameters that control its ‘propensity’ to innovate (create new categories) and to interact with other agents. These parameters define the ‘attitude’ of an agent. Note that the weights in the target functions defined above can be interpreted as attitude parameters as well.

14.4.3 Shared Beliefs

The model intends to take into account cognitive attributes that characterize groups of agents (or the whole collectivity of agents). Although such properties take on their full expression only at the level of the collectivity (macro-level), they have a relevant effect on (and are affected by) the behavior of the individual agents. Such collective factors are represented by a system of beliefs shared by all the members of the collectivity. These shared beliefs influence (and are influenced by) both goals and actions of the agents; in the model they are represented as the ‘average targets’ of a group of users and producers. Namely,

$$F^{(p)}(c) = \left\langle c, \frac{1}{k} \sum_{j=1}^k w_j^{(p)} \right\rangle \quad (14.1)$$

and

$$F^{(u)}(c) = \left\langle c, \frac{1}{K-k} \sum_{j=k+1}^K w_j^{(u)} \right\rangle \quad (14.2)$$

where c is a category and $\frac{1}{k} \sum_{j=1}^k w_j^{(p)}$ and $\frac{1}{K-k} \sum_{j=k+1}^K w_j^{(u)}$ are average weight vectors (i.e., their elements are the average of the elements of the individual agents’ weight vectors).

14.4.4 Environment

The agents in the model can act in a self-sustaining mode and do not necessarily require the explicit intervention of external forces that impact their behavior. However, we allow for the possible inclusion of such external conditions. Temporary or permanent constraints on the whole system could be easily modeled by variables that act on the shared system of beliefs (or on agents’ targets). Such constraints could be interpreted as external groups of agents or other conditions not specifically included in the model.

14.5 Dynamics

14.5.1 Micro-Level: Generation of New Categories

In our model, the agents construct new categories as a consequence of two processes: social interaction and invention. The latter is implemented in the model by a mechanism that allows for creation of new categories through blending cognitive features of existing categories, following recent work by Fauconnier and Turner (2002). They formulated a cognitive theory of concept integration, named conceptual blending.⁵ Although the general idea of conceptual blending has been very popular in cognitive sciences, so far there is no general agreement on the formal definition of the blending process itself. Some interesting work on this topic is provided Veale (1997) and Falkenhainer (Falkenhainer, Kenneth, & Gentner, 1989). Nevertheless, many authors agree on the fact that such a mechanism produces a new cognitive object (category) encompassing at least the following requirements:

- a. a new category is built starting from one or more available categories;
- b. the new category should keep most of the features of the original categories;
- c. the new category must be distinguishable from the original categories; and
- d. the selection process of the categories to blend is guided by the intentions of the subject performing the blending.

Since the details of the dynamics of the cognitive processes regulating the creation of new categories are beyond the goals of this work, a stochastic-type representation of the combination mechanism seems to be appropriate for our purposes. A genetic algorithm appears to be a reasonable first-approximation to the formation of new categories that satisfies the above requirements with the level of detail in which we are interested. There are two basic types of operators that are classically employed in genetic algorithm literature: crossover and mutation. We formulate a straightforward application of crossover and mutation to our representation of categories as binary vectors as follows. The crossover operator takes two categories, $c = (c_1, \dots, c_n)$ and $v = (v_1, \dots, v_n)$, and produces two new categories in the following manner. A number r is picked at random from the set of integers $\{1, \dots, n\}$ and two new categories \bar{c} and \bar{v} are created from c and v according to the following equations:

⁵ This theory is based on basic ideas advanced by George Lakoff in his book *Women, Fire and Dangerous Things* (1987). Although there is no cognitive theory that has yet been able to cover even a significant fraction of the phenomena of human cognition, some claim that conceptual blending is rising in prominence among such theories. Blending is generally described as involving two cognitive objects (e.g., categories) that, according to a given structural mapping, will generate a third cognitive object, called Blend. This new cognitive object is generated by a cognitive process called 'selective projection', which maintains the partial structure from the original cognitive objects and adds emergent structure of its own.

$$\bar{c}_i = \begin{cases} c_i & \text{if } i < r \\ v_i & \text{otherwise} \end{cases} \quad (14.3)$$

$$\bar{v}_i = \begin{cases} v_i & \text{if } i < r \\ c_i & \text{otherwise} \end{cases} \quad (14.4)$$

Next, one of the two resulting categories is randomly discarded, say, \bar{v} . Finally, mutation is then applied to \bar{c} in the form of binary mutation. Binary mutation flips at random one or more elements of \bar{c} to its other value.

$$\bar{c}_i = \begin{cases} 1 - c_i & \text{if element } i \text{ is picked} \\ c'_i & \text{otherwise} \end{cases} \quad (14.5)$$

Cross-over and mutation satisfy requirements (a), (b) and (c). In the next sections, we discuss requirement (d) more specifically from the user's and producer's points of view.

14.5.2 Meso and Macro Level: Social Interactions

14.5.2.1 Producers' Actions

The producers' behavior is characterized by two basic types of action: (i) creating new categories (invention) and (ii) instantiating categories in the form of artifacts that will be used by one or more users (production).⁶

Invention. The producer chooses two categories based on the evaluation of the target function, the functionality score previously recorded from the user and the system of beliefs to which the producer belongs. The categories are chosen according to a roulette wheel mechanism (Holland, 1975). The probability that producer j picks the category c_i is defined to be:

$$P_j(c_i) = \frac{\alpha^{(p)} f_j^{(p)}(c_i) + \beta^{(p)} F^{(p)}(c_i) + \gamma^{(p)} s_i}{\sum_{k=1}^n \alpha^{(p)} f_j^{(p)}(c_k) + \beta^{(p)} F^{(p)}(c_k) + \gamma^{(p)} s_k} \quad (14.6)$$

where $f^{(p)}$ is the producer's target function, s_i is the functionality score associated with the category c_i (to be described in more detail in the next section when discussing the user's actions), $F^{(p)}$ is the collective target function for producers and n is the number of categories available to producers at the time of the choice. In addition $\alpha^{(p)}$, $\beta^{(p)}$, and $\gamma^{(p)}$ are positive constants that control the relative influence of $f^{(p)}$, $F^{(p)}$ and s_i on the choice of the category. In particular, we assume that the shared system of beliefs affects agents on a slower time scale than 'quick' events

⁶ See Read 19ss, 200× for a discussion of instantiation as a process for translating between the ideational and the phenomenological domains.

such as interaction processes with other agents and creation of new categories. Therefore, the parameter β is usually set to a small value in comparison to α and γ .

Next, the crossover and binary mutation operators described in the previous section are applied to the two selected categories, and a new category \bar{c} is obtained. The likelihood for the crossover and mutation operations to be applied is controlled by a parameter $m^{(p)} \geq 0$ (producers' innovation propensity). Finally, the target function is evaluated at \bar{c} and the weights of the producer fitness function are updated as follows,

$$w_j^{(p)} \leftarrow (1 - q)w_j^{(p)} + q \frac{\bar{c}}{\sum_i \bar{c}_i}, \quad (14.7)$$

and then normalized. Here $0 \leq q \leq 1$ is a parameter controlling the degree of influence of the new category \bar{c} (normalized) on the fitness function. Usually, we set q to a small and fixed value, e.g., $q = 0.05$; but it seems more realistic to assume that larger values of the target function correspond to more significant shifts.⁷

Production. When a new artifact is produced, one or more existing categories are 'packaged' in a real object/artifact that is then passed to the user. However, the artifact is able to convey the user only some of the cognitive features that characterize the original category \bar{c} , for which it is an instantiation. Therefore, production is treated as involving partial loss of information from the producer's category due to the passage from the cognitive domain of categories to the physical domain of matter and energy.

More specifically, the user acquires only a subset of the elements from the original vector \bar{c} , while the rest of the elements are suppressed. We call such a partial representation an instance, and the product from the category \bar{c} is denoted with \bar{c}^* (e.g., if the original category is $\bar{c} = (c_1, \dots, c_i, c_{i+1}, \dots, c_n)$, after the production, the user gets to know $\bar{c}^* = (c_{i+1}, \dots, c_n)$ but not (c_1, \dots, c_i)). This information reduction step will be referred in the proceeding of the discussion as *instantiation*.

14.5.2.2 Users' Actions

The users' behavior is summarized by the following activities: (i) exploiting the artifacts/judging their functionality (utilization), and (ii) creating new ways to utilize the products artifacts (invention of use).

Utilization. After receiving an instance \bar{c}^* from the producer, the user selects a category \tilde{c} from those available, while optimizing two criteria. The first one is that the cognitive features of the category to be chosen must be as close as possible to the corresponding cognitive features of the incoming instance \bar{c}^* . The distance between the two sets of cognitive features \bar{c}^* , \tilde{c}^* is computed using a metric d (here defined

⁷ We could model q as a function of $f_j^{(p)}(\bar{c})$. For example, $q(\bar{c}) = ae^{-b_j^{(p)(\bar{c})}}$ for some $a, b \in \mathbb{R}^+$.

to be the Manhattan distance⁸). The second criterion is that the chosen category must correspond to a large value of $f^{(u)}(\tilde{c})$. Next, provided that a suitable category \tilde{c} is selected, the value of its functionality score, $S_{\tilde{c}} = f^{(u)}(\tilde{c})$, is sent back to the producer, who associates it with the original category \tilde{c} .

Invention. In our model, the producer is not the only kind of agent that can generate new categories. From time to time, the user attempts to create a new category before evaluating the functionality score of a product received from the producer.⁹ For simplicity, we describe this mechanism by using the same genetic operators as we used for the producer: simple crossover and binary mutation. The selection function, however, takes into account the closeness of the categories to be selected with respect to the given instance. More precisely, given an instance c^* , the probability that the user chooses the category c_i is

$$P_j(c_i) = \frac{\alpha^{(u)} f_j^{(u)}(c_i) + \beta^{(u)} F^{(u)}(c_i) + \gamma^{(u)} d(c_i^*, c^*)}{\sum_{k=1}^n \alpha^{(u)} f_j^{(u)}(c_k) + \beta^{(u)} F^{(u)}(c_k) + \gamma^{(u)} d(c_i^*, c^*)} \quad (14.8)$$

where $f^{(u)}$ is the user's fitness function, $F^{(u)}$ is the collective fitness function for users, and $d(c_i^*, c^*)$ is the distance between c_i^* and c_i . In the formula above, $\alpha^{(u)}$, $\beta^{(u)}$ and $\gamma^{(u)}$ are positive constants that allow one to tune the relative influence of $f^{(u)}$, $F^{(u)}$ and d in the choice of the category. The user's innovation propensity is regulated by a parameter $m^{(u)} \geq 0$, which controls the occurrence of cross-overs and mutations.

Finally, the user evaluates its target function at \tilde{c} , and, subsequently, the weights of the producer fitness function are updated similarly to the producer's case, following the updating rule:

$$w_j^{(u)} \leftarrow (1 - q)w_j^{(u)} + q \frac{\tilde{c}}{\sum_i \tilde{c}_i} \quad (14.9)$$

where $0 \leq q \leq 1$ controls the influence of the new category \tilde{c} on the user's target function.

14.6 The Simulations

In order to test the behavior of the model, we coded a preliminary implementation in the Matlab environment. In this section, we provide a general description of the implementation and discuss some interesting features of the model that emerged from numerical simulations. In the simulation, three main tasks are performed:

⁸ The Manhattan (or city-block) distance between two vectors x and y is defined as $d(x, y) = \sum |x_i - y_i|$.

⁹ This has to be interpreted as an attempt to create a new way to find a more beneficial use for a given artifact (a new use that fits better to the user's goals).

(i) setting the initial conditions (initial set of agents, the collection of categories owned by each agent), (ii) performing the actual interaction among agents (simulation) and (iii) presenting graphical and numerical summaries of the simulations. Although we could have introduced a larger number of parameters¹⁰ and modifications to enrich the behavior of the agents, in this first version we adopt some simplifications. Namely,

- the number of dimensions of the cognitive space (number of cognitive features) is the same for each agent;
- the number of agents is fixed (no new agents are born or die after the interactions begin);
- the number of categories for each agent at each given time step is fixed (unexploited categories are ‘forgotten’);
- the number of producers equals the number of users;
- at each time step each agent can create at most one new category; and
- each agent interacts with every other agent at a given time step.

We set the initial conditions for the simulations as follows: the number of features in the cognitive space is $d = 1,000$, the number of categories owned at each time step is $n = 10$ (as explained in the previous sections we assume that the number of categories is much lower than the number of dimensions in the cognitive space). We chose initial target functions with weights concentrated around a limited number of cognitive features; in particular the initial weights were assigned at random using the following binomial expressions:¹¹

$$w_i^{(p)} = \binom{d}{i} (\Pi^{(p)})^i (1 - \Pi^{(p)})^{d-i}, \quad (14.10)$$

$$w_i^{(u)} = \binom{d}{i} (\Pi^{(u)})^i (1 - \Pi^{(u)})^{d-i}, \quad i = 1, \dots, d, \text{ and} \quad (14.11)$$

$$(\Pi^{(p)}, \Pi^{(u)}) \in (0, 1) \quad (14.12)$$

In the following experimental settings, only a given fraction of the cognitive features of a category can be communicated to the user by the producer via instantiation; for simplicity, we set the fraction to one half and we keep it fixed for each time step. Furthermore, the fitness level of an agent is measured by the ‘best’ category available at a given time (highest value of the target function evaluated for the categories available at a given time).

¹⁰ It can be argued that overparametrization of a model can affect significantly its validity. In fact, introducing many parameters that need to be fitted from data and/or threshold parameters creates room for interpretative uncertainty. At the current stage of our work, it seems more reasonable to focus our analysis on a limited number of parameters characterized by straightforward interpretation. In particular, the value of a parameter can be set and interpreted as a proportion of the other parameters included in the model and for which the effect of the parameter on model outcome can be evaluated.

¹¹ Other choices for the initialization of the weights of the target functions are possible.

In the next sections, different experimental settings are studied. First, for illustrative purposes, we consider the simple case when only one producer and one user interact in the absence of a shared system of beliefs. Next, we study the behavior of the model in a multi-agent setting.

14.6.1 *Single User and Single Producer*

In this section, we examine the case when only two agents (a producer and a user) interact on a recurrent basis. In this example, the two agents act in a secluded system and thus their shared beliefs coincide with individual goals. Therefore, producer and user pick categories to combine and exchange exclusively based on their individual objectives and on feedback signals from the other in response to their actions.

The user and the producer begin their reciprocal interaction with quite different cognitive equipment; in fact, we set both target functions and sets of categories to be dissimilar and sparse at the initial stage of the simulation. Moreover, their actions are limited by incomplete information exchange as only a fraction of the cognitive features of a category in the producer's mind is transmitted to the user (instantiating the category results in information loss as discussed above).

Nevertheless, by continuous (and secluded) interaction we observe a convergence process described by two peculiar qualities. The first concerns the goals of producer and user; as the interactions continue, we observe a progressive convergence of the objective functions.¹² The weights that producer and user assign to their cognitive dimensions gradually 'align.' This can be understood as 'learning by doing,' under conditions of partial information. The second quality is related to the composition of the emerging categories. The categories that emerge after a sufficiently large number of time steps are more 'beneficial' than the older ones, and the emergent categories tend to fit better and better with the agents' respective objectives. This means that the agents, not only learn from each other by interacting, but also achieve increasingly rewarding benefits from the result of the learning process. It is important to note, however, that in order to make their objectives converge, the simulations showed that increasingly more cognitive features must be progressively involved. As a result, the categories emerging at the end of the process contain a larger number of features than the ones with which the agents begin.

Beneficial categories and alignment of agents' intentions (targets) emerge through micro-level fluctuations (selecting and combining categories) that are strengthened by positive feedbacks and stabilized by negative feedbacks. The amplification of the micro-level fluctuations is controlled in various ways by the value of the 'innovation propensity' parameters α , m and q of an agent. In particular, when $\alpha^{(p)} = \alpha^{(u)}$ and $m^{(p)} = m^{(u)}$, we observe that different values of q (the parameter

¹² The convergence of objectives is measured by the distance between the vectors of weights of the target functions.

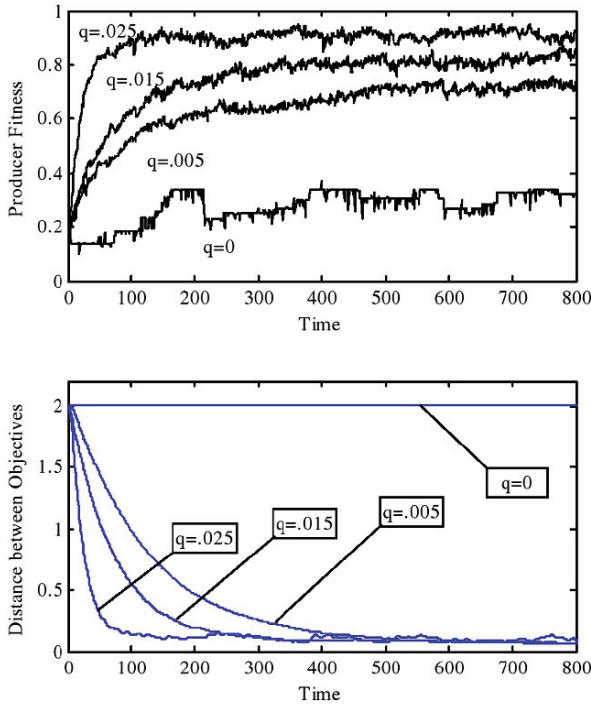


Fig. 14.1 Single user-single producer system. Producer's fitness and distance between producer and user's objectives for choices of q . Producer's fitness is measured in terms of the largest value of the target function, evaluated for all available categories (or 'best available category'). The distance between user's and producer's categories is computed using Euclidean distance between the two 'best' categories

q models the agents' inclination to modify their intentions) affect significantly the processes of convergence. Figure 14.1 illustrates examples of the convergence process for choices of $q = \{0, 0.005, 0.015 \text{ and } 0.025\}$; for this simulation we set $\alpha^{(p)} = \alpha^{(u)} = 0.5$, $\gamma^{(p)} = \gamma^{(u)} = 0.5$ and $m^{(p)} = m^{(u)} = 1$. Larger values of q are associated with faster convergence dynamics. This means that the agents can profit from slightly adjusting their objective functions when a newly created category turns out to be more valuable than the older ones.

Conversely, if their intentions are unalterable (i.e. their objective functions do not change), the gain derived from the generation of new categories occurs only shortly. In fact, if the parameter q is set at its minimum value ($q = 0$), then the alignment is only shortly experienced. As a result, the convergence process between producer and user is unstable and there is no consistent increase either in the convergence of objectives or the value of the newly generated categories.

It is also important to notice that the level of q plays a role in the final level of fitness at which user and producer 'lock in.' The faster the convergence process, the higher the potential fitness level of the two agents involved in the interaction.

On the other hand, the strength of the feedback dynamics (leading to amplification or stabilization of the feedback loops) between user and producer is controlled by the parameter γ . This turns out to be a key parameter in the model; in particular the simulations show that small values of the parameter γ (small, relatively to α , i.e. $\gamma = 0.1\alpha$) generate unstable paths where short periods of increased fitness are followed by descents.

14.6.2 Multiple Users and Producers

In this section, we are concerned with the behavior of the model when more than one user and one producer interact on a recurrent basis. Besides the features already discussed for the case of a single user and producer, this simulation introduces the shared beliefs.

The multi-agent simulations performed so far suggest some interesting qualitative properties of the model emerging at the meso- (agents) or macro- (collectivity) level. In particular, these concern three factors: the density of the population of agents, the innovation propensity of the agents, and the relevance of shared beliefs during the creation of new categories.

14.6.2.1 Agents' Density and Innovation Propensity

A series of simulations was performed to study the effect of the number of agents on convergence (or divergence) dynamics under conditions of *sufficiently high* (and homogeneous) adaptability of objectives for all the agents. We set $q_k = q > q^*$ for all $k = 1, \dots, K$ and the other parameters as $m = m_k$, $\alpha = \alpha_k$, $\beta = \beta_k$ and $\gamma = \gamma_k$ for all $k = 1, \dots, K$. For instance, Fig. 14.2 illustrates the path of two particular producers in the model. The plot displays two sets of curves, corresponding to two choices of the number of agents ($K = 4$ and $K = 8$).

When more agents act in the model, both the increase in agents' fitness and the alignment of their objectives tend to happen faster. Nevertheless, the condition required for such behavior is that the innovation propensity q and m must be large enough (in the case shown in Fig. 14.2, we used $q = 0.05$ and $m = 1$).

Next, we consider the case when the population of agents is heterogeneous with respect to their innovation propensity. The population is composed of two types of agents, characterized by lower or higher innovation propensity respectively. In that situation, the simulations show two recurrent qualities. First, as expected, the agents with lower innovation propensity exhibit irregular convergence paths where phases of alignment with other agents (and increasing fitness) are followed by declines. Usually, after an initial positive jump at the beginning of the process, they reach neither a high level of fitness nor a good alignment of objectives. Conversely, the more innovative agents show more stable growth and alignment paths and tend to reach higher levels of fitness than do the agents in the other group. In Fig. 14.3, we present a simulation with $K = 8$ agents, where four of them (2 producers and 2 users) are more innovative than the rest. The innovation propensity parameters for

Fig. 14.2 Producer’s fitness in two systems with $K = 4$ and $K = 8$ agents

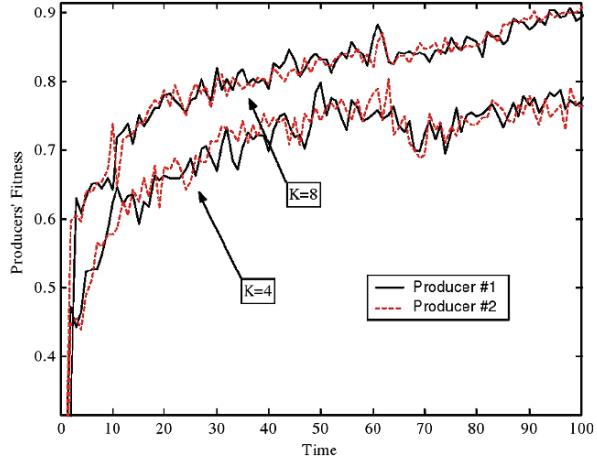
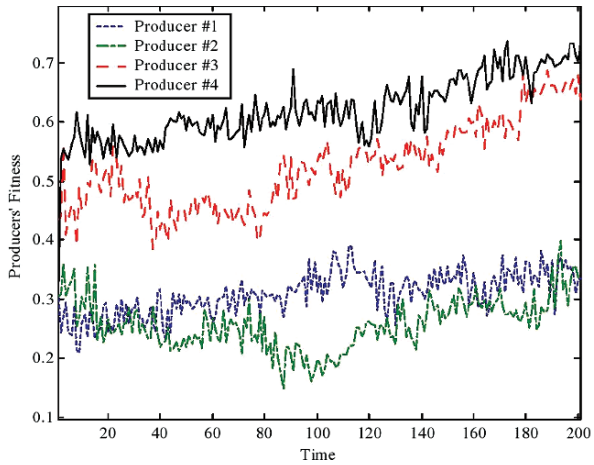


Fig. 14.3 Producer’s fitness and innovation propensity. For producers #3 and #4 we set higher innovation propensity ($q = 0.05$ and $m = 1$) than for producers #1 and #2 ($q = 0.01$ and $m = 0.5$)



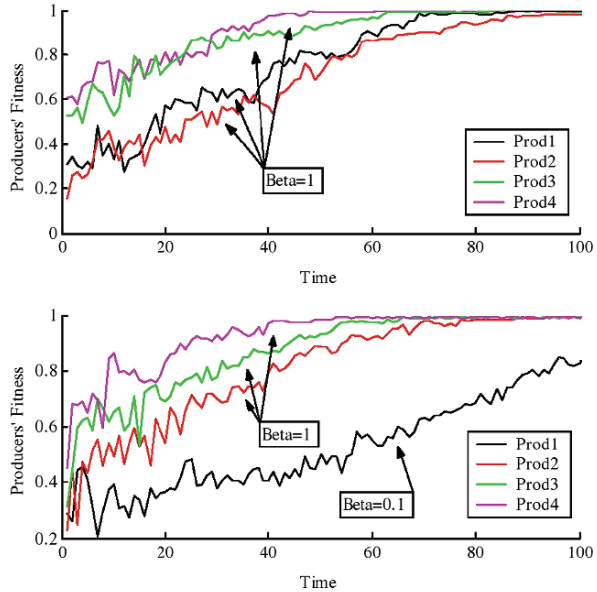
the innovative group are $q = 0.05$ and $m = 1$, while for the others we set $q = 0.01$ and $m = 0.5$.

14.6.2.2 The Role of Shared Beliefs

Finally, we produced a series of simulations to survey the impact of the shared beliefs on the individual agents. In a situation where each agent possesses similar innovation characteristics, we consider a subset of agents who were assigned a smaller weight to their shared beliefs during the process of creating new categories. In general, the shared beliefs function as an adhesive for the group of agents involved in the interaction process: they ease the alignment among agents of the group and favor the attainment of higher fitness level for each member.

The plots in Fig. 14.4 illustrate the typical behavior encountered in the course of our experiments. The two plots display simulations with $K = 8$ agents (4 producers

Fig. 14.4 Fitness for producers belonging to groups with different shared beliefs: for producers #1 and #2 we set $\beta = 1$, while for producers #3 and #4 we set $\beta = 0.1$



and 4 users). The initial conditions and the parameters are identical except for β . While in the first simulation all the agents assign the same weight to the shared beliefs ($\beta = \beta_k = 1$ for all k), in the second case, one of them (producer #1) is distinguished by a smaller weight ($\beta_1 = 0.1$). The subsets of agents that ‘disregard’ the shared beliefs would need higher innovation propensity in order to keep up with the rest of the agents.

14.7 Conclusions

In this paper, we have discussed a new way to model information exchange in societies, which takes into account innovation at the level of the individual agent and the circulation of newly generated information at the group level. In particular, the model provides a connection between cognitive phenomena occurring at different levels of social aggregation and the concomitant time scales: individual agents (micro-level), local interaction among agents (meso-level), and collectivity (macro-level).

One of the distinctive features of this model is the differential representation of information according to the roles of the agents. The information belonging to the agents is modeled in terms of sets of categories and the individuals can create new categories starting from existing ones in response to information obtained from other agents through social interaction processes. Individuals act in accordance with their social role (as either producer or user of artifacts), and the information is exchanged in terms of ‘packages of categories’ instantiated by the artifacts that are produced and circulated in the society. It is important to note that the agents do

not generate novelty randomly. In the model, a specific set of factors guides the agents in creating new categories according to the individual agent's objectives. The intentions of the agents are explicitly modeled as goals (i.e. features in their respective cognitive spaces that are regarded as important) that can evolve in response to interactions with other agents. As the dynamics of the agents' interactions proceeds, the exchange of information combined with the search for new and more rewarding categories (better satisfying agents' goals) may lead to alignment processes among certain groups of agents. The alignment can be measured in terms of similarity of the objectives owned by the agents.

The results of the simulations presented in Section 14.5 indicate that the degree of adaptability of agents' objectives to new and more 'valuable' categories affects (1) the speed of the alignment processes among agents and (2) the potential benefits that may be derived from further interaction. In particular, the parameter that tunes the strength of the feedback dynamics between user and producer (e.g., user's response about the functionality of a certain artifact) turns out to be crucial. For small values of this parameter, the simulations exhibit unstable paths where short periods of cognitive alignment of the relevant agents are then followed by divergences.

Other factors playing a significant role in both the generation of new categories and in the alignment processes are the number of the population of agents and the influence of the system of beliefs on the convergence process. Most simulations conducted under conditions of homogeneous innovation propensity showed that when sufficiently large subgroups of individuals assign bigger weights to the shared beliefs, the alignment processes happen faster and can potentially reach a higher level of fitness than in the case of other subgroups of agents. In this sense, shared beliefs seem to ease the occurrence and the speed of the convergence processes. Advantages in being disconnected from the shared system of beliefs can be enjoyed only by sufficiently innovative agents.

Substantial improvements to the current model can be obtained by working on two aspects: (i) a more accurate representation of the cognitive processes leading to creation of new categories at the level of the individual, and (ii) a richer description of the pattern of social interactions (depending on the purpose of the analysis at hand). As far as concerns (i), although models for category innovation such as conceptual blending have become increasingly popular, at this stage they appear to be too general to provide insight into how we model the generation of new categories in contexts leading to innovation. On the other hand, aspect (ii) has been considerably simplified in our preliminary version of the model and future improvements are certainly possible. In the simulations performed so far, each agent was allowed to interact with all other agents at any time step; however, in the real world, a given agent usually interacts with a subset of particular agents according to a more complex set of conditions. Lane, Malerba, Maxfield and Orsenigo (1996) identify factors necessary for beneficial interaction paths among agents (generative relationships) such as common intentions, heterogeneity of competencies and attributions, and mutual directedness in feeding a recurrent path of interactions. Including these and other conditions in the model may allow for the study of a richer variety of real situations in different social contexts.

References

- Anderies, J. (2006). Robustness, institutions, and large-scale change in social-ecological systems: The Hohokam of the Phoenix basin. *Journal of Institutional Economics*, 2, 133–155.
- Ballot, G., & Weisbuch, G. (Eds.). (2000). *Applications of simulation to social sciences*. Paris, France: Hermes Science Publications.
- Böeda, E. (1986). Approche technologique du concept Levallois et évaluation de son champs d'application: étude de trois gisements saaliens et weichseliens de la France septentrionale. Unpublished Ph.D. Thesis, University of Paris X, France.
- Böeda, E. (1990). De la surface au volume : analyse des conceptions des débitages Levallois et laminaires. In C. Farizy (Ed.), *Paléolithique moyen récent et Paléolithique supérieur ancien en Europe. Ruptures et transitions : examen critique des documents archéologiques -Actes du Colloque international de Nemours, 1988* (pp. 63–68). Nemours: Mémoires du Musée de Préhistoire d'Ile de France.
- Bruner, J. S., Goodnow, J. J., & Austin, G. A. (1956). *A study of thinking*. New York, NY: Wiley.
- Falkenhainer, B., Kenneth, F., & Gentner, D. (1989). The structure-mapping engine: Algorithms and Examples. *Artificial Intelligence*, 41, 1–63.
- Fauconnier, G., & Turner, M. (2002). *The way we think: Conceptual blending and the mind's hidden complexities*. New York, NY: Basic Books.
- Foster, G. W. (1948). *Empire's children. The people of Tzintzuntzan*. Washington, DC: Smithsonian Institution, (Institute of Anthropology Publication).
- Gärdenfors, P. (2000). *Conceptual spaces: The geometry of thought*. Cambridge, MA: MIT Press.
- Gilbert, N. (1991). *Artificial Societies*. Guildford, UK: University of Surrey.
- Holland, J. (1975). *Adaptation in natural and artificial systems*. Ann Arbor, MI: University of Michigan Press.
- Janssen, M. A., & Jager, W. (2000). The human actor in ecological-economic models. *Ecological Economics* (Special Issue), 35, 3.
- Janssen, M.A. (2002). *Complexity and ecosystem management: the theory and practice of multi-agent systems*. Cheltenham UK/Northampton, MA: Edward Elgar Publishers.
- Kohler, T., & Gumerman, G. (Eds.). (2000). *Dynamics in human and primate societies: Agent-based modeling of social and spatial processes*. New York: Santa Fe Institute and Oxford University Press.
- Kohler, T., & van der Leeuw, S. E. (Eds.). (2007). *The model-based archaeology of socionatural systems*. Santa Fe, NM: School of American Research.
- Lakoff, G. (1987). *Women, fire and dangerous things*. Chicago, IL: University of Chicago Press.
- Lane, D.A., Malerba, F., Maxfield, R., & Orsenigo, L. (1996). Choice and action. *Journal of Evolutionary Economics*, 6, 43–76.
- Macrae, N., & Bodenhausen, G. (2000). Social cognition: Thinking categorically about others. *Annual Review of Psychology*, 51, 93–120.
- Markman, A. B., & Gentner, D. (2001). Thinking. *Annual Review of Psychology*, 52, 223–247.
- Papousek, D. A. (1989). Technological change as social rebellion. In: S. E. van der Leeuw & R. Torrence (Eds.), *What's new: A closer look at the process of innovation* (pp. 140–166). London, UK: Unwin and Hyman.
- Pigeot, N. (1992). Reflexions sur l'histoire technique de l'Homme: de l'évolution cognitive à l'évolution culturelle. *Paléo*, 3, 167–200.
- Selby, H. A., & El Guindi, F. (1976). Dialectics in Zapotec thinking. In K. Basso & H. A. Selby (Eds.), *Meaning in anthropology* (pp. 181–196). Albuquerque, NM: University of New Mexico Press.
- Veale, T. (1997). Creativity as pastiche: A computational model of dynamic blending and textual collage, with special reference to the use of blending in cinematic narratives. Proceedings of ICLC'97, the 1997 conference of The International Cognitive Linguistics Association, Amsterdam, The Netherlands.

- van der Leeuw, S. E. (1976). *Studies in the technology of ancient pottery* (Vols. 2). Amsterdam, The Netherlands: University Printing Office.
- van der Leeuw, S. E. (1990). Archaeology, material culture and innovation. *SubStance* 62–63, 92–109.
- van der Leeuw, S. E. (1993). Giving the potter a choice: conceptual aspects of pottery techniques. In P. Lemonnier (Ed.), *Technological choices: Transformation in material culture from the neolithic to modern high tech* (pp. 238–288). London, UK: Routledge Kegan Paul.
- van der Leeuw, S. E. (2000). Making tools from stone and clay. In T. Murray, & A. Anderson (Eds.), *Australian archaeologist. Collected Papers in Honour of J. Allen* (pp. 69–88). Canberra, Australia: ANU Press.
- van der Leeuw, S. E. (2005). Why model? *Cybernetics and Systems* 35(2–3), 117–128.
- van der Leeuw, S. E., & McGlade, J. (Eds.). (1997). *Time, process and structural transformations in archaeology*. London, UK: Routledge.
- van der Leeuw, S. E., & Papousek, D. A. (1992). Tradition and innovation. In F. Audouze, A. Gallay, & V. Roux (Eds.), *Ethnoarchéologie: Justification, problèmes, limites* (pp. 135–158). Antibes, France: A.P.D.C.A.
- van der Leeuw, S. E., Papousek, D. A., & Coudart, A. (1992). Technological traditions and unquestioned assumptions. *Techniques et Culture* 16–17, 145–173.

Chapter 15

Exaptive Processes: An Agent Based Model

Marco Villani, Stefano Bonacini, Davide Ferrari and Roberto Serra

15.1 Introduction

This chapter introduces an agent-based model designed to investigate the dynamics of some aspects of exaptation that have been discussed previously in this volume. It is strongly related to the model introduced in the previous chapter. Indeed, in the model described here, cognitive categories represent the main tools that the producers and users of artifacts employ in order to interpret their environment, as in the case discussed in Chapter 14. The main addition provided by the current model, however, is the explicit introduction of artifacts.

As stressed in Chapter 1, artifacts are a key component of human organizations and activities. Artifacts are entities constructed by an organization to enhance its or other organizations' functionalities functionality. One of their main properties of interest to us is their capability to convey information, although they may not be explicitly designed for this purpose. In addition, there are artifacts specifically designed to store and carry information, like e.g. books, radios, televisions, including the very special kind of artifact represented by computers, which are able to process information at a very high level of abstraction.

Since artifacts convey information, their explicit representation eases the understanding of the exaptation phenomenon, seen in this context as a shift in terms of "leading attributions." Actually, their introduction is important in order to characterize the ontology necessary to identify exaptation events. Here we focus on phenomena occurring at the micro-level (how individuals collect information about the external world, categorize it, and combine existing categories in order to create new ones) and meso-level (the exchange of information among individuals). However, we do not explicitly include the details concerning the macro-level events (the shared system of beliefs and the common physical and technological resources); which are left for further research.

M. Villani (✉)
Department of Social, Cognitive and Quantitative Sciences, University of Modena
and Reggio Emilia, Modena and Reggio Emilia, Italy
e-mail: mvillani@unimore.it

In the first two sections, we introduce the notion of exaptation. The third and fourth sections describe the model that we developed in order to explore some aspects of exaptation and its dynamics. Finally, we discuss the results of our first simulations and identify some elements able to favor the emergence of exaptation phenomena.

15.2 Exaptation

Recently, the concept of exaptation has been introduced to explain the changes resulting from innovation processes and the rise of new technologies. Exaptation originates from the domain of biology, where it appeared in Gould and Verba (1982) who referred to species evolution as the mechanism complementary to Darwinian adaptation. The following definition (Ceruti, 1995) gives insight into the main idea of exaptation: “. . .the processes whereby an organ, a part, a characteristic (behavioral, morphologic, biochemical) of an organism, which was originally developed for a certain task, is employed for carrying out tasks that are completely different from the original one.” The typical example provided by Gould (2002) is represented by a line of feathered dinosaurs, arboreal or runners who developed the capability to take advantage of feathers for flying, when originally they were adapted for thermoregulation.

Furthermore, exaptation can provide a key to interpret the serendipity that characterizes the generation of new products. Exaptation emphasizes that the functionalities for which a technology has been developed are only a subset of the consequences generated by its introduction. In many cases, there can be several different consequences generated by a new technology, a product, or a process and thus its exaptive potential can be very large. Hence, exaptation is to be interpreted as a central idea connecting technological progress and emergence of recurrent patterns of interaction.

Mokyr (1998) states that exaptation “refers to cases in which an entity was selected for one trait, but eventually ended up carrying out a related but different function.” Such a definition captures the idea that exaptations are those characteristics of a certain technology that are recruited for a purpose different from the original one (a process that may lead in turn to a cascade of further changes). Different from *adaptations*, which present functions for which they are selected, exaptations generate effects that are not subject to pressures from the current selections, but potentially relevant later on.

A classical example of technical innovation illustrating both adaptation and exaptation is the Compact Disk (CD). The CD was originally developed in 1960 in the Pacific Northwest National Laboratory in Richland, WA and it was designed for a specific task: to solve the problem of the sound quality deterioration of the classical vinyl records. Its inventor, J.T. Russell, developed a system based on the idea of using light to carry information, avoiding the usual contact with mechanical parts of the recording device. The CD-ROM was patented in 1970 as a digital-optic system for recording and reproducing sound. Later, researchers exapted the technology of

the CD-ROM for a different purpose: storage media for computer data. Although the latter represented a function not originally intended for the CD-ROM, it became clear that it was indeed effective. As a result, during the 1970's, the Laboratory refined the CD-ROM technology, selling a product that could be usefully employed for different purposes and improving some of its characteristics (increase of memory capacity, recording speed, sound quality).

Another important aspect of such a phenomenon is represented by what Gould defines as the "exaptive pool." The exaptive pool represents the potential allowed for future selection episodes (at all levels). There are two categories of potential: (i) intrinsic potential and (ii) real instantiations.

To understand intrinsic potential, let us recall that the smallest and lightest among the US coins, the dime, despite its small purchasing value, is still in use. An unforeseen consequence of its technical characteristics is that the dime can be employed as an occasional screwdriver. Hence, the dime is an adaptation if considered as money and exaptation if considered as a screwdriver. Note that the potential for the additional functionality as a screwdriver is an intrinsic characteristic of the dimensions and the shape of the coin and, thus, cannot be considered as disjoint or separated from it.

The second category of the exaptive pool includes real entities, material things that become part of the item under consideration, due to various reasons. They are not currently associated with a particular use and at the same time do not generate substantial damages, thus avoiding elimination by selection. The members of this category can be generated in various ways as structures that previously were considered useful, or as neutral characteristics introduced "incognito" with respect to the selection process.

15.3 The Origin and Peculiar Features of Exaptation

To complete a picture of exaptation as a phenomenon, we consider two questions:

1. What are the factors that give rise to exaptation?
2. What are the traits that distinguish it from other processes?

In order to answer the first question, we refer to the idea of decomposability introduced by Simon (1969) to study the architecture of complex systems. An artifact is seen as a hierarchical structure composed of subparts that are approximately independent in the short term, but connected by a global behavior in the long term. Therefore, the subparts are selected, and they proliferate as a consequence of being only one among many aspects of the whole artifact. While some such subparts can have an important role with respect to the current functionality of the whole artifact, others remain latent awaiting future activation.

Based on the above observations, we can divide the possible origins of an exaptation into three groups: (1) the case where a subpart is already providing a positive contribution to the functionality for which the technology was designed. Only in a

later moment and after changing the context, the subpart becomes the main component. (2) The case where the subpart has no role in the overall performance of the system. (3) Finally, the case where the subpart provides a negative contribution to the overall performance of the system.

An example of case (1) is represented by the phonograph invented by Edison in 1877. The innovative technology for amplifying and reproducing sound suggested the commercialization of the invention as an office dictaphone. Only a dramatic change of the context led to the exaptation of the phonograph, later re-named a gramophone, and nowadays considered one of the most popular inventions made during the 19th century. Case (2) is well identifiable in the process of waste vitrification, originally developed to reduce environmental pollution (safe waste disposition and limited radioactive waste), which is now mainly for biological hazards, e.g. the elimination of biochemical weapons. Finally, case, (3) can be exemplified by the innovation in plastic production during the Second World War, based on subproducts (formerly not used) deriving from oil refineries (Dew, Sarasvathy, & Venkataraman, 2004).

Let us now address the question of which traits distinguish exaptation from other processes. An exaptation is usually identifiable with respect to a change of context, which generates a change in the utility of the technology, and not simply to a new mixture of existing elements, as for example in Schumpeter's processes (Schumpeter, 1934, 1976). One could think that an exaptation is simply an unintended consequence of technology. However, this statement ignores the fact that the act of exapting requires the intentional activation of a technology that otherwise would remain latent.

15.4 A Model of Exaptation

The idea that exaptation plays a key role in innovations is convincing and appealing, and can be supported by case studies, like those briefly mentioned in the previous sections. Yet the concept has some subtleties, and in order to explore the behavior of a system which can undergo exaptation it is interesting to take a modeling approach. Key questions to be addressed include the very possibility to develop a model that shows exaptation, the growth dynamics of artifact space generated in this way, and its possible unbounded expansion.

The aim of the EMIS system (*Exaptation Model in Innovation Studies*) is to highlight the factors that contribute to the occurrences of exaptation. EMIS is an agent-based model characterized by the presence of two kinds of social agents (*A*), producers and users. The agents have only partial knowledge of the world and each agent owns: (i) a set of categories (*C*), utilized to interpret a certain set of artifacts and (ii) a set of weights representing the importance that the agent assigns to each cognitive characteristic of an artifact (represented by means of the correspondent categories). The model postulates a continuous interaction between producers and users: the artifacts are transferred from the producers to the users and subsequent

feedback messages are sent from the users to the producers. In this section, we introduce the general architecture of the model, while later we will introduce some simplifications that make the model more manageable (and which could be relaxed in future research).

15.4.1 Agents

The agents A_i ($i \in \{1, \dots, g\}$, $g \in \mathbb{N}^0$) are the model's main active units. They have limited system knowledge and are distinguishable by means of their identifier ID. Their total number g does not change in time, and they are grouped in two different classes:

- producers (represented by means of the symbol A_p , where $p \in \{1, \dots, l\}$)
- users (represented by means of the symbol A_u , where $u \in \{1, \dots, h\}$)

Thus, the total number of agents is $g = l + h$.

Each agent owns a given number of categories, which can be different for agents belonging to the two different classes. We denote the categories belonging to producers and users by C_p and C_u , respectively¹. The number of categories belonging to a given agent is constant through time. Moreover, each agent is characterized by a weight vector (different for agents belonging to different classes, respectively W_p and W_u).

Note that in our representation, only the producers are able to build and modify the artifacts (one artifact for each category owned by the producer). Conversely, only the users can evaluate the artifacts.

15.4.2 Artifacts

The artifacts are “goods,” built by producers and utilized by users. Each artifact art^s_p is identified by an identification variable, IDs, and corresponds to only one category (each category belonging to a particular producer ID p). The artifacts art^s_p are characterized by an extremely simple representation: they are D-dimensional vectors, whose elements (indicated by the words “characteristic”, or “feature”) take values in $\{0, 1\}$, where 1 indicates the presence of a given characteristic (feature) and 0 its absence.

Despite this simple representation, the artifacts are suitably defined for conveying information, and can be successfully “interpreted” by users; this fact allows the system to produce interesting behaviors.

¹ A more correct notation would be $C_{k,p}$ and $C_{k,u}$, with $k \in [1, \dots, s]$, s being different for producers and users. Nevertheless, for reason of clarity, in the following we omit such an index: simply, the reader should remember that both users and producers possess more than one category (the producers owning one category for each artifact).

15.4.3 Categories

Categories are the tools the agents use in order to interpret their environment; in particular, the users have to evaluate the artifacts and the producers have to assemble them. In this context, each category is a D -dimensional vector, whose elements $C^{(j)}_x$ (again indicated by the words “characteristic”, or “feature”) are discrete random variables taking values “1” and “0”, respectively with probability τ and $(1 - \tau)$, $\tau \in [0, 1]$. Further, we assume that the number of relevant characteristics (that is, the characteristics corresponding to symbols “1”) is only a fraction η of the total number of features. A category composed only of relevant characteristics would entail the unrealistic assumption that agents attribute relevance to every single perceptible detail.

In general, the agents use the categories with several objectives; in this work, we wish to highlight two of these:

- the agents use the categories to evaluate artifacts, or
- the agents use the categories to produce new artifacts (which could be new kinds of artifacts, or simply copies of already existing ones).

In general, the categories that the agents use in order to interpret the artifacts could be similar, but not identical, to the categories used to produce the same artifacts. This is an interesting interplay, but it is out of the scope of this work: as previously noticed, in the following we make the simplifying hypothesis that a typical agent uses different categories to interpret artifacts or to produce them. The other assumption we make is that we are dealing with two specialized class of agents, the users (who use the categories in order to evaluate the artifacts’ functionality) and the producers (who use the categories in order to process the artifacts).

15.4.3.1 Artifact Evaluation Phase

The evaluation phase is a complex process that could be decomposed into the following simpler steps:

- identification of the artifact’s relevant characteristics
- definition of their positive or negative impact on the artifact evaluation
- quantitative computation of their influence on the global evaluation

With the word “identification,” we mean the process that individuates the features of the artifact the agent suppose interesting; these features correspond to the “1” tags on the category, and in such a way considers the category acts as a “filter.”

Once the interesting features are identified, the agent has to make a first evaluation defining the features’ positive or negative impact. Moreover, for each characteristic the agent could or could not use the information coming from the “context:” for example, an agent could like a sport car such as a Ferrari, regardless of its color, or it could like only a red Ferrari. In the first case it is enough that the feature “Ferrari” is present, whereas for a positive evaluation the second case needs also the “red” tag present. Another possibility is that the agent likes Ferrari cars,

whatever their color, unless the color is yellow: it hates yellow Ferraris (like, for example, one of the authors): in this case, the simultaneous presence of a “1” tag in the “Ferrari” and “Yellow” fields determines a negative contribution. Eventually, for this agent, the fact that a Ferrari has no wings never enters into consideration: not all the possible features are relevant, and in such a way the only feasible “context” is determined by relevant characteristics. The process we described takes into consideration only relevant characteristics, and requires just three levels: “1” for positive contributions, “-1” for negative contributions, and “0” for contributions neither positive nor negative. In order to summarize, for each features the final process takes the shape of an input-output table, where the Boolean inputs coming from the features composing the context determine the three-level outcome (see Fig. 15.1).

In this paper, we use the term “*functional attribution*,” or simply “*attribution*,” referring to: “the functionality carried out by the corresponding feature of the artifact I’m evaluating – given the context I’m considering – is useful/indifferent/damaging.”

This is a very complex step, and, in general, it is very hard to insert realistic details. However, a first approximation, we suppose here that the context of each feature is random (the other features involved in the evaluation are randomly chosen), and that the outcome of the evaluation is random (for each “input” combination there is a given probability that the outcome be “1”, “0” or “-1”). This approach is by now a standard one and, for example, was successfully used in classical work in

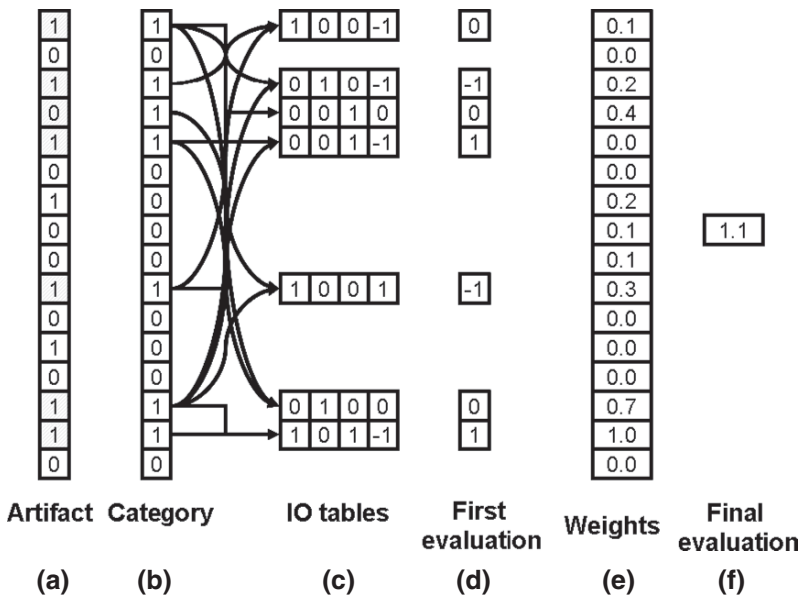


Fig. 15.1 The artifact evaluation process. The category identifies the artifact’s relevant features (in colored background in -a-); the values of these characteristics are filtered through the category (b) and passed to the IO tables (c); the outcomes of this table constitutes the first feature evaluations (d); the scalar product of the outcomes of the IO tables and the weight vector (e) constitutes the artifact evaluation (f)

theoretical biology by Stuart Kauffman (1993). In this work and its descendants, a common practice is that of identifying a unique value, common to all the relationships, for the number of the inputs constituting the context. In this chapter we are following the same procedure, identifying therefore with the variable k the number of inputs of the input-output tables (IO table in the following part of this chapter).

During the last step, the agent provides a quantitative evaluation by multiplying the result of this first evaluation by a weighting factor that depends upon the particular agents' propensities. For this purpose, each agent owns a D-dimensional vector of weights, whose elements, $W_p^{(i)}$ and $W_u^{(j)}$ take values in the interval $[0,1]$. Each element represents the importance that the agent assigns to the corresponding cognitive feature of one of its categories. For example, we could appreciate a car for several reasons: power, style, size, price, color, maintenance costs, practicality, and so on; nevertheless, not all these characteristics have the same influence on our judgment. $W_x^{(i)}$ ($x \in \{u, p\}$) represents the weight that the agent assigns to the i -th cognitive feature of the D-dimensional cognitive space (color is more important than maintenance costs, or class more than price, etc.). In the present version of the model, these vectors are built during the initialization phase at $t = 0$ and are constant in time.

To summarize this important process, in EMIS we define the "functionality of an artifact with respect to a particular category," as an index measuring the level of user satisfaction with the artifact. The index appraises the satisfaction received from the artifact, when the user evaluates the artifact by filtering it by means of the corresponding category. In order to evaluate an artifact, the agent must "interpret" it by means of its IO tables. We can define the artifact functionality (with respect to a given category) simply as the scalar product between the vector resulting from the IO tables and the agent's weight vector W_u .

15.4.3.2 Artifact Building Phase

In this section, we introduce the production of artifacts. Typically, an agent tries to build an artifact as much similar as possible to the prototype memorized in one of its categories, by processing already existing artifacts. It tries to add some desirable characteristics to the artifact, but despite its efforts it has to deal with errors and physical/technical constraints.

The artifact building phase is a process that may be decomposed in a series of steps very similar to the evaluation phase:

- identification of the artifact's relevant characteristics
- definition of their positive or negative impact on the artifact's evaluation
- artifact processing

Typically, the agent already used the category in order to build an artifact; therefore, the agent takes into consideration the same artifact in order to ameliorate it. As during the evaluation phase, the agent identifies the relevant features by means of the "1" tags of the category it uses to process the artifact, and defines their positive

or negative impact by means of input-output tables. We remember that in general these latter tables are not identical to the first ones: the evaluations of a user could be different from the needs and the considerations of an artifact producer.

However, at this point, we can introduce a suitable simplification. The aim of the producer is to build an artifact useful to the users: but to do that, it does not need to reproduce exactly the details of the set of (category + IO tables) present on the users' minds. As we will see below, the users communicate to the producers the position of the features of the artifacts they like/dislike. In this case, for the producers it is enough to integrate in their categories the users' information, and subsequently to try to introduce on the artifact this same information, without further modifications. In other words, if we assume that the particular producer's desires (or idiosyncrasies) have only second-order effects, we can assume that the producer's IO tables have the simple form of an identity (only one entry, with the outcome identical to the entry value) for a desired characteristic, and of a negation (only one entry, with the outcome reverse with respect to the entry value) for an undesired characteristic.

As a last step, the producer attaches the tag "1" to the artifact in correspondence to the outcome "1" of the IO tables, and a "0" tag in correspondence to a "-1". Their relevant features, or the "0" outcomes of the IO tables, mean that in this positions the producer is not interested in changes, and that therefore the already present values will not altered (see Table 15.1).

The feature correspondence represented in Table 15.1 can be summarized as follows:

- if the values are identical, there is no change in art^s_p
- if $C^{(j)}_p = 1$, then $art^s_p = 1$
- if $C^{(j)}_p = 0$, then art^s_p is unchanged
- if $C^{(j)}_p = -1$, then $art^s_p = 0$

If some additional restrictions were not imposed, this process would lead to the "perfect" artifact, where all the desired characteristics are present at the maximum level (e.g., think of a car that can fly, navigate, interact with human beings, produce and translate documents, and make excellent coffees!). In order to take into account these constraints, we decided simply to limit the number of characteristics present simultaneously in the same artifact. If, after producer processing, an artifact has a number of 1's exceeding the given threshold σ , a stochastic removal process eliminates a subset of the current characteristics.

Finally, we remark that since the final goal is to simulate exaptation phenomena, we can disregard an overly detailed description of production processes and costs.

Table 15.1 Example: artifact production/innovation

Considered vector	Features value					
Artifact – A_p	1	0	1	0	1	0
Outcome from the IO tables	1	1	0	0	-1	-1
New artifact	1	1	1	0	0	0

15.4.4 Model Dynamics

The most critical interactions for the outcome of the model take place between producers and users. Two distinct parts compose the interaction process: (i) when the user receives and evaluates an artifact built by a producer; and (ii) when the user provides feedback evaluation to the producer about the satisfaction level reached by the artifact (the artifact functionality).

First, we focus attention on the delivery and subsequent evaluation of an artifact. In order to evaluate the artifact the user: filters the artifact with respect to all its categories; computes its functionality (already described in the previous paragraphs); and finally communicates the best result to the producer. Moreover, the user can deliver to the producer some additional information, which can be useful for future artifact innovations. Specifically, the user can transmit to the producer particular subsets of the two categories that give the highest functionality values. These subset are composed of

- I_{jq} ($q = 1$), *Actual Information (AI)*. Groups of features of the selected artifact that correspond to the inputs of the characteristics that contributed highly to the determination of the functionality value; at the same time, the user communicates their positive or negative contribution to the global evaluation;
- I_{jq} ($q = 2$), *Desired Information (DI)*. Groups of features of the selected category that potentially have the highest positive contribution power (that is, among the interesting features, the input combinations that have as outcome a “1” and that correspond to the highest values of the weights W_u).

These two different kinds of information allow the user to make explicit its requests. In particular,

- AI represents the features of the current artifact that make a positive or negative contribution to the functionality; and
- DI expresses what the agent likes or dislikes about the artifact.

Sometimes, it is possible that the two subsets have a non-empty intersection: for example, it is possible that an artifact feature (or a combination of features) is important for a category, and, at the same time, it gives a negative contribution to the functionality of another category. In such a case, we impose that either AI or DI is randomly left out in the transmission.

The producer uses the transmitted features to modify the features $C_p^{(i)}$ of the category employed to build the artifact. The producers have to integrate the new information coming from the users with the information already present on their categories. In order to do this action, the producer

- Copies, in its category, the feature (or the combination of features) giving a positive contribution to the functionality and sets to identity the corresponding IO tables; and

- Copies, in its category, the feature (or the combination of features) giving a negative contribution to the functionality, and sets to negation the corresponding IO tables.

The goal of the above calculation is to transform the features communicated to the producer by the user. All the remaining features are left unchanged, except some random “noise” (with a given probability, some “1” becomes “0” or some “0” becomes “1”, and similarly some identity becomes a negation or some negation becomes identity). Finally, in order to limit the number of relevant features, the new vector is filtered by the removal process described in the previous section, according to the given threshold σ . The final vector represents the new category that the producer employs with the final IO tables to build the subsequent artifacts’ generation.

15.4.5 EMIS and the Study of Exaptation

In the initial paragraph, we have defined exaptation as an emergent phenomenon in evolution dynamics. Where does exaptation appear in our model? Recall that EMIS simulates exchanges of products (artifacts) between producers and users, the users evaluating the artifacts by means of their tables and categories. In this context, an exaptation is a category change in interpreting the artifact. For example, after hundreds of steps the category that systematically was returning the best functionality might no longer be the best one: the last innovation(s) has (have) increased the functionality of another category, which in such a way becomes the new reference category for the selected artifact.

Recall that the producer supplies the user only with the best functionality value among all the values computed using all the categories it owns. The category that furnishes such a best value during one interaction is likely to have a large value also in the next interaction, and so forth. In a sense, the best category is, for the user, the “leading” category for this particular artifact. Sometimes, but quite rarely, another category reaches a functionality value larger than that of the leading category. In a sense, this can be interpreted as a variation of the utilization context of the artifact under consideration; in this case, we are observing an exaptation event.

15.4.5.1 The $k = 1$ Model

As we previously stated, in order to perform a first analysis of the model behaviors, we introduce some simplifications. The IO tables represent a flexible but complex feature of the model, and sometime this machinery is over-dimensioned with respect to the evaluation task. A straightforward simplification is that of setting to 1 the number of entries of IO tables: in such a way, we are neglecting the complex internal relations some categories can have, in order to focus our attention on the information exchange among agents. A very interesting result is that in this context exaptation can be present.

The $k = 1$ simplification allows an easier realization of several model steps. In particular:

Table 15.2 Example of semantic distance calculation between artifact and category

Considered vector	Features value					
Artifact – A_p	1	0	1	0	1	0
Category – A_u	1	1	0	0	-1	-1
Distance	0	1	0	0	1	0

- the categories can include directly the first artifact evaluation, composed of elements $\{1, 0, -1\}$ and eliminating in such a way the IO tables
- the information provided by the user to the producer is composed only of features and not of a group of features
- in order to incorporate the users' information, the producer simply computes feature by feature the average among its values and the incoming ones
- the production of the artifact is straightforward

It is easy to find the largest functionality value F_{max} of an artifact with respect to a particular category: it is simply given by the scalar product of the weight vector and the category itself, where all the “-1”s are set to “0”. Eventually, it is possible to define in an intuitive manner a semantic distance between 2 categories: in this context, the distance describes the discrepancy between the artifact and the category referred to that particular artifact. Table 15.2 is based on the two following principles:

- a distance between artifact and category occurs when the agent’s “desires” a feature that is not currently present
- a distance between artifact and category occurs when the agent “does not desire” the present feature

A detailed description of this realization can be found in Villani, Bonacini, Ferrari, Serra, and Lane (2007); in the following part of this chapter we comment on the results obtained using this $k = 1$ model and compare them with its full version with $k > 1$.

15.5 Model Dynamics

15.5.1 $k = 1$ Model: the Initialization

In order to test the model, simple simulations for which we reduce the number of the actors, maintaining only one producer and one user:

$$l = 1 \quad h = 1 \quad g = 2$$

The knowledge space of the categories involved is $D = 1,000$ features. The user utilizes five categories, which guarantee sufficient diversity, while the producer owns only one category, corresponding to the artifact that it is building. The other parameters we have to fix in order to create the initial categories are:

1. the threshold η , limiting the number of 1s in $C^{(j)}_i$, is set to 100
2. the initial fraction of 1s in $C^{(j)}_i$, is set to 0.05
3. the initial fraction of -1 s in $C^{(j)}_i$, is set to 0.05

The value of η is relatively low with respect to the D value and indicates that the space of all the possible characteristics of a category is very large with respect to the really imagined ones. Each feature of the initial artifact (at time $t = 0$) takes the value 1 with probability

$$P_{art} = \frac{\sigma}{D} \quad (15.1)$$

where σ is the threshold defining the maximum number of 1s that can be present on an artifact, and D is the total number of features composing the knowledge space. Thus, initially we have a “raw” artifact that is able to contain a large number of details. This raw artifact is processed a first time by the producer, to correctly link each artifact to its referring category. Because of the fact that $\sigma > \eta$, the artifact carries out a number of functions larger than the number of functionalities for which it has been selected (exaptive pool of possibilities).

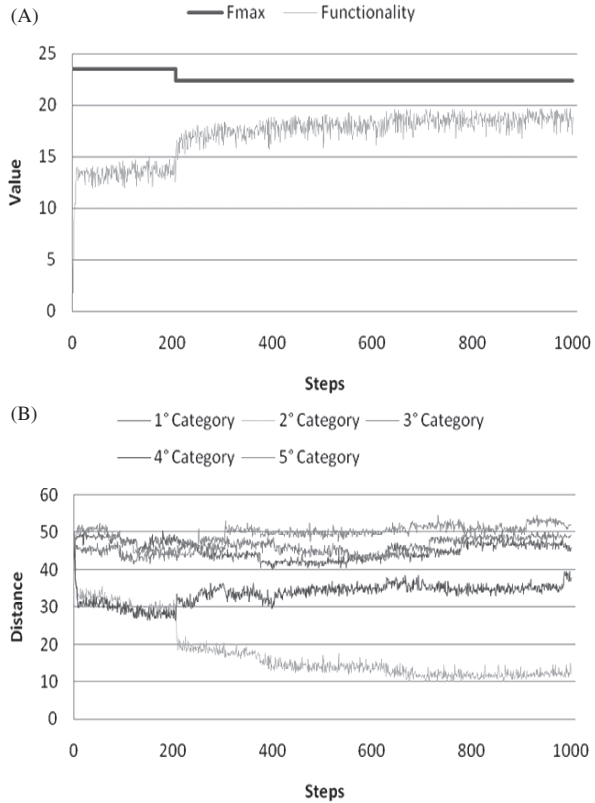
15.5.2 A Typical Run

Figure 15.2 shows some of the main variables simulated by EMIS. Figure 15.2a shows the best functionality value of the artifact at a given time step, and its F_{max} value. Note that at step 207, the user changes the leading category (exaptation) and F_{max} . In this case, the upper limit for the achievable satisfaction decreases, despite that an increase of the actual functionality. Figure 15.2b shows the distance among user’s categories and the artifact actually built by the producer.

At each step, the user analyses the producer’s artifacts, which are interpreted using its own categories and furnishes to the producer the corresponding (more elevated) value of functionality. At step 207, the functionality value of category 2 outperforms the functionality value of category 4. During the rest of the simulations, category 2 maintains its superiority and we do not observe other exaptation events.

The exaptations occur mostly during the first steps of the simulations, whereas they are rarely observed during the long stabilization period (that is – by definition – after 50 steps). Actually, during the first interactions the producer can easily modify its artifacts and simultaneously increase the functionality values with respect to several different categories. However, when the artifact is highly specialized the simultaneous satisfaction of several requirements is a challenging task. Conversely, symmetric information and large bandwidth values can support best performances of the model in terms of attainment of F_{max} , allowing the producer to satisfy the user’s requests.

Fig. 15.2 (A) The best functionality value of an artifact in time, and its F_{max} value. At step 207, the user changes the referring category (exaptation); similarly the F_{max} value changes. In this case, the maximum reachable satisfaction decreases despite the actual functionality value increases. **(B)** The distance among the user categories and the artifact built by the producer



15.6 Experiments

We individuate three main factors that are able to favor the emergence of exaptation phenomena:

1. communication among different agents,
2. communication and production noise, and
3. evolution of the users' categories.

In all the situations, the user can utilize two different modalities of communication:

1. symmetrical: both communicated categories provide the same amount of information, “actual” and “desired” or
2. asymmetrical: the category that better interprets the artifact provides “actual” information, while the second one provides “desired” information

This specification has been introduced to verify whether the presence of communication of symmetry/asymmetry can favor the occurrence of exaptation phenomena. Recall that the “actual” information represents the objective user’s

evaluation, whereas the “desired” information expresses what the user wants. Therefore, symmetric communication means that the user treats the categories without any bias, whereas asymmetric communication means that the user transmits an objective report about its first category, and then communicates some desires corresponding to a second one.

15.6.1 Communication

Typically, the user transmits to the producer a (small) subset of the features extracted from the two categories that are returning the best functionality values. In this paragraph, we analyze the behavior of the $k = 1$ model by varying the number of transmitted features, or bandwidth (B). In particular, in this set of experiments the total number of transmitted characteristics is set to be $B = \{20, 40, 60, 80, 100, 200\}$.

The data are noisy, but we can propose two considerations:

- exaptation phenomena have a weak link with the bandwidth; nevertheless, it seems that larger values of the bandwidth correspond to more numerous late exaptation events
- the number of exaptations found by using the asymmetric modality is larger than the number of exaptation found by using the symmetric modality

These facts suggest that an unbiased communication modality is not able to favor a context change, whereas a qualitatively asymmetric communication modality could effectively support the success of categories that are not favored initially. Furthermore, to favor exaptation phenomena, it seems to be helpful to transmit a large amount of information, creating the potential for the discovery of new and previously disregarded solutions.

15.6.2 Noise

A second study concerns the analysis of two types of noise that could be specified in the model:

- *Communication noise.* The value of the features communicated by the user agents is changed with probability α , (1 or -1 from a 0 tag, 0 or -1 from a 1 tag, and 1 or 0 from a -1 tag).
- *Production noise.* The value of the characteristics of the artifact built by the producer is changed with probability β (0 from 1, and 1 from 0).

The communication noise does not affect the main behavior of the model, although its presence requires more time to reach higher functionality values. The production noise worsens such a tendency, but at the same time increases the frequency of exaptation occurrences, both in the long period and in presence of low bandwidth

communications. Some innovations, obtained because of this type of error, are able to foster a change of context for the whole artifact.

15.6.3 Learning

For the $k = 1$ model, it is possible to propose a simple and interesting learning modality. Agents can modify their categories by learning from the environment: at each time step there is the probability f of selecting a category for a modification; in this case, a fraction P_{ch} of the category's features change their value (in all our trials $P_{ch} = 0.02$). The general schema of this learning is that the differences between artifact and category tend to decrease and that the category is plastic. Therefore: (a) if the artifact feature is "1", the corresponding feature of the new category becomes "1", unless the involved characteristic be "not desired" (in this case, it reduces the distance and becomes "0"); (b) if the artifact feature is "0", the corresponding feature of the new category becomes "0", unless the involved characteristic be "not desired" (in this case, it maintains the old value).

With these settings, we find that the exaptation frequency increases almost linearly with the growth of the adjournment rate f (see also Fig. 15.4a) and that the F_{max} level of satisfaction increases with the adjournment rate (Fig. 15.4b). A real example of "learned" exaptation could be that of the SMS (the Short Message Service), initially introduced to send brief official messages from the telephone company, which subsequently became a new means for social communication (especially among the young!). In this case the users succeeded in understanding the communicative potentiality of this system, overcoming limits of space by means of the creation of a particular language, more "assembled" and intuitive. As a result, the phone companies initiated a market strategy based upon this new functionality.

15.6.4 $k = 1$ First Conclusion

From these first simulations, we can conclude that some of the most important elements favoring the emergence of exaptation events are:

1. an asymmetrical communication, where evaluations and desires are differently expressed for different categories;
2. a high number of cognitive features (characteristics) communicated among the agents;
3. a high level of production noise; and
4. the plasticity of the users' categories.

15.6.5 Higher k

Now we can insert in the system the relationships among the different features belonging to the same category, following the same stages as in the previous analysis.

The only exception concerns the learning experiments, because the presence of the feature contexts deepens the problem and increases considerably the cognitive dimensions involved; these aspects will be analyzed in future work.

We performed simulations with $k = 2$, $k = 3$, and $k = 4$, generally similar to the ones depicted in Figs. 15.2 and 15.3. During this first session of trials, we assumed a uniform distribution of “1”, “-1” and “0” inside the IO tables of the users’ categories; of course, in the near future, we are planning a huge simulation campaign in order to deepen the first conclusions presented in this chapter. The performed simulations confirm that the exaptation phenomena have a weak link with the bandwidth, and do not reveal important differences between symmetric and asymmetric modalities.

The influence of noise is similar in both systems, that is, the communication noise has no effects whereas the growth of the production noise increases the exaptation occurrences (Fig. 15.5a); Fig. 15.5b confirms that the differences between the symmetric and asymmetric modalities are not so great. Figure 15.5c shows that as production noise increases, this also augments in a significant manner the long term exaptations (the exaptations that happen after the threshold of 50 steps); this phenomenon is not so evident in the case of the $k = 1$ model. Last, but not least, the increase of the number of inputs of the IO tables does not have significant consequences (Fig. 15.5d).

To summarize the previous paragraphs, as we supposed, the simple case of $k = 1$ is able to capture the main exaptation phenomenon traits, and it could be used profitably to scout for new behaviors, which can be confirmed successively by the

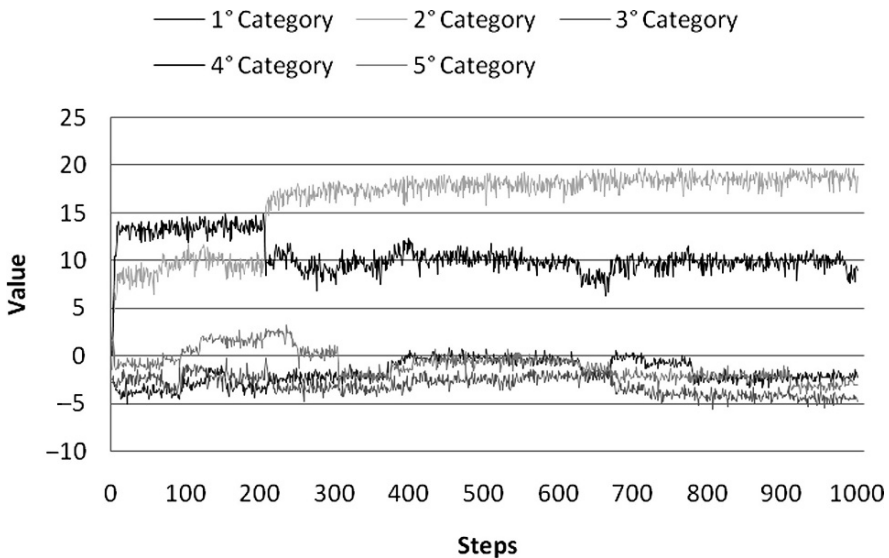
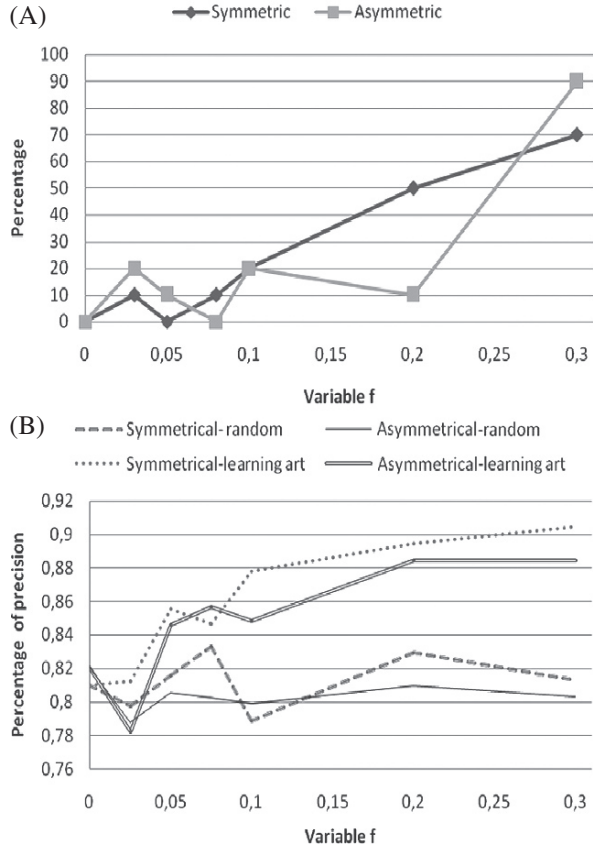


Fig. 15.3 The figure shows the evaluation of the artifact received from the agent’s categories, during the same simulation of Fig. 15.2. Please note that, at step 207, the second category suddenly increases its value and becomes the referring category

Fig. 15.4 (A) Learning modality: the percentage of simulations (over 10 runs) with at least one exaptation, by varying the adjournment rate f . (B) The fraction of F_{max} reached by the artifact functionality, by varying the adjournment rate f for both random and learning modality



general model with $k > 1$. The presence of IO tables allows a higher probability of exaptation occurrences in long term sequences, and we have some indications that it is enough to use a limited context in order to have complex situations.

15.7 Conclusion

In this chapter, we have introduced an agent-based model designed to investigate the dynamics of some aspects of exaptation in a world populated by agents, whose activity is organized around production and utilization of artifacts. The model, EMIS (*Exaptation Model in Innovation Studies*), not only explicitly includes agents and artifacts, but also encompasses the agents’ subjective representation of the artifacts by means of cognitive categories. In our model, the agents build (as producers) and interpret (as users) the artifacts using cognitive categories.

One of the main features of EMIS is the description of the information exchange dynamics among agents. Two main processes characterize these dynamics: (i) interpretation and storage of information by each agent, and (ii) circulation of

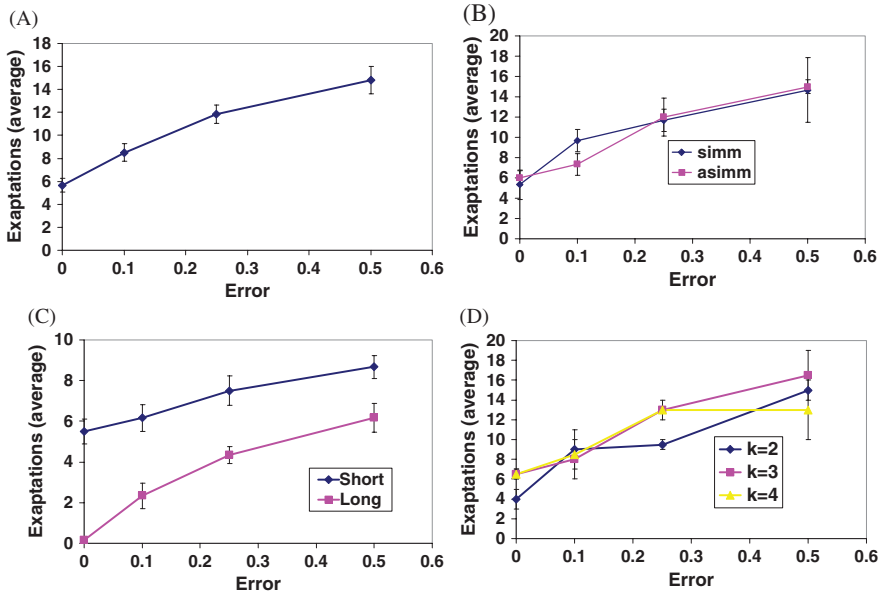


Fig. 15.5 The effect of production noise. (A) The average of exaptation occurrences as function of the production error; (B) the average of exaptation occurrences for symmetrical and asymmetrical communication modalities; (C) the average of exaptation occurrences in short and long runs; (D) the average of exaptation occurrences for systems with different k values

information through the exchange of artifacts. The first process takes place in the cognitive domain of each agent. In particular, at any given time, each agent owns a cognitive representation of a number of real objects (artifacts) in terms of a set of cognitive features (categories). The second process involves the communication among agents of a fraction of the information stored in categories. As a last step, the producers employ their sets of cognitive categories to make artifacts that are in turn submitted to the users. The users evaluate the functionality of such artifacts by means of their cognitive categories and next send signals to producers about their “satisfaction” with such artifacts.

The agents in the model are able to attribute “functionalities” to the artifacts in terms of categories and IO tables; a given attribution can generate a certain reward associated with the artifact of reference. Thus, the type of representation proposed is suitable to take into account exaptation events, which are understood as shifts in terms of the “leading attributions” (attributions corresponding to highest reward) that the agents assign to the artifacts.

The model has been implemented in a computer environment. The main goal of the computer simulations is to determine which factors play a central role in the information exchange processes, with particular attention to the study of exaptation. From our first simulations, we can conclude that some of the most important elements favoring the emergence of exaptation events are

1. a high level of production noise and
2. the plasticity of the users' categories.

There are also indications that an asymmetrical communication (where evaluations and desires are differently expressed for different categories, and a high number of cognitive features, or characteristics, are communicated among the agents) could be helpful to increase the exaptation occurrences.

While the present version represents a highly simplified model, the results obtained so far appear to encourage the further development of EMIS. Future research improvements to the model should be aimed at expanding the exploration of the synergistic effects among the different features belonging to the same category, and simulating more complex interactions among categories and artifacts.

References

- Ceruti, M. (1995). *Evoluzione senza fondamenti*. Laterza, Italy: Bari.
- Dew, N., Sarasvathy, S. D., & Venkataraman, S. (2004). The economic implications of exaptation. *Journal of Evolutionary Economics*, 14(1), 69–84.
- Gould, S. J. (2002). *The structure of evolutionary theory*. Cambridge, MA: The Belknap Press of Harvard University Press.
- Gould, S. J., & Verba, E., (1982). Exaptation, a missing term in the science of form. *Paleobiology*, 8(1), 4–15.
- Kauffman, S. A. (1993). *The origins of order*. Oxford, UK: Oxford University Press.
- Mokyr, J. (1998). Induced technical innovation and medical history: an evolutionary approach. *Journal of Evolutionary Economics*, 8(2), 119–137.
- Schumpeter, J. A. (1934). *The theory of economic development*. Cambridge, MA: Harvard University Press.
- Schumpeter, J. A. (1976). *Capitalism, socialism and democracy*. New York, NY: Harper and Row.
- Simon, H. A. (1996). The architecture of complexity: hierarchic systems. In *Sciences of the artificial* (3rd ed., pp. 183–216). Cambridge, MA: MIT Press.
- Villani, M., Bonacini, S., Ferrari, D., Serra, R., & Lane, D. (2007). An agent based model of exaptative processes. *European Management Review*, 4, 141–151.

Chapter 16

Power Laws in Urban Supply Networks, Social Systems, and Dense Pedestrian Crowds

Dirk Helbing, Christian Kühnert, Stefan Lämmer, Anders Johansson, Björn Gehlsen, Hendrik Ammoser and Geoffrey B. West

16.1 Scaling Laws in Urban Supply Networks

The classical view of the spatio-temporal evolution of cities in developed countries is that urban spaces are the result of (centralized) urban planning. After the advent of complex systems' theory, however, people have started to interpret city structures as a result of self-organization processes. In fact, although the dynamics of urban agglomerations is a consequence of many human decisions, these are often guided by optimization goals, requirements, constraints, or boundary conditions (such as topographic ones). Therefore, it appears promising to view urban planning decisions as results of the existing structures and upcoming ones (e.g. when a new freeway will lead close by in the near future). Within such an approach, it would not be surprising anymore if urban evolution could be understood as a result of self-organization (Batty & Longley, 1994; Frankhauser, 1994; Schweitzer, 1997).

Comparison with biological systems promises further insight. Quantities like metabolic rates, population growth, life-span, etc. have been discovered to scale with the average body mass of biological species over about 20 orders of magnitude (West, Brown, & Enquist, 1997; Enquist, Brown, & West, 1998). The corresponding power laws reflect the underlying function, structure, and organization of biological species and even extend to the realm of ecological systems such as natural forests with different sized trees. For example, it turns out that all trees of one size class consume the same amount of solar energy as trees of a different size class (Enquist et al., 1998).

It would be interesting to find out, whether a system of cities could be viewed as an ecological system with similar relationships. In this connection, it is useful to remember Zipf's (1949) law, according to which the population sizes of cities are inversely proportional to their rank. This implies the relationship $n_k \propto 1/N_k$ for the number n_k of cities of size class k (e.g. with more than 5×10^k but less than $5 \times 10^{k+1}$ inhabitants). Therefore, as the energy usage E_i by the population of a city i (when

D. Helbing (✉)

Institute for Transport and Economics, TU Dresden, Andreas-Schubert-Str. 23, 01062 Dresden, Germany; Collegium Budapest – Institute for Advanced Study, Szentháromság u. 2, 1014 Budapest, Hungary

Table 16.1 Scaling exponents and their 95% confidence intervals for different variables of electric energy supply in Germany as function of population size. For details, (see Kühnert et al., 2006)

Variable	Exponent	95% Confidence interval
Usable electric energy	1.1	[1.04, 1.13]
Electric energy delivery to households	1.0	[0.96, 1.06]
Length of low-voltage cables	0.9	[0.82, 0.92]

neglecting the energy consumption by industrial production) grows linearly with the population N_i (see the entry “electric delivery to households” in Table 16.1), the number of cities of size k times their energy usage is constant. In other words, the inhabitants of all cities of one size class k consume the same energy as the inhabitants of all cities of any other size class, similar to the ecological example of trees in a forest.

Among the many different approaches trying to explain Zipf’s law (e.g., Simon, 1955; Steindl, 1965; Schweitzer, 2003), the one by Gabaix (1999) is surprising because of its simplicity. According to Gabaix, the simplest stochastic model with multiplicative noise $\xi_i(t)$, namely

$$\frac{dN_i}{dt} = [A + \xi_i(t)] N_i(t), \quad (16.1)$$

is able to generate Zipf’s distribution. In agreement with “Gibrat’s law” (Gibrat, 1931; Sutton, 1997), it assumes that the growth rates $A_i(t) = A + \xi_i(t)$ are stochastically distributed and varying around a characteristic value A independent of the (population) size $N_i(t)$ of a city i . Note, however, that the exponent of Zipf’s law seems to be different from 1 in some countries (Pumain, Paulus, Vacchiani, & Lobo, 2006).

Therefore, let us discuss the consequences if the deterministic part of the growth law would be slightly different from Equation (16.1), namely of the form

$$\underbrace{\frac{dN_i}{dt}}_{\text{Growth}} = \underbrace{BN_i(t)^\beta}_{\text{Resource Generation}} - \underbrace{CN_i(t)^\gamma}_{\text{Maintenance}} \quad (16.2)$$

This equation reflects that the difference between the generation of resources N_i of system i and its maintenance determines its growth dN_i/dt in time. B and C are treated as constants. The powers β and γ allow one to take into account scaling exponents different from 1. While the case $\beta = \gamma = 1$ corresponds to Equation (16.1), any difference of one of the exponents β or γ from 1 would have dramatic consequences.

Equation (16.2) has a surprisingly rich variety of solutions (see Fig. 16.1). If $BN_i(0)^\beta - CN_i(0)^\gamma < 0$, we have either a decay to a finite value, a decay to zero, or an unexpected, delayed decay to zero, depending on whether β is smaller than, equal to, or greater than γ . In the case $BN_i(0)^\beta - CN_i(0)^\gamma > 0$, we find a limited growth for $\beta < \gamma$, and an exponential growth for $\beta = \gamma$, as for the deterministic version

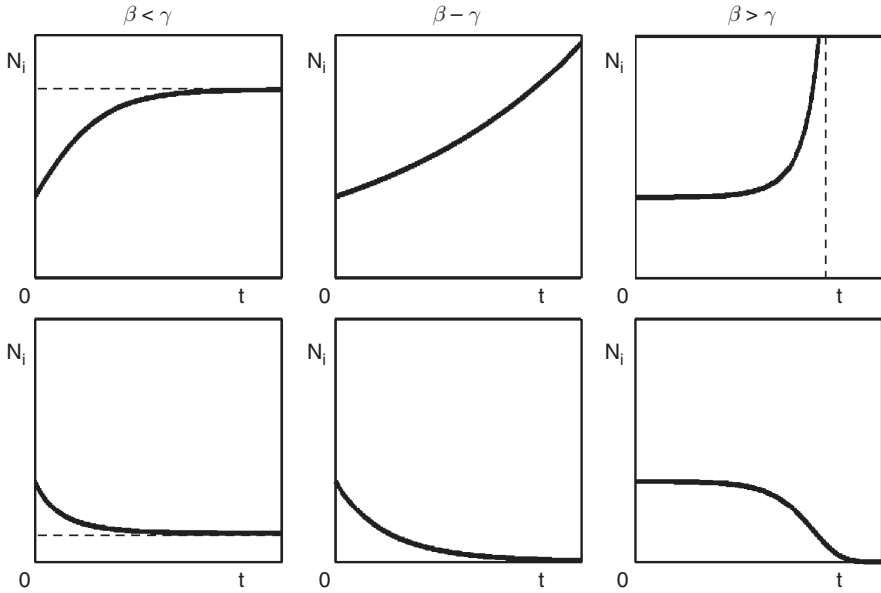


Fig. 16.1 Schematic illustration of the different possible solutions of the growth Equation (16.2). The course of the growth behavior depends on the relations between the parameters β , γ , and on the initial value $N_i(0)$. The *top row* is for $BN_i(0)^\beta > CN_i(0)^\gamma$, while we have $BN_i(0)^\beta < CN_i(0)^\gamma$ in the *bottom row*

of Equation 16.1. However, if β were greater than γ , the growth curve would have a singularity, i.e. it would increase without limits within finite time. This possibility would have dramatic implications for urban systems, as the system would sooner or later go out of control. It would also be a distinguishing feature from biological systems, as these are usually characterized by scaling exponents β smaller than 1 and $\gamma = 1$. Moreover, growth processes of biological species sooner or later saturate similar to the curve displayed in Fig. 16.1a.

To determine the nature of urban growth processes, we have analyzed data of European cities to reveal some of the fundamental forces at play in the formation and development of urban organization. Our empirical results show that, in spite of the enormous variation of particular features (climate, economic specialization, age), cities are unified by mechanisms that are on average simple scaling functions of their population size. Our data sets of urban supply systems for European cities i larger than 50,000 inhabitants contained information about the local energy consumption in German cities and Western European points of interest collected by TeleAtlas[©] for route guidance and geo-information systems. We have evaluated variables X_i and countries for which the data sets were either close to complete or a good statistical representation. The underlying rationale for our empirical investigation was to collect measures of resource production and consumption as a function of urban size, measured in terms of population N_i .

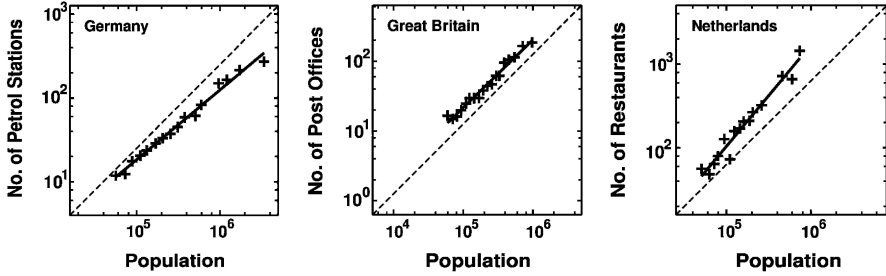


Fig. 16.2 Examples of supply systems with (a) sublinear, (b) linear, and (c) superlinear scaling. The figures show the number of supply stations as a function of the respective population sizes of cities in double-logarithmic representation, using the logarithmic binning method. For details, see Kühnert et al. (2006)

The first empirical fact to emphasize is that scaling is a wide-spread property of urban organization. For most countries, we found power-law scaling relations over two orders of magnitude in population size N^i (see, for example, Fig. 16.2). The scaling relations have the simple form

$$X_i = X_0 N_i^\beta \quad (16.3)$$

where N_i is the population size, X_0 a normalization constant independent of N_i , and β the scaling exponent. Our results are summarized in Table 16.1 and Fig. 16.3. For details, see Kühnert, Helbing, D., and West (2006).

Despite the width of the confidence intervals, one can draw several interesting conclusions:

1. The scaling exponents of different countries are consistent, i.e. of the same order. In fact, the 95% confidence intervals tend to have a common subset, which may be used for a more precise determination of the respective scaling exponent, if universality (i.e. country-independence) is assumed. Statistical analysis of variance tests support this picture.
2. A proportionality of the number of “supply stations” to the population size corresponding to a scaling exponent of 1 is only found for some supply systems. This includes hospitals and hospital beds, post offices, and pharmacies.
3. There are also cases of sub- or superlinear relationships. For example, the scaling exponents for the number of car dealers and petrol stations are smaller than 1 (sublinear case), while the scaling exponent for restaurants is larger than 1 (superlinear case).

What are the reasons for observed differences in the scaling exponents? The proportionality of post offices, pharmacies, and doctors to the population size is probably dictated by comparable individual demands, combined with the requirement of a certain level of reachability (by foot). Moreover, it is often regulated by government. As a consequence, each size class of cities offers approximately the same number of these “supply stations.”

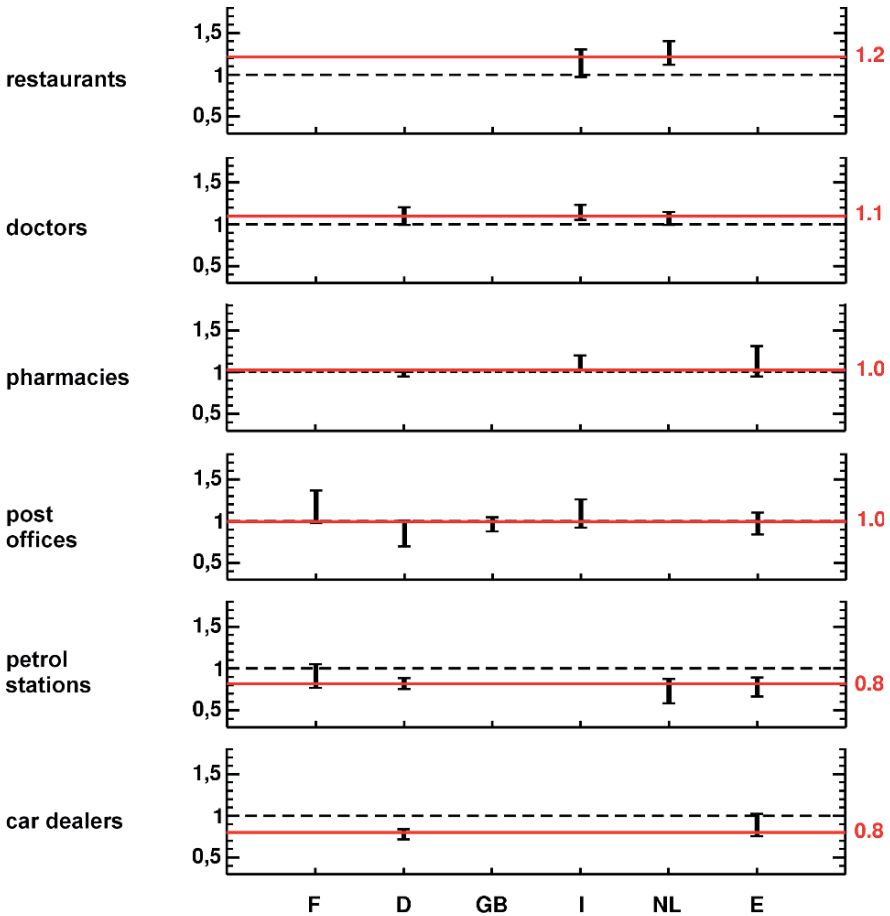


Fig. 16.3 Scaling exponents and confidence intervals for different supply systems and countries: France (F), Germany (D), Great Britain (GB), Italy (I), The Netherlands (NL), and Spain (E) (after Kühnert et al. 2006)

Sublinearly scaling quantities, such as the number of petrol stations or car dealers, indicate an “economy of scales,” i.e. efficiency gains by serving larger agglomerations. In other words, one such “supply station” serves more people in a larger town and distributes larger quantities (e.g., sells more fuel per month). This is certainly reasonable and typical for material supply systems, which are more profitable for larger population sizes or governed by a free market. Therefore, sublinear scaling supply systems profit from higher population densities and a more efficient usage of capacities in larger service units (e.g., by better utilization or reduction of the relative statistical variation etc.). Sublinear scaling is also expected for the number of shopping centers or for polyclinics.

But why do some supply systems scale superlinearly? This question concerns, for example, the number of restaurants, but a similar relations seem to hold for

museums, theaters, colleges, etc. We recognize that these supply systems satisfy social and communicative needs. That is, information exchange seems to increase overproportional with the number of inhabitants in a town. The number of patents, as a function of the population size (Strumsky, Lobo, & Fleming, 2005), and other variables (Pumain et al., 2006) confirm this conclusion. The same applies to other non-conserved variables such as money or wealth (Bettencourt, Lobo, Helbing, Kühnert, & West, 2007). If these variable would determine city growth, a finite time singularity would be expected (which would have to be avoided by increasing innovation or friction).

16.2 Scaling Laws in Urban Road Networks

Let us now turn to the questions of how the scaling laws we identified relate to spatial structures of urban organization and to fractal features of supply and transportation systems. In biological systems, the power laws mentioned in Section 16.1 can be explained by a minimization of energy losses in the respective biological supply system with the constraint that the supply system is space-filling, as all elements (e.g., all cells in the body) must be reached (West et al., 1997). This organization principle implies hierarchical and self-similar structures such as the system of blood vessels.

Therefore, are urban transportation networks also organized in a hierarchical, self-similar way? Self-similar, fractal features have, in fact, been found in the organization of cities, according to Christaller's (1980) theory of central places, and in the structure of public transportation systems (Frankhauser, 1994; Hołyst, Sienkiewicz, Fronczak, & Suchecki, 2005). The same applies to urban boundaries and urban sprawl (Batty & Longley, 1994; Frankhauser, 1994; Makse, Havlin, & Stanley, 1995; Schweitzer, 1997, 2003). But what about urban road networks?

Despite distances being very crucial for logistic, geographical, and transportation networks, surprisingly little attention has been paid to the spatial structure of urban networks in the past. Urban road networks with links and nodes representing road segments and junctions, respectively, exhibit unique features different from other classes of networks (Newman, 2002; Jiang & Claramunt, 2004; Buhl et al., 2006; Crucitti, Latora, & Porta, 2006; Gastner & Newman, 2006; Porta et al., 2006). As they are almost planar, they show a very limited range of node degrees. Thus, they can never be scale-free like airline networks or the internet (Gastner & Newman, 2006) Nevertheless, road and airline networks can both be viewed as solutions of an optimization process minimizing average travel costs, if the travel costs for one airline connection are approximately equal, but the travel costs of road traffic are proportional to the length of links. Hence, a small-world network with a hub-and-spoke structure results for air traffic, while a Poisson node distribution is typical for road networks (Gastner & Newman, 2006).

For an empirical analysis, we have extracted road network data of the administrative areas of the 20 largest German cities from the geographical database Tele Atlas

MultiNet™, which is typically used for real-time navigation systems or urban planning and management. The data provide a geo-coded polygon for each road segment as well as a series of properties, e.g., the length, average expected travel-time, speed limit, driving direction, etc. Since the road network of Hanover, which ranked 11th, could not be extracted unambiguously, it was excluded from our analysis.

Note that, according to human perception, the effort of traveling is not measured in distances, but in terms of the energy consumption by the body required to perform the travel activity (Kölbl & Helbing, 2003). This means that travel times are the relevant quantities for the destination and route choice of car drivers. This implies that routes along faster roads appear “shorter” than along slower ones. A distant, but well accessible, destination is virtually closer than a near one with a longer access time. The heterogeneity of road speeds also has an impact on the distribution of vehicular traffic over the road network. Faster roads are more attractive for human drivers, resulting in a concentration of traffic along these “arterial” roads, see Fig. 16.4a.

The importance of a road or a junction can be characterized by the number of cars passing through it within some time interval. This can roughly be approximated with the measure of link betweenness centrality, b_e , and node betweenness centrality, b_v . It is given by the number of shortest paths, with respect to travel-time, between all pairs of nodes in the corresponding graph, of which, the particular link e or node v is a part (Albert & Barabási, 2002; Newman, 2002; Brandes & Erlebach, 2005; Costa & da Rocha, 2006; Porta et al., 2006). The road networks of Germany’s largest cities show an extremely high node betweenness centrality b_v at only a small number of nodes, while its values are very low at the majority of nodes. As a consequence, the distribution of its frequency density distribution $p(b_v)$ follows a power law $p(b_v) \sim b_v^{-\delta}$ with exponent $\delta \approx 1.4$ (see Fig. 16.4b and Table 16.2). Note that values of $\delta > 1$ indicate a high concentration of traffic over a few important intersections. In Dresden, for example, 50% of all road meters carry as little as 0.2% of the total traffic volume only, while almost 80% of the total traffic volume are concentrated on no more than 10% of the roads. Most interestingly, half of the total traffic volume is handled by only 3.2% of the roads in the network.

The bundling of traffic streams on a few arterial roads (see Fig. 16.4a) reflects the clear hierarchical structure of the roads (Levinson & Yerra, 2006). However, the usage pattern does not display the regularity of hierarchical networks such as Cayley trees, in contrast to many supply networks in biology, such as vascular systems (Brown & West, 2000).

Is this result just an effect of the respective urban topography? Or is it a result of the fact that the time span of urban evolution is short compared to biological evolution, so that deviations from optimal (resource-efficient) structures occur? Or do the boundary conditions of urban growth change so fast that urban systems are always in a transient state?

The apparent universality of the scaling exponent, $\delta \approx 1.4$, suggests that there must be other reasons for the irregular and not strictly hierarchical structure of urban road networks. Universal scaling laws are, in fact, a very surprising feature of urban road networks in view of all the particularities of cities regarding their history, climate, economic specialization, etc. At least in Germany, universal power laws are

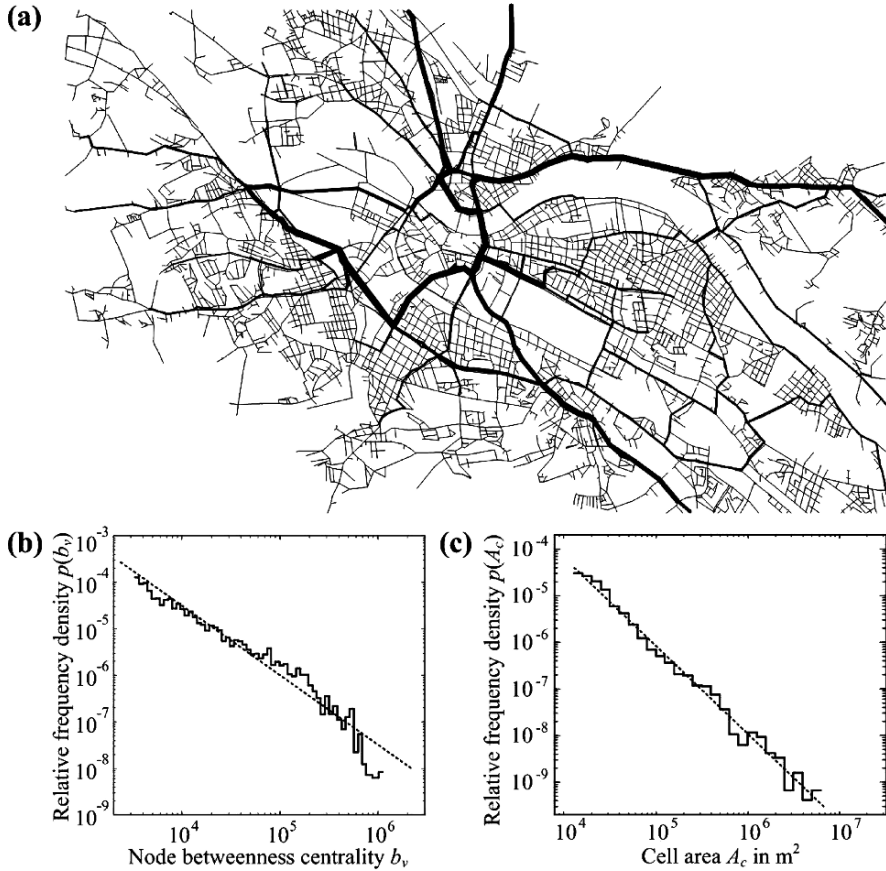


Fig. 16.4 (a) Street network of Dresden, Germany. The width of the links represents the respective betweenness centrality b_v , which is a simple measure of the estimated amount of traffic on the roads. (b) The corresponding distribution of the node betweenness centrality b_v obeys the power-law $p(b_v) \sim b_v^{-\delta}$ with exponent $\delta = 1.36$ (dotted line). (c) The distribution of surface areas enclosed by roads is also power-law distributed with an exponent of 1.89. (After Lämmer et al., 2006)

also found for the size distribution of the areas enclosed by roads (see Fig. 16.4c and the cell size exponent in Table 16.2) and other quantities like the effective dimension or the Gini index (see Lämmer, Gehlsen, & Helbing, 2006 for details).

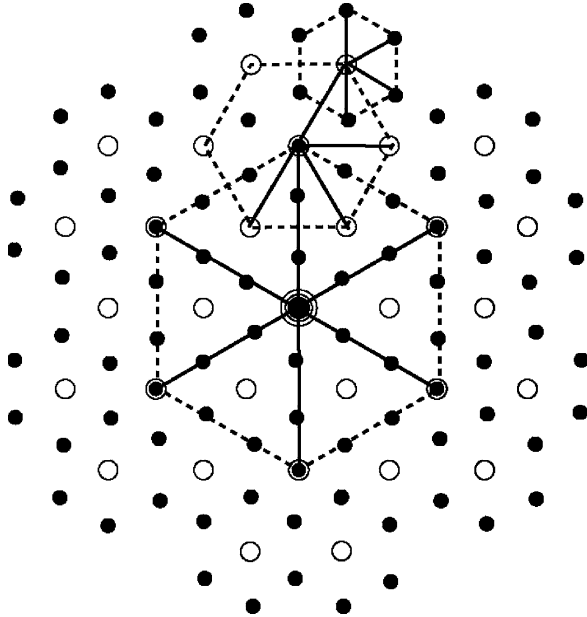
16.3 Deficiencies of Strictly Hierarchical Organizations

Nodes (intersections) and links (roads) of urban networks are often blocked by building sites, accidents, or congestion. This restricts the reliability of nodes and links considerably. We believe that this is a strong reason for network structures

Table 16.2 The 20 largest cities of Germany and their characteristic coefficients; see main text (after Lammer et al., 2006)

City rank	City name	Population in 1,000	Area in km ²	No. of nodes	No. of links	Betweenness exponent δ	Gini index g	Cell size exponent β
1	Berlin	3,392	891	37,020	87,795	1.481	0.871	2.158
2	Hamburg	1,729	753	19,717	43,819	1.469	0.869	1.890
3	Munich	1,235	311	21,393	49,521	1.486	0.869	2.114
4	Cologne	969	405	14,553	29,359	1.384	0.875	1.922
5	Frankfurt	644	249	9,728	18,104	1.406	0.873	2.009
6	Dortmund	591	281	10,326	22,579	1.340	0.875	1.809
7	Stuttgart	588	208	10,302	21,934	1.377	0.894	1.901
8	Essen	585	210	11,387	24,537	1.368	0.892	1.932
9	Düsseldorf	572	218	8,237	16,773	1.380	0.849	1.964
10	Bremen	543	318	10,227	21,702	1.351	0.909	1.931
11	Hanover	517	204	1,589	3,463	–	–	–
12	Duisburg	509	233	6,300	14,333	1.480	0.900	1.924
13	Leipzig	495	293	9,071	21,199	1.320	0.880	1.926
14	Nuremberg	493	187	8,768	18,639	1.420	0.854	1.831
15	Dresden	480	328	9,643	22,307	1.355	0.870	1.892
16	Bochum	389	146	6,970	15,091	1.337	0.847	1.829
17	Wuppertal	364	168	5,681	11,847	1.279	0.881	1.883
18	Bielefeld	325	259	8,259	18,280	1.337	0.872	1.735
19	Bonn	309	141	6,365	13,746	1.374	0.889	2.018
20	Mannheim	309	145	5,819	12,581	1.455	0.897	1.959

Fig. 16.5 Illustration of a strictly hierarchical “arterial” road network capable of connecting all cities, assuming a spatial organization according to Christaller’s theory of central places. Note that not all links are shown here in order to avoid an overloaded picture



that are not organized in a strictly hierarchical manner (in contrast to Fig. 16.5). As will be illustrated for the example of information flows in organizations, strict hierarchies are only optimal under certain conditions, particularly a high reliability of nodes and links.

Experimental results on the problem solving performance of groups (Ulschak, 1981; Tubbs, 2003) show that small groups can find solutions to difficult problems faster than any of their constituting individuals, because groups profit from complementary knowledge and ideas. Small groups also have a potential to assess situations and future developments better than their single members (Chen, Fine, & Huberman, 2003). The actual performance, however, sensitively depends on the organization of information flows, i.e., on who can communicate with whom. If communication is unidirectional, for example, this can reduce performance. However, it may also be inefficient if everybody can talk to everyone else. This is, because the number of potential (bidirectional) communicative links grows like $N(N - 1)/2$, where N denotes the number of group members. As a consequence, the number of information flows explodes with the group size, which may easily overwhelm the communication and information processing capacity of individuals. This explains the slow speed of group decision making, i.e. the inefficiency of committees. It is also responsible for the fact that, after some transient time, (communication) activities in large (discussion) groups often concentrate on a few members only, which reminds of the bundling of traffic flows discussed in the last section. A similar effect is observed in insect societies such as bee hives. When a critical colony size is exceeded, a few members develop hyperactivity, while most colony members become lazy (Gautrais, Theraulaz, Deneubourg, & Anderson, 2002).

These findings indicate that there may be an optimal size of companies and organizations (Huberman & Loch, 1996). Considering the limited communication and information processing capacities of individuals, the optimal number of group members seems to be seven (or less) (Miller, 1956; Baddeley, 1994). This implies the need for bundling and compressing information flows, which is, for example, satisfied by hierarchical organizations. But are there better forms of organization than strictly hierarchical ones? Some of the relevant questions are:

- How robust is the communication or organization network with respect to failure of nodes (due to illness, holidays, quitting the job) or links (due to difficulty personal relationships)?
- How suitable is the organization for crisis management?
- How well does an organization interconnect interrelated activities?
- What is the degree of information loss when communication within an organization network is imperfect?

Similar to road networks and biological supply networks (such as the respiratory system), organizations must be organized space-filling in their covered competence field with staff members playing the role of terminal units. For matters of illustration, we will focus on regular, two-dimensional space-filling kinds of subdivision, as they are particularly suited for a modular organization structure. They share some properties with urban road networks, while the tree-like organization of arterial, water or respiratory supply systems in biological species is three-dimensional (West et al., 1997).

Regular area-filling kinds of subdivision can be triangular, quadratic, or hexagonal. These subdivisions are all compatible with a strictly hierarchical organization, see Fig. 16.6. If the top level consists of one individual (the CoE) and each member of a certain level, except for the lowest one, has N_D subordinates, the number of staff members in a system with L hierarchies is given by

$$N = \sum_{l=1}^L N_D^{l-1} = \frac{N_D^L - 1}{N_D - 1} \quad (16.4)$$

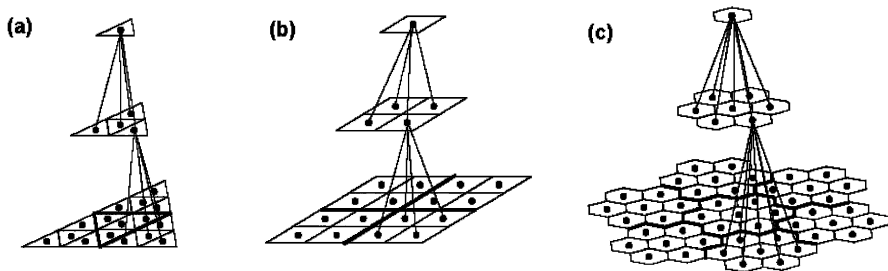


Fig. 16.6 Examples of strict hierarchies based on (a) a triangular, (b) a quadratic, or (c) a hexagonal area-covering organization. Dots represent staff members, while the links indicate the communication pathways (after Helbing, Johansson, et al., 2006)

Table 16.3 Number of members in a hierarchical organization as a function of the hierarchy levels, when everyone (apart from the lowest organizational level) has four or seven subordinates, respectively (after Helbing, Johansson, et al., 2006)

Levels	Four subordinates		Seven subordinates	
	No. of members	Cumulative no.	No. of members	Cumulative no.
1	1	1	1	1
2	4	5	7	8
3	16	21	9	57
4	64	85	343	400
5	256	341	2,401	2,801
6	1,024	1,365	16,807	19,608
7	4,048	5,413	117,649	137,257
8	16,192	21,605	823,543	960,800

(see Table 16.3). While triangular and quadratic structures correspond to $N_D = 4$, subordinates, hexagonal structures are compatible with $N_D = 4, 5, 6$, or 7 (Helbing, Johansson, Mathiesen, Jensen, & Hansen, 2006). As a consequence, for a given number N of individuals, the number of hierarchical levels can be reduced by a hexagonal kind of organization (see Table 16.3). Note, however, that a strictly hierarchical organization of the road system for Christaller's (1980) hexagonal system of central places corresponds to $N_D = 4$ (see Fig. 16.5); otherwise some cities would have multiple access routes.

As the number of hierarchical levels reflects the number of intermediate steps from the bottom level to the top (and vice versa), on the one hand, it is desirable to have a small number of hierarchical levels ("flat hierarchies") to minimize information delays. On the other hand, assuming that the information compression is roughly proportional to the inverse $1/N_D$ of the number N_D of subordinates, flat hierarchies have a higher degree of information loss from one hierarchical level to the next higher one. (This assumes that a fixed amount of communication and information processing capacity is basically divided among the number of subordinate. Given a fixed number N of members of an organization, let us calculate the probability P that certain information from the basis (i.e., the lowest hierarchical level) reaches the top level (or vice versa). Considering that information is compressed by a factor of $1/N_D$ from one level to the other and lost with some probability $p > 0$, we get

$$P = (1 - p)^{L-1} \left(\frac{1}{N_D} \right)^{L-1} = (1 - p)^{L-1} \left(\frac{1}{N_D^{L-1}} \right) \approx (1 - p)^{L-1} \frac{1}{N} \quad (16.5)$$

because of $N \approx N_D^{L-1}$. That is, the larger the number of hierarchy levels, the greater the chance that some potentially relevant information from the bottom level will never reach the top level.

Let us now discuss how the value of P can be improved by redundant information flows. In disaster response management, strictly hierarchical organizations tend to show certain weaknesses with potentially serious consequences:

- Important information is lost due to information compression.
- Information takes too much time to get from the sender to the intended receiver because of too many hierarchical levels to be crossed.
- Information never reaches its destination, because some information node or link does not function.

These weaknesses can be mitigated by additional side links (information flows within the same hierarchy level) and shortcuts between different hierarchy levels (see Fig. 16.7). Since the information flow, over an existing communication link, is basically proportional to $(1 - p)$, redundant links will always reduce the probability that information is lost, as additional information channels are available. However, establishing and maintaining additional information channels is costly, at least it

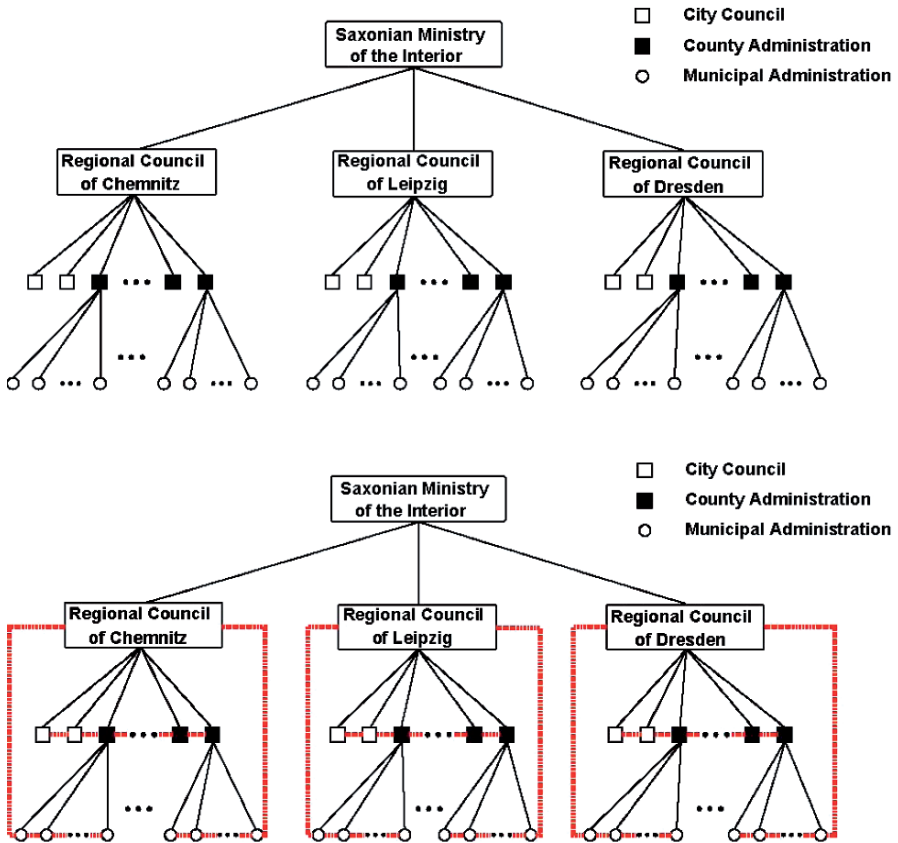


Fig. 16.7 *Top*: Illustration of the hierarchical information flows in the disaster response management during the floodings in Saxony, Germany, in August 2002 (after Helbing, Johansson, et al., 2006). *Bottom*: Improved information flows can be reached by additional side links and shortcuts (dashed thick lines)

requires time. Therefore, their optimum number depends on the reliability of nodes and links. Generally speaking, it increases with the failure rate.

One interesting question is how to establish the most urgently needed links. If some information link breaks down or does not function properly, an alternative information link should be used or established. The same applies, if the capacity of some information channel is not high enough (e.g., due to a lack of communication time or communication ability). That is, information channels should be adaptively strengthened, when needed. This can either be done by extending or redistributing communication times or by establishing additional links.

Regarding the identification of missing links, it is interesting to see how Amazon (www.amazon.com) recommends books to customers, based on their previous purchase decisions and those of other customers. This method is based on an evaluation of correlations among different purchasing activities. A similar method has been recently suggested by Adamic and Adar (2003), who have identified missing links by analysis of e-mail communication. In some sense, it is recommended to establish a new link (a “shortcut”), if it would reduce information flows via many nodes, i.e. if it would reduce detours.

16.4 Spontaneous Self-Organization of Hierarchies

Note that it can be difficult to establish a hierarchical organization. Social systems are complex systems in which the non-linear interactions between its individuals can dominate efforts to control their behavior. However, a hierarchical organization can often emerge by self-organization of its elements. One example is the phenomenon of “crowd turbulence.” Fruin, 1993) reports:

At occupancies of about seven persons per square meter the crowd becomes almost a fluid mass. Shock waves can be propagated through the mass, sufficient to . . . propel them distances of three meters or more. . . . People may be literally lifted out of their shoes, and have clothing torn off. Intense crowd pressures, exacerbated by anxiety, make it difficult to breathe, which may finally cause compressive asphyxia. The heat and the thermal insulation of surrounding bodies cause some to be weakened and faint. Access to those who fall is impossible. Removal of those in distress can only be accomplished by lifting them up and passing them overhead to the exterior of the crowd.

This drastic picture visualizes the conditions in extremely dense crowds quite well, but it does not provide a scientific analysis and interpretation.

Our detailed analysis of video recordings of the pilgrimage in Mecca has now revealed how extremely dense crowds, after a previous transition from laminar flows to stop-and-go waves (Helbing, Ammoser, & Kühnert, 2006), develop a turbulent dynamics characterized by random displacements of pedestrians into all possible directions (see Fig. 16.8a). These displacements are measured as the distance moved between two successive stops of a pedestrian and can reach magnitudes up to 12 m or more (see Fig. 16.8b).

We suggest comparing extreme crowding with driven granular media of high density. Under quasi-static conditions (Radjai & Roux, 2002), these are building up

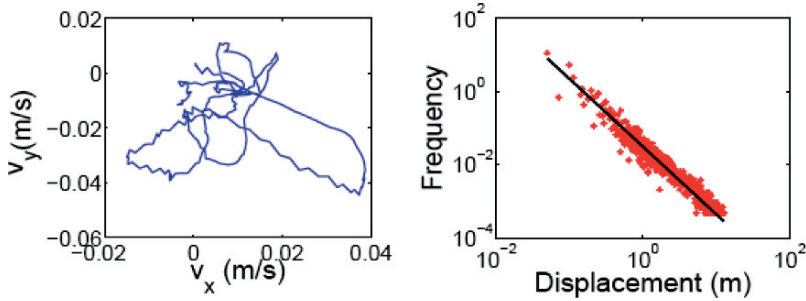


Fig. 16.8 *Left:* Typical time-dependence of both components of velocity in the course of time during turbulent crowd motion. One can clearly see the motion into all possible directions, including the change from forward to backward motion. *Right:* The double-logarithmic representation of the frequency of differently sized displacements between stopping events reveals a power law. (After Helbing, Johansson, & Al-Abideen, (2007))

force chains (Cates, Wittmer, Bouchaud, & Claudin, 1998), which are characterized by strong fluctuations in the strengths and directions of the local forces. As in earthquakes (Bak, Christenson, Danon, & Scanlon, 2002; Johnson & Jiz, 2005), this can lead to events of sudden, uncontrollable stress release with power-law distributed displacements. Such a power-law has, in fact, been discovered by our video-based crowd analysis (Fig. 16.8b). It indicates a self-similar behavior. However, instead of the vortex cascades in turbulent fluids, one observes another kind of hierarchical organization: at extreme densities, individual motion is replaced by collective motion. That is, there are groups of people moving at the same speed. These groups form clusters moving at similar speeds, which again form larger clusters, etc.

Note that the spontaneous formation of hierarchies is quite typical in social systems: individuals form groups, which form companies, organizations, and parties, which make up a society or nation. A similar situation can be found in biology, where organelles form cells, cells form tissues, tissues form organs, and organs form bodies. Another example is well-known from physics, where elementary particles form nuclei, which combine to atoms with electrons. The atoms form chemical molecules, which organize themselves as solids. These make up celestial bodies, which form solar systems, which again establish galaxies.

Obviously, the non-linear interactions between the different elements of the system give rise to a formation of different levels, which are hierarchically ordered one below another. While changes on the lowest hierarchical level are fastest, changes on the highest level are slow.

On the lowest level, we find the strongest interactions among elements. This is obviously the reason for the fast changes on the lowest hierarchical level. If the interactions are attractive, bonds will arise. These cause the elements to behave no longer completely individually, but to form units representing the elements of the next level. Since the attractive interactions are more or less “saturated” by the bonds, the interactions within these units are stronger than the interactions between them. The relatively weak residual interactions between the formed units induce

their relatively slow dynamics. Consequently, a general interdependence between the interaction strength, the changing rate, and the formation of hierarchical levels can be found.

16.5 Summary and Conclusions

In this chapter, we have started with an empirical study of urban supply networks. We have found various power laws: While a linear scaling with the population size was found for the number of doctors or pharmacies in a city, quantities like petrol stations, supermarkets or hospitals scale sublinearly, indicating an economy of scales. Non-material quantities such as information, money, or social interactions, however, scale superlinearly. If these factors determine the speed of urban growth, this implies a finite-time singularity which can only be avoided by friction or innovation (Bettencourt et al., 2007).

Furthermore, we have compared urban systems with biological and ecological systems. Despite some interesting analogies, the differences are significant. For example, there is no strict hierarchical organization of urban transport networks. Nevertheless, we find power-laws for the distribution of traffic flows and the distribution of areas enclosed between urban roads. The power-law exponents are universal, at least for Germany's 20 largest cities.

A deviation from a strictly hierarchical organization is reasonable when the functioning of the nodes or links of a network is not reliable (e.g., due to failures). In such cases, redundant links (side links and shortcuts) increase the robustness of the system. However, the possibilities to design an urban or social system are limited, as their organization and dynamics is, to a large extent, the result of self-organization. Nevertheless, hierarchies result quite naturally based on non-linear interactions among the system elements. As an example, we have discussed the case of "turbulence" in extremely dense pedestrian crowds.

Acknowledgments This study has been partially supported by the EU projects ISCOM and MM-COMNET and the DFG projects He2789/5-1, 6-1, and 7-1. The authors would like to thank Luis Bettencourt, Jose Lobo, Denise Pumain, and Janusz Holyst for inspiring discussions. They are also grateful to the VDEW for the data of the German electricity suppliers. Moreover, they acknowledge support of our data analyses by HE Habib Zein Al-Abideen, Salim Al Bosta, staff by the company Stesa, and the students Markus Winkelmann, Winnie Pohl, Kristin Meier, and Peter Felten. S.L. is grateful for a temporary scholarship by the "Studienstiftung des Deutschen Volkes."

References

- Adamic L. A., & Adar, E. (2003). Friends and neighbors on the web, *Social Networks*, 25(3), 211–230.
- Albert, R., & Barabási, A. -L. (2002). Statistical mechanics of complex networks. *Reviews of Modern Physics*, 74, 47–97.
- Baddeley, A. 1994. The magical number 7 – still magic after all these years. *Psychological Review*, 101(2), 353–356.

- Bak, P., Christensen, K., Danon, L., & Scanlon, T. (2002). Unified scaling law for earthquakes. *Physical Review Letters*, 88, Article number 178501.
- Batty, M., & Longley, P. (1994). *Fractal cities: A geometry of form and function*. London, UK: Academic Press.
- Bettencourt, L. M. A., Lobo, J., Helbing, D., Kühnert, C., & West, G. B. (2007). Growth, innovation, scaling, and the pace of life in cities. *Proceedings of the National Academy of Sciences*, 104(17), 7301–7306.
- Brandes, U., & Erlebach, T. (Eds.). (2005). *Networks analysis*. Berlin, Germany: Springer.
- Brown, J. H., & West, G. B. (2000). *Scaling in biology*. Oxford, UK: Oxford University Press.
- Buhl, J., Gautrais, J., Reeves, N., Solé, R. V., Valverde, S., et al. (2006). Topological patterns in street networks of self-organized urban settlements. *The European Physical Journal B*, 49, 513–522.
- Cates, M. E., Wittmer, J. P., Bouchaud, J. P., & Claudin, P. (1998). Jamming, force chains, and fragile matter. *Physical Review Letters*, 81(9), 1841–1844.
- Chen, K. -Y., Fine, L. R., & Huberman, B. A. (2003). Predicting the future. *Information Systems Frontiers*, 5, 47–61.
- Christaller, W. (1980). *Die zentralen Orte in Süddeutschland* (3rd ed.). Darmstadt, Germany: Wissenschaftliche Buchgesellschaft.
- Costa, L. da F., & da Rocha, L. E. C. (2006). A generalized approach to complex networks. *The European Physical Journal B*, 50, 237–242.
- Crucitti, P., Latora, V., & Porta, S. (2006). Centrality measures in spatial networks of urban streets. *Physical Review E, Statistical, Nonlinear, and Soft Matter Physics*, 73(3) 036125–036129.
- Enquist, B. J., Brown, J. H., & West, G. B. (1998). Allometric scaling of plant energetics and population density. *Nature* 395, 163–165.
- Frankhauser, P. (1994). *La fractalité des structures urbaines*. Paris, France: Anthropos.
- Fruin, J. J. (1993). The causes and prevention of crowd disasters. In R. A. Smith, & J. F. Dickie (Eds.), *Engineering for crowd safety* (pp 99–108). Amsterdam, The Netherlands: Elsevier.
- Gabaix, X. (1999). Zipf's law for cities: An explanation. *Quarterly Journal of Economics*, 114(3), 739–767.
- Gastner, M., & Newman, M. (2006). The spatial structure of networks. *The European Physical Journal B*, 49, 247–252.
- Gautrais, J., Theraulaz, G., Deneubourg, J. -L., & Anderson, C. (2002). Emergent polyethism as a consequence of increased colony size in insect societies. *Journal of Theoretical Biology*, 215, 363–373.
- Gibrat, R. (1931). *Les Inégalités Economiques*. Paris, France: Librairie du Recueil Sirey.
- Helbing, D., Ammoser, H., & Kühnert, C. (2006) Information flows in hierarchical networks and the capability of organizations to successfully respond to failures, crises, and disasters. *Physica A – Statistical Mechanics and its Applications*, 363(1), 141–150.
- Helbing, D., Johansson, A., Mathiesen, J., Jensen, M. H., & Hansen, A. (2006). Analytical approach to continuous and intermittent bottleneck flows. *Physical Review Letters*, 97, Article number 168001.
- Helbing, D., Johansson, A., & Al-Abideen, H. Z. (2007). The dynamics of crowd disasters: an empirical study. *Physical Review E*, 75, Article number 046109, part 2.
- Holyst, J. A., Sienkiewicz, J., Fronczak, A., & Suchecki, K. (2005). Scaling of distances in correlated complex networks. *Physica A – Statistical Mechanics and its Applications*, 351, 167–174.
- Huberman, B. A., & Loch, C. H. (1996). Collaboration, motivation and the size of organizations. *Journal of Organizational Computing*, 6, 109–130.
- Jiang, B., & Claramunt, C. (2004). A structural approach to the model generalization of an urban street network. *Geoinformatica*, 8(2), 157–171.
- Johnson, P. A., & Jiz, X. (2005). Nonlinear dynamics, granular media and dynamic earthquake triggering. *Nature*, 437, 871–874.
- Kölbl, R., & Helbing, D. (2003). Energy laws in human travel behaviour. *New Journal of Physics*, 5, Article number 48.

- Kühnert, C., Helbing, D., & West, G. B. (2006). Scaling laws in urban supply networks. *Physica A – Statistical Mechanics and its Applications*, 363, 96–103.
- Lämmer, S., Gehlsen, B., & Helbing, D. (2006). Scaling laws in the spatial structure of urban road networks. *Physica A – Statistical Mechanics and its Applications*, 363, 89–95.
- Levinson, D., & Yerra, B. (2006). Self-organization of surface transportation networks. *Transportation Science* 40(2), 179–188.
- Makse, H.A., Havlin, S., & Stanley, H.E. (1995). Modeling urban-growth patterns. *Nature*, 377, 608–612.
- Miller, G. A. (1956). The magical number 7, plus or minus 2 – some limits on our capacity for processing information. *The Psychological Review*, 63, 81–97.
- Newman, M. E. J. (2002). Assortative mixing in networks. *Physical Review Letters*, 89, Article number 208701.
- Porta, S., Crucitti, P., & Latora, V. (2006). The network analysis of urban streets: a primal approach. *Environment and Planning B – Planning and Design*, 33(5), 705–725.
- Pumain, D., Paulus, F., Vacchiani, C., & Lobo, J. (2006). An evolutionary theory for interpreting urban scaling laws. *Cybergeo*, 343, 20p.
- Radjai, F., & Roux, S. (2002). Turbulentlike fluctuations in quasistatic flow of granular media. *Physical Review Letters*, 89(6), Article number 064302.
- Schweitzer, F. (Ed.). (1997). *Self-organization of complex structures: From individual to collective dynamics*. London, UK: Gordon and Breach.
- Schweitzer, F. (2003). *Brownian agents and active particles*. Berlin, Germany: Springer.
- Simon, H. (1955). On a class of skew distribution functions. *Biometrika*, 42(3–4), 425–440.
- Steindl, J. (1965). *Random processes and the growth of firms*. New York, NY: Hafner.
- Strumsky, D., Lobo, J., & Fleming, L. (2005). Metropolitan patenting, inventor agglomeration and social networks: a tale of two effects. *Santa Fe Institute Working Paper*, 05-02-004.
- Sutton, J. (1997). Gibrat's legacy. *Journal of Economics Literature*, 35(1), 40–59.
- Tubbs, S. L. (2003). *A systems approach to small group interaction*. Boston, MA: McGraw-Hill.
- Ulschak, F. L. (1981). *Small group problem solving: An aid to organizational effectiveness*. Cambridge, MA: Addison-Wesley.
- West, G. B., Brown, J. H., & Enquist, B. J. (1997). A general model for the origin of allometric scaling laws in biology. *Science*, 276, 122–126.
- Zipf, G. K. (1949). *Human behaviour and the principle of least effort: An introduction to human ecology*. (1st ed.). Cambridge, MA: Addison-Wesley.

Chapter 17

Using Statistical Physics to Understand Relational Space: A Case Study from Mediterranean Prehistory

Tim Evans, Carl Knappett and Ray Rivers

17.1 Introduction

Spatial relationships among entities across a range of scales are fundamental in many of the social sciences, not least in human geography, physical geography, and archaeology. However, space has received a surprisingly uneven treatment; in archaeology, for example, spatial analysis only really came to the fore in the 1960s and 70s, through the influence of the ‘New Geography’ (Haggett, 1965; Chorley & Haggett, 1967). David Clarke (1968, 1977), one of the principal exponents of spatial analysis in New Archaeology, described three levels of resolution in spatial archaeology: the micro level, the semi-micro or meso level and the macro level, the last of these representing relationships between sites (Clarke, 1977, p. 13). Yet, despite the clear implication that these levels should articulate, many subsequent studies have tended to aim at just one level. World-systems theory, for example, forms the basis for core-periphery models that examine macro-level spatial relationships (e.g., Schortman & Urban, 1992; Peregrine, 1996; McGuire, 1996; Chase-Dunn & Hall, 1997; Stein, 1998; Kardulias, 1999).

While Clarke’s emphasis on different spatial scales has the advantage of clarity, his general approach, and indeed that of much spatial analysis of this kind, has been criticized for its overly deterministic approach to space. The idea that space has absolute geometric properties has been increasingly challenged by scholars arguing that space is a relative construct, a process that emerges out of social practices. In geography this critique already has a long history (e.g., Harvey, 1973), which has been taken up by increasingly diverse and influential voices (e.g., Lefebvre, 1991; Harvey, 1996; Soja, 1996; Thrift, 1996; Hetherington, 1997; Murdoch, 2006). This ‘spatial turn’ has also been experienced in archaeology, where approaches to space were ‘relationalized’ through the influence of phenomenology in landscape studies (Bender, 1993; Tilley, 1994; Knapp & Ashmore, 1999; Smith, 2003; Blake, 2002, 2004).

T. Evans (✉)
Department of Physics, Imperial College London, London, UK

In both disciplines, however, and perhaps across the social sciences more broadly, the move towards relational conceptions of space and away from geometric determinism has arguably created a dualism between relational and physical space. It is our aim, in this paper, to develop a methodology that can go some way toward bridging the gap that has opened up between them (cf. Hillier, 2005). This requires an approach that incorporates the fundamental notion that humans create space through social practices, while also acknowledging that physical parameters are not entirely redundant in this process. One of the misconceptions hindering this rapprochement has been that spatial analysis is deemed to be bound to Euclidean geometry. Yet, recent advances in complexity science and in the study of complex networks in particular, give the lie to this idea (e.g., Batty, 2005, on networks in geography). These advances also allow for the evolution of spatial dynamics from the bottom-up, in ways seemingly unimaginable to central place theory or core-periphery models. Finally, it fulfils the ideal of being able to link together different scales, such as the micro, meso and macro levels mentioned above in relation to the work of Clarke (1977).

While complexity science has certainly had a major influence on our approach, we believe that some of the problems with spatial analysis can actually be worked through at a more basic level. A fundamental problem is that in much spatial analysis, even in the more sophisticated forms of GIS, interactions between points are seen as secondary to the existence of those points. Batty (2005, p. 149) has described this as ‘the geography of locations, not relations.’ The same criticism has been leveled in very similar terms by Doel, who bemoans the fact that ‘in geography, the fundamental illusion is the autonomy and primacy of the point’ (Doel, 1999, p. 32). The equivalent to this, in the archaeological analysis of regional systems, is that the sites are thought to emerge and gain their character on largely local grounds, and any interactions with other communities in the region follow on from that. The connections between sites are simply drawn as lines, without weight or orientation. Such ‘site-centrism’ makes it difficult to entertain the thought that between-site interactions might themselves have contributed to the size and status of the sites in question.

How, then, might we redress the balance and look at what happened from the perspective of interactions as a step on the way take account of both sites and their interactions across different scales. With this in mind, how can we achieve what we might dub, borrowing from Batty, an ‘archaeology of relations?’ This question is precisely what we will now explore for our case study area, the Bronze Age Aegean. We have chosen this area partly because of the specialization of one of the authors (CK), but also because there already is an excellent example of just the kind of interactionist perspective that we seek to develop: the work of Broodbank on the Early Bronze Age Cyclades (Broodbank, 2000). It is the only systematic attempt thus far, for any period of the prehistoric Aegean, to explain the growth of certain sites in terms of their interactions. This approach perhaps was encouraged by the fact that some important sites in the area and period in question – the Cycladic islands during the Early Bronze Age (c. 3000–2000 BC) – are very hard to explain in terms of local resources, occurring on small rocky islands with limited agricultural

or mineral resources. Indeed, some are only inhabited for the first time in the Late Neolithic, in contrast to the more agriculturally viable mainland and larger islands (such as Crete), where settlement stretches back to the Early Neolithic. Thus, it seemed likely that relative regional location had a substantial role to play in a site's importance.

17.2 From EBA to MBA Networks

It should be emphasized that our involvement in this project began not with the Early Bronze Age (EBA) Cyclades, but with a larger area in a later period – the whole of the southern Aegean in the later Middle Bronze Age (MBA). This period is well bounded in time, as the record shows significant gaps at its temporal boundaries. Furthermore, the sail appears c. 2000 BC, which facilitates the study of the transition towards new levels of inter-regional interaction and a metamorphosis in the character of regional exchange networks.

One key question we wanted to answer was why some sites, such as Knossos on Crete, grew to be so large and influential. The size of these sites is usually explained in local terms of surplus and growth, which enabled exchange with other sites. We were interested in reversing this equation, exploring the possibility that some characteristics of the larger interaction networks contributed to the growth of the sites (see also Rihl & Wilson, 1991, in relation to the growth of the city-states of Geometric Greece, c. 800–700 BC).

But before we can do this, we need to step back, and ask: 'what are the fundamental characteristics of the EBA Cycladic network?', and 'how might the more complex networks of the MBA differ?' The principle behind Broodbank's treatment of the EBA Cyclades as a network is straightforward, taking sites as the vertices and their connections as the edges, and thus transforming the Cyclades into a simple mathematical graph. He then adopts a basic technique from graph theory¹ known as 'Proximal Point Analysis' (PPA), already used effectively in archaeology and anthropology for interaction studies in other archipelagos, notably in Oceania (Terrell, 1977; Irwin, 1983; Hage & Harary, 1991, 1996). In PPA, edges are drawn from each hypothetical site to its three nearest neighbors in geographical space. Some sites emerge as more connected than others, with five or six edges to other sites. These sites possess greater 'centrality' in the network, meaning that they might be expected to have a more prominent role in regional interactions. When certain parameters such as site density are altered, simulating population increase over time, the texture of the network changes and other sites can emerge as central. When Broodbank compared the results of his PPA with the archaeological data, he found that it did indeed predict that a site on Keros, for example, would possess centrality in such a network. Of the five major Early Cycladic sites, three were 'central' in the PPA.

¹ See Evans (2005) for a review of basic graph theory and bibliography of exemplary applications in a variety of fields.

Of course, Broodbank also had to suggest some motivation for these interactions – communities do not just interact without motives or goals. The EBA Cyclades are agriculturally marginal and not self-sufficient and he cited basic demographic processes and the need for social storage networks (Broodbank, 2000, pp. 81–96), with power and prestige emerging consequentially out of network interactions.

Whilst inspired by Broodbank's approach, we realized that his networks, while appropriate for the EBA Cyclades, could not be translated to the MBA Aegean for reasons that can be summarized as geographical, technological and organizational.

17.2.1 'Sails' Change Behavioral Scales

We are familiar with the fact that changes in transport technology, which expand the distances over which individuals can easily travel, lead to new behavioral patterns. In our case, there appear to have been substantial changes in transport technology between the EBA and the later MBA, i.e. the advent of the sail c. 2000 BC as a replacement or supplement for rowing technology. As a result, the distances traveled could easily increase by an order of magnitude. Even if such long trips were still not the norm, sail technology may have made them just significant enough that they could form the basis of important, if weak, links in the sense of Granovetter (1973, 1983). Whereas Broodbank has argued that the EBA Cyclades can form a consistent network within themselves (with some external linkages to the mainland) as in Fig. 17.1, the introduction of the sail renders feasible the MBA interaction networks spanning the whole Aegean (Fig. 17.2), including not only the Cyclades but also the Dodecanese, Crete, and the landmasses of Asia Minor and mainland Greece.

In the EBA Cyclades it was plausible for Broodbank to allocate equal site size in his analysis, and also equal connections, in terms of weight and directionality. When we come to the interaction networks appropriate to the Aegean of Fig. 17.2, which emerge in the late MBA, a very different picture confronts us:

1. Vertices: we know that there are sites of substantially differing sizes and roles, quite unlike the situation in the EBA. Note the assumption that large sites developed due to local internal processes (e.g. access to agricultural surplus).
2. Edges: we can also see that there are very different kinds of links existing simultaneously, varying in orientation, length and weight.

With this scalar change, the main dynamic to concern us is the emergence of 'Minoanization' at the end of the MBA (see Broodbank, 2004 for a recent review). In this process, a number of sites across the south Aegean, on both islands and mainland, developed increasingly complex exchange links and shared cultural traits. The driving force behind this was the large island of Crete, on which certain central sites, Knossos in particular, seem most involved. The similarities in material culture between sites on and off Crete are so pronounced that some have been led to speak of

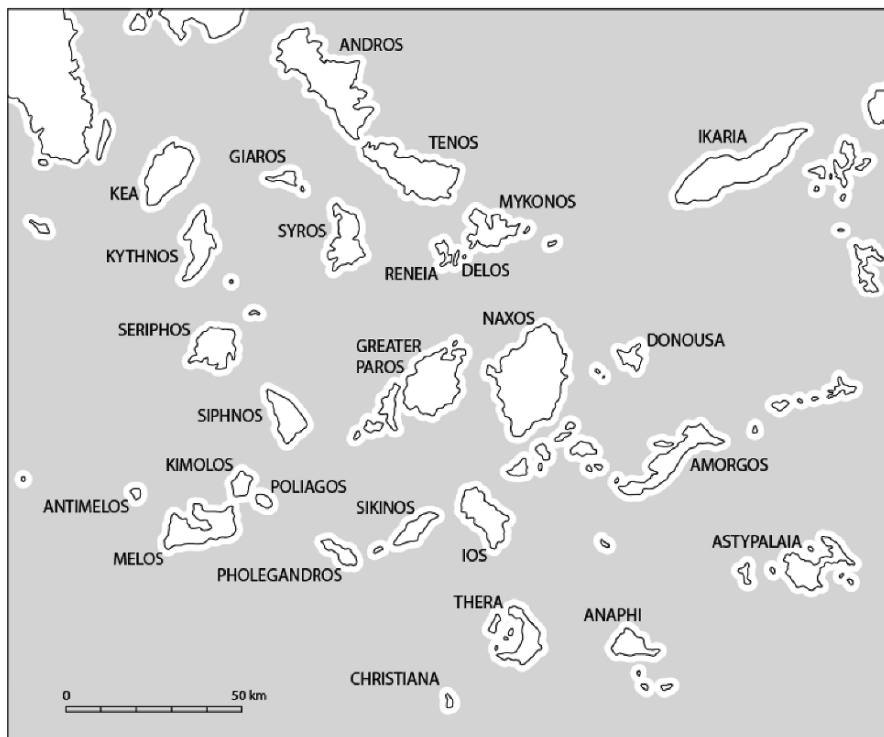


Fig. 17.1 The Cyclades

colonization. This interpretation is connected with the idea of a Minoan sea-empire ('thalassocracy'), but there is no direct evidence that the fleet needed to maintain such an empire actually existed. The source of the thalassocracy idea can actually be traced back to Thucydides, who was of course commenting more than 1000 years later than the period described. A thalassocracy may never have existed and, even if there were some form of 'empire' centered on Crete, it may have been established through economic 'emporía' rather than political colonies.

Whether through direct colonization or indirect acculturation, the Cretan palaces capitalized on their regional dominance and extended their influence beyond the island. Essentially, this represents the earliest ever occurrence of 'state-led' expansionism (because Knossos was probably the center of a regional polity or early state) in the prehistoric Aegean (and, by extension, Europe). Present interpretations are, however, inadequate at many levels, not least the general tendency to explain first the growth of individual sites in local terms (good land, resources, etc.), and then to extrapolate connections between sites from there. In other words, the 'vertices' (sites) always precede the 'edges' (links). There are, naturally, some exceptions to this, with Davis' work on the 'Western String' route through the Cyclades linking Crete to the mainland (Davis, 1979), and Berg's assessment, using world-systems theory, of Southern Aegean interactions in the Middle to early Late Bronze Age

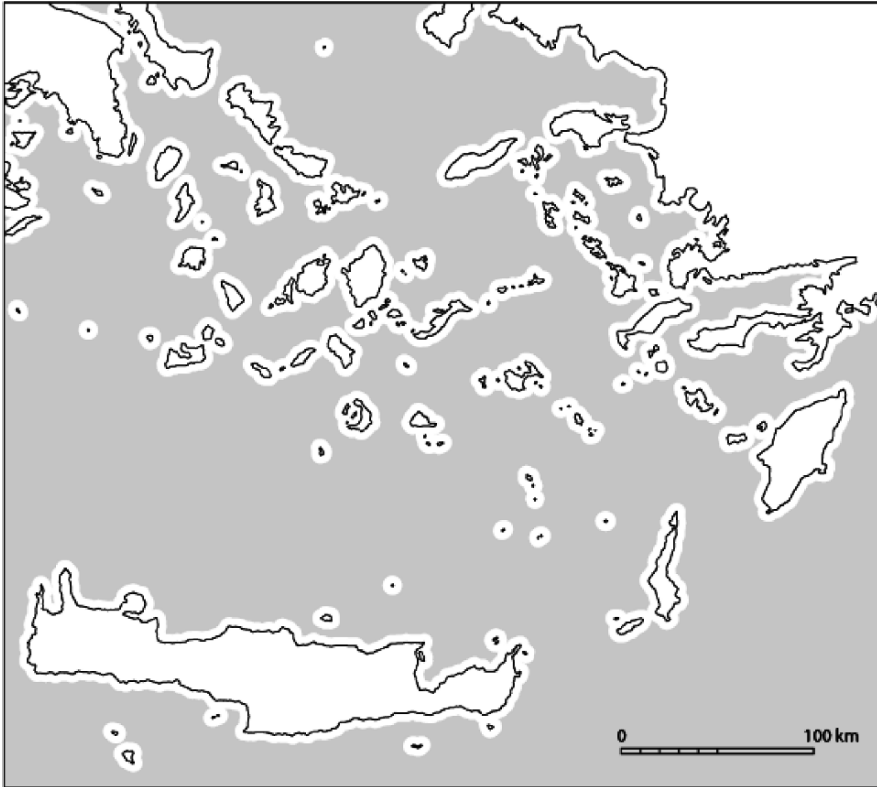


Fig. 17.2 The Aegean

(Berg, 1999). However, these and other studies, while focusing on interactions, have tended not to use explicit network models composed of nodes and links (notably, in these cases, the nodes are undefined).

The two key aspects of these Minoan networks are:

1. An evolution from exchange to affiliation – initially the connections between islands involve exchange of goods, but eventually these are supplemented by actual imitation of artifact styles and technologies, suggestive of some process of cultural affiliation but not necessarily colonization (Knappett & Nikolakopoulou, 2005). It is interesting that this latter process appears to correspond in time with the probable emergence of a single political center on Crete – i.e. Knossos (although this is debated – see Adams, 2006). This central site can be regarded as a hub, and, so, we intend to investigate the possible link between hubs and strong ties in networks of this kind.
2. A relatively rapid emergence and collapse. These interaction networks only endure for a mere two hundred years or so. The EBA networks are followed by a gap in occupation at many sites, if not total abandonment; some of the

most important ‘vertices.’ such as Chalandriani on Syros or Dhaskaleio-Kavos on Keros, are never again occupied (Broodbank, 2000). This does not indicate a particularly resilient system. Furthermore, the MBA network is followed by others in the Late Bronze Age (LBA), based on the Mycenaean mainland, especially the Argolid (Mycenae). Like the Minoan networks, each only lasts a few hundred years, ending cataclysmically (in the sense that there is significant population decline, abandonment of many settlements, and a major decrease in exchange contacts with the east Mediterranean) with the onset of the so-called ‘Dark Ages.’

With these observations in mind, we shall argue that the way in which the meso-level of *intra-island* site activity is accommodated in the macro-level of the *inter-island* network as a whole plays a crucial role in how the whole system functions. This has profound consequences for the way in which we make and understand our models.

17.3 Incorporating the Archaeological Record in Network Models

17.3.1 Meso- and Macro-Levels

The main characteristics of the EBA Cyclades are that they are agriculturally marginal, with small populations that are not self-sufficient. Broodbank’s model is one of punctuated exogenous evolution, in which population growth within an island leads to new communities budding off from the old, maintaining an approximately common size. See Fig. 17.3, taken from Broodbank, in which the increasing number of vertices corresponds to an increasing population.

In Fig. 17.3, each vertex corresponds to a definite population/unit of resource (e.g. 50 people per site at the site density shown is commensurate with the estimated total Cycladic populations in the EBA (Broodbank, 2000)). On attaching each vertex to its three nearest neighbors, we see that, as population increases, the islands become more self-sufficient, and contacts between them become less necessary. In this model for the EBA Cyclades, it is plainly the meso-level that drives the macro-level.

As we have already indicated, the situation for the MBA southern Aegean is very different. Geographically, our area of intended study of Fig. 17.2 is highly heterogeneous, with many small islands and some large islands, such as Crete and Rhodes, and areas of mainland. In contrast, Broodbank’s study area (Fig. 17.1) consisted solely of islands of roughly equivalent size. In fact, if we were to lay a regular grid over the network of the EBA Cyclades with grid size that of the largest islands, we would see that it can, in fact, be conceived as a lattice whose symmetry has been distorted. Although the islands are irregularly spaced, with some squares in the grid more heavily populated and others less so, the edge length does not vary that strongly and very few squares are empty. This means that clustering is unlikely to be all that pronounced. Without being quantitative, we see, in Fig. 17.2, that the

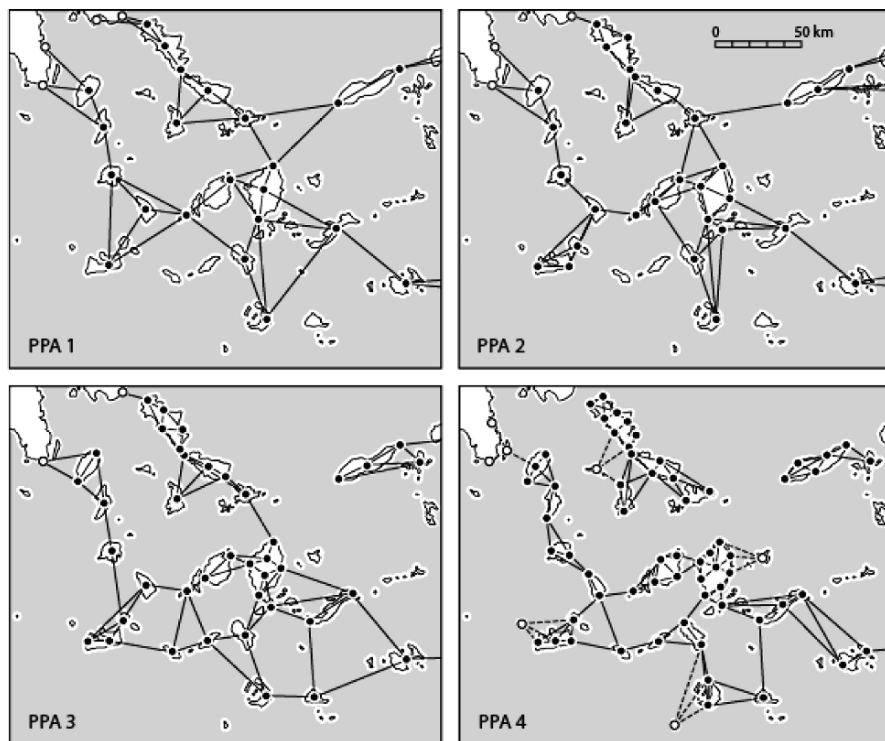


Fig. 17.3 Broodbank's PPA for the Cycladic PPA (after Broodbank, 2000, Fig. 53)

topology is quite different when we increase the scale of the network to include Crete, the Dodecanese, and parts of the Greek mainland and Asia Minor. There are many more open squares if a lattice of the same size is overlain on this space, and there is much greater scope for clustering. In the earlier Cycladic cluster, one island can connect with any other island through a series of relatively short hops. This is not the case for the larger network, in which some long-distance edges are unavoidable if the network is to remain fully connected.

Returning to the larger southern Aegean network, what might it take to overcome the clustering that comes with such an asymmetrical topology? Surely these clusters will remain separated unless there is a strong incentive to connect. In a patchy resource environment there might very well be such an incentive, with interconnectivity providing a safety net against annual fluctuations in resource provision from region to region. Another parameter to consider is the expense of maintaining a long-distance edge. Presumably, a large site with more resources has a greater chance of maintaining such an edge than does a small site. And this leads us to the realization that, if a network is indeed created over a large asymmetrical grid of this kind, then large sites are likely to feature. Such sites are likely to be islands (or their mainland equivalents) treated as single entities, rather than the meso-scale

communities that make up their populations. Furthermore, large sites searching for information about resource availability are much more likely to target to other large sites in that quest, reinforcing the dominance of the inter-island over the intra-island links. It is now the macro-level that appears to drive the meso-level.

17.3.2 The Role of the Archaeological Record

The different emphases between the meso- and the macro-scales in the EBA and MBA require different approaches to the archaeological record. There are two different aspects to the record, one of which belongs to the *input* of our models and one to the *output*. Ideally, we would like to predict where sites should be and what their importance is, but we are unable to do both. The best that we can do is to treat the location of the network nodes as input, and then use the output of our model to determine their relative significance (e.g. population size). This output can then be checked against settlement size as a proxy for population size.

Even a small island will have several significant meso-level communities, only some of which will correspond to archaeologically significant sites, because of the incompleteness of the record. The question is, does it matter how the population is distributed between the sites, both known and unknown archaeologically, on a given island? In the light of the above, this is related to how many nodes we allocated to the island *ab initio*. If it does matter, then, given our lack of knowledge about them, our conclusions can depend on unknown data to the extent that the usefulness of the model is limited to very general statements.

Broodbank's EBA model highlights this dramatically. As the network jumps from one pattern of sites to the next, both their number and positions change. The sites are not chosen through direct archaeological evidence, but are assigned hypothetically in a simple geometrical way, more or less equidistantly throughout each island on the basis of population estimates derived from archaeological survey data. As such, the vertices serve as a proxy for the archaeological record, but any attempts to relate them to significant meso-level sites are doomed to failure on two counts. First, as we have said, not all significant sites are known. Secondly, even when we have reasonable knowledge of sites, the geometric algorithm is too rigid to replicate them, as we see from Fig. 17.4, also taken from Broodbank, where significant EBA sites are displayed. Inevitably, the same potential problem of an incomplete archaeological record also applies to our MBA models. However, our matching the record is different in that, insofar as the macro-scale dominates the meso-scale, we integrate over the meso-level communities to replace them by macro-level aggregates.

17.3.3 Robustness

This MBA characterization of aggregating meso-level sites seems to show an as much indifference to the details of the archaeological record as do Broodbank's proxy sites. However, as far as the models are concerned, the stability of the

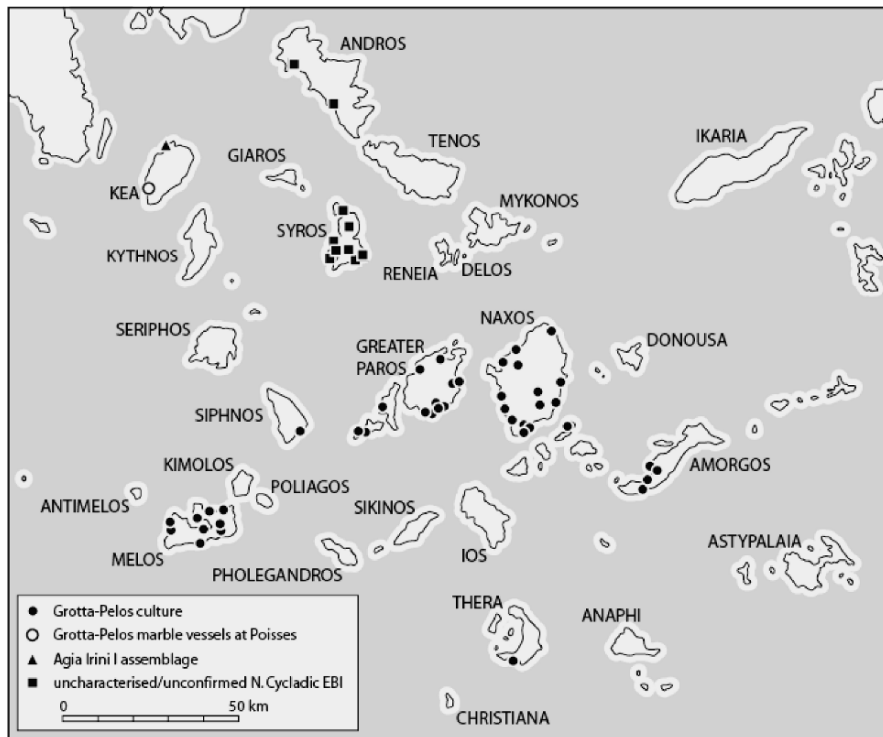


Fig. 17.4 Various EBA sites (after Broodbank, 2000, Fig. 43)

conclusions that we would draw is very different. This is an important issue, and before we present our model we need to think about model-making, in general.

We wish to construct models that are robust. There are several ways to define robustness, but by it we mean that models should have the following two attributes. The first is that they should minimize the effects of our ignorance about the archaeological record as *input*. That is, for a *given model*, different assumptions as to the distribution of sites should have little effect. The second requirement is that, for the *same input*, different models which are alike should have *similar outputs*.

As to the first, Broodbank's algorithm is not robust. As a simple example, we can accommodate increasing population by increasing community size proportionally, without increasing the number of sites. In this way, there is no evolution, and the network is static. This is not to disparage Broodbank's groundbreaking work, since the importance of inter-island links is an essential ingredient in the development of intra-island communities, and vice-versa. However, in our attempts to develop more sophisticated measures of influence (e.g. frequency of interactions, cultural transmission, in-betweenness, etc.) we need robustness if we are to make progress.

The requirement that related models have related outputs is not as simple as it looks since our models, unlike Broodbank's, are not deterministic. The same model, with the same initial conditions, may give different outcomes. Nonetheless, we

assume that these different outcomes share generic features so that we can talk about the typical outcome of a given model with specified initial conditions. Our second requirement for robustness is then primarily a question of morphology. Rather than treat all models as individual, we assume that we can divide them into sets, each set having its own characteristic behavior. It is then a question of determining which set, or *universality class*, gives a better description.

These two requirements for robustness have very different implications and impose strong constraints on our models. For the case in point, we choose to minimize the effects of our ignorance by insisting that the meso-level does not determine the macro-level for the MBA Aegean. Given our earlier discussion of the MBA period, this assumption is not unreasonable. Then, *as a consequence*, we shall see the tendency for like to seek out like. To explain this further, consider two maps of the southern Aegean, one at a large enough scale to show all major (known) sites of MBA habitation, from which Fig. 17.6 (showing Naxos and Mykonos) is a part. Known, or hypothesized, sites are listed but, by definition, unknown sites cannot be. However, the carrying capacity (overall resources available to the island's inhabitants) can be estimated, and used as a proxy for the probability that other, unknown, sites exist on an island.

If the primary dynamic is the affiliation between islands and major centers, and not the detailed interactions within each island, this suggests that we can replace all the individual sites of Naxos in Fig. 17.5 by the single effective supersite (represented by the single vertex of the large circle) of Fig. 17.6, whose attributes are the sum of those of the individual sites. When we subsequently attach a single coordinate to this supersite, we adopt, in effect, a 'center-of-mass' approach. Such aggregation can be applied island by island, or local region of influence by local region.

In the large scale, this approach then permits us to ignore the mesoscopic details and to just incorporate the supersites, one per small island, and more for Crete and the mainland. This process is often performed in physics, where it is known as 'coarse graining.'

The underlying assumption, that the whole is equal to the sum of its parts, is not the only one. The assumption that affiliation can be viewed as an inter-island process not only suggests that local sites can be aggregated into supersites, but the multiple links between individual sites can be replaced by superlinks between supersites (see Fig. 17.7). That is, we would infer the same relative affiliation strengths irrespective of the scale of the map. To pursue the analogies with Newtonian mechanics further, this *forces* us to adopt 'gravitational' models, since the attribute of gravitational energy is that it is the same, whether it is calculated from the centers of mass, or from the individual constituents of those masses.

Of course, for this to be the case *exactly*, we need a definite power-law behavior describing the fall-off of the linkage-'potential' with distance. We cannot justify such a specific requirement, as there will be some dependence on the cost of a community breaking into two. Nonetheless, on changing the scale of the map to a larger scale (coarse-graining) so that the meso-level detail becomes indistinct as it merges into a macroscopic network across the Aegean, the broad patterns of influence remain the same. When we turn to model-making that incorporates both

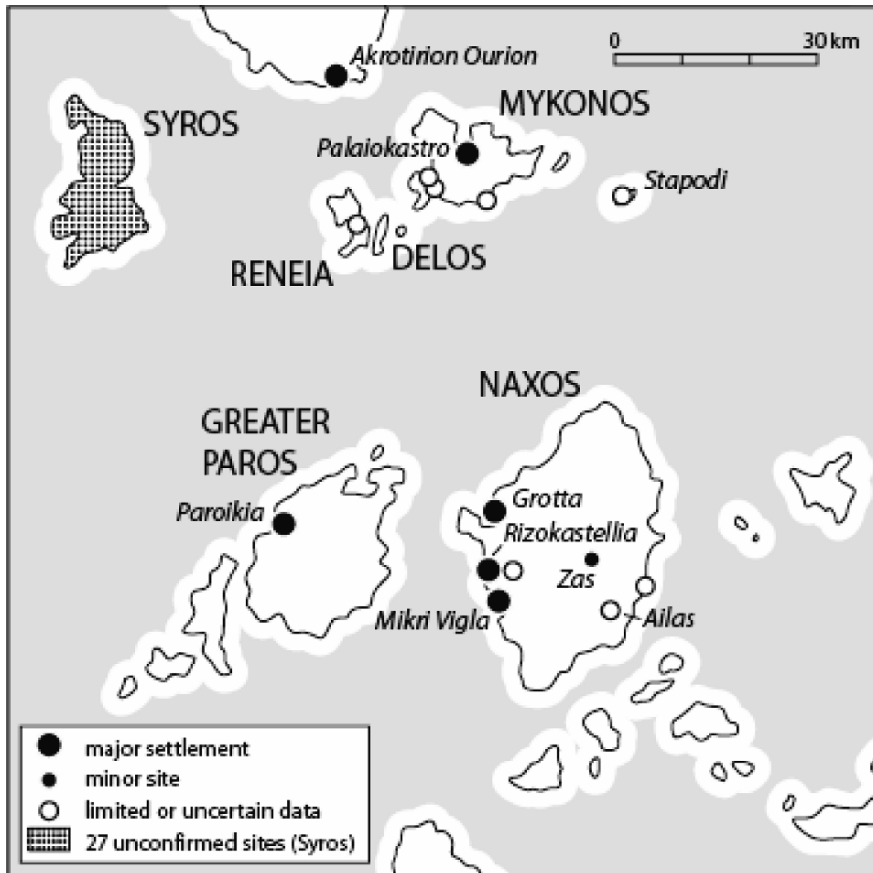


Fig. 17.5 Some MBA sites (after Broodbank, 2000, Fig. 109)

the center-of-mass and the gravitational effects, we shall see that this enthusiasm for like-to-like in establishing links has general implications, such as a tendency to instability in the networks.

Henceforth, all our sites are supersites and our links are superlinks, and we shall revert to calling them ‘sites’ and ‘links’ respectively. With the positions of the vertices determined by ‘center-of-mass,’ they are not directly correlated to a single site from the archaeological record, just as the vertices of Broodbank’s networks were not, although in each case they are informed by it.

17.4 Model-Making for the MBA Aegean

While there are models within social network analysis that use graph theory in increasingly complex ways (Carrington, Scott, & Wasserman, 2005; de Nooy, Mrvar, & Batagelj, 2005), we have decided to go a step further and combine some of these insights with techniques from statistical physics. Statistical physics shows how large

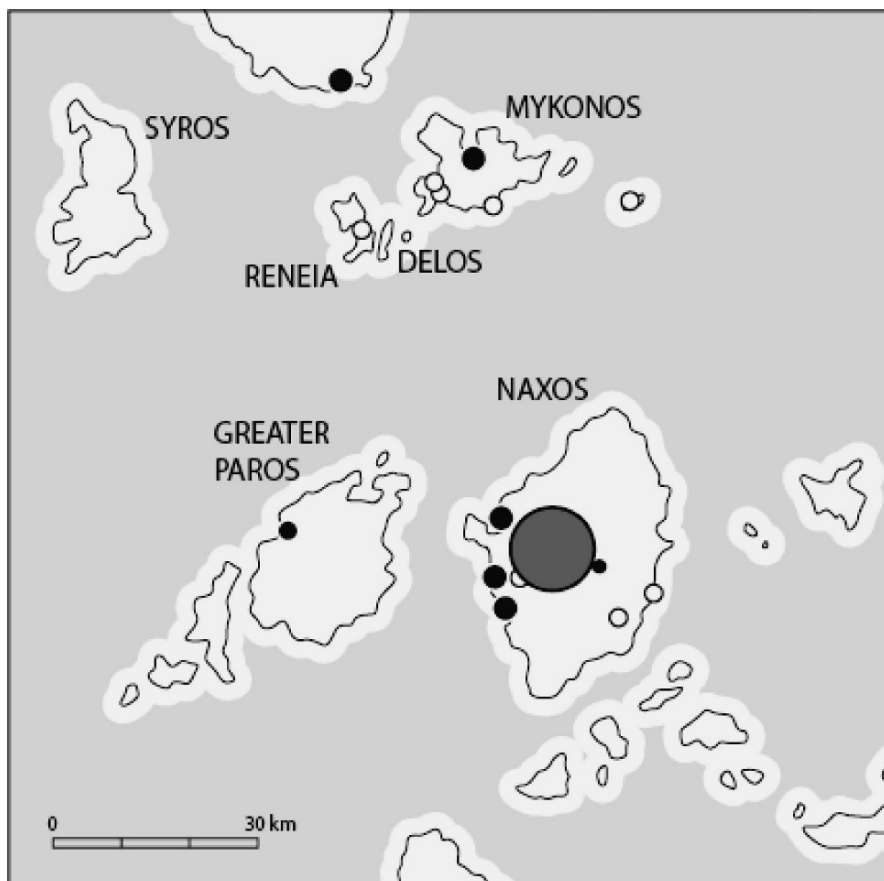


Fig. 17.6 A supersite on Naxos: There is some ambiguity in the positioning of this vertex, when its spatial coordinates are needed later, but if a natural harbor exists this provides an obvious bias (after Broodbank, 2000, Fig. 109)

numbers of interacting entities often have relatively simple generic behavior on large scales regardless of the details of their interactions. Network theory shows how specific behavior is embedded within this. This approach can help us develop an explicit focus on the dynamics of interaction in complex networks and on the interface in such networks between local and global behaviors. Such a recourse to ‘harder science’ than that of Broodbank may seem reminiscent of the (ultimately unsuccessful) application of systems theory to archaeology in the late ‘60’s and early ‘70’s, in which order was imposed top-down, and the ensuing models therefore were unable to take into account that small changes in one subsystem may lead to substantial change at the overall system level. However, while the ‘70’s endeavor was hamstrung by a number of factors, not least the tendency to prescribe the character of the subsystems and their interactions mechanistically, we believe that the new generation of network analysis enables us to conceive of order emerging from the bottom-up, in a far more fluid and contingent manner.

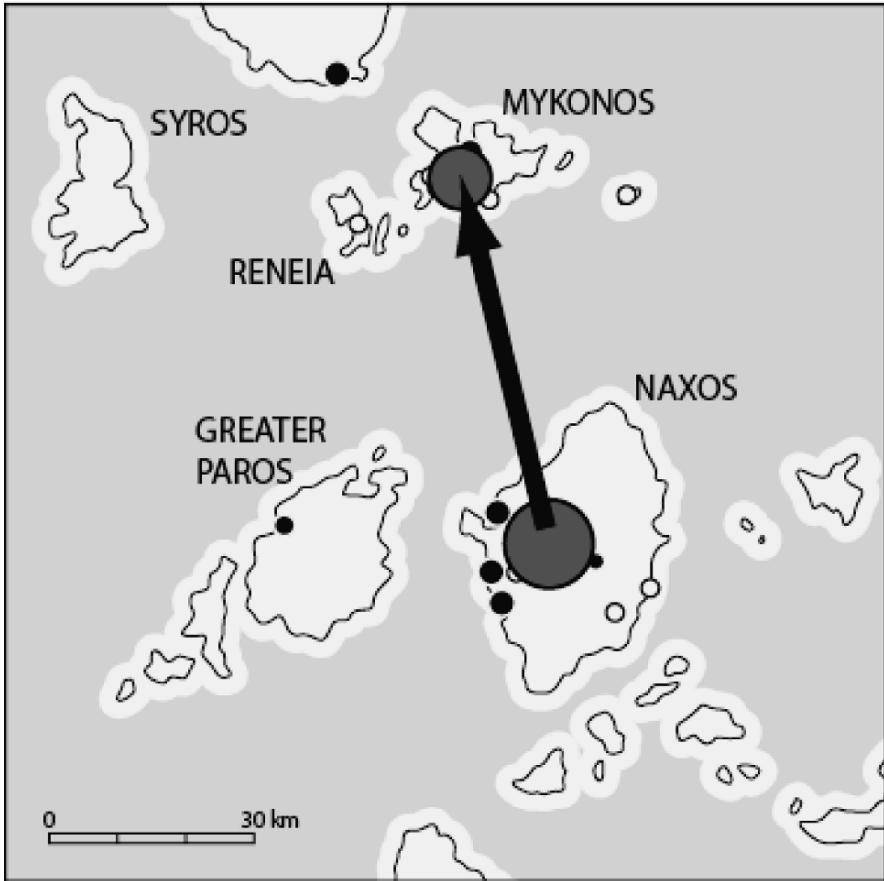


Fig. 17.7 Two supersites connected with a single superlink (after Broodbank, 2000, Fig. 109)

To be specific, our analysis is performed on 34 site vertices in the southern Aegean, as shown in Fig. 17.8. The overall approach we have chosen to adopt portrays the MBA network based on these vertices, with its complicated constraints and interactions, as explicable in terms of an ‘energy landscape’ through which the system moves (hence conceiving of the system as having agency or behavior of some kind). In a contemporary context we would think of this as a ‘cost-benefit’ analysis in which the ‘energy’ or ‘benefit’ is a function H of the state of the network. A system with low energy H is close to some optimal solution in which all the different constraints and interactions are balanced.²

² In his seminal work on spatial analysis in archaeology, Clarke identifies four general theories underlying most of the detailed spatial studies in archaeology that have attempted to move beyond description. These are anthropological spatial theory; economic spatial theory; social physics theory; and statistical mechanics theory (Clarke 1977, p. 18). Considering just the last two for



Fig. 17.8 The Aegean with the location of the 34 sites used for initial investigations. Site 1 is Knossos, site 2 is Malia, site 10 is on Thera, and site 27 is Mycenae. The key to the remaining sites is in the table below. In the following figures all the sites are given equal weight ($S_i = 1$) and distances (d_{ij}) are as the crow flies. Both aspects would need to be improved in a real study but these figures illustrate the concepts and are not to be taken as the basis for detailed archaeological discussions

Rather like the stock market, evolution has both long-term and short-term characteristics that, most simply, can be thought of as a smooth general trend on which volatile short-term fluctuations are superimposed. Although the optimal solution is rarely, if ever, reached, numerous different solutions may exist that approach the optimal. So, one of these solutions may have Knossos as a key central place. Some small changes in certain parameters might then jog the system and cause it to fall into another configuration, equally close to optimal, but in which Knossos is no longer central. Or perhaps, it might transpire that Knossos is again central, but for different reasons and with a different set of connections. It is in this sense that our conclusions are statistical.

present purposes, Clarke notes that the physical analogies from particle physics have proved helpful in formulating models, as in gravity models for example, but that they have been conceptually unsatisfactory, able only to describe rather than explain. As for statistical mechanics theory, it “represents an interesting elaboration of the missing statistical and stochastic background behind the social physics approach” (1977, p. 20). Its basis is that “the most probable state of any system at a given time is the one which satisfies the known constraints and which maximises its entropy” (1977, p. 20).

This is an ideal approach for systems in quasi-equilibrium over long periods of time in which evolution is smooth. However, it can also indicate the onset of rapid change due to a shift in external circumstances. That there are dramatic ‘jogs’ to Aegean interaction systems seems quite clear – the innovation of the sail at the beginning of the MBA could be one, and the destruction caused by the Thera eruption might be another. In our approach such dramatic changes might be modeled by making a major change in the nature of the ‘landscape.’ In such circumstances, what were stable site exploitations in the valleys of this ‘landscape’ can become unstable configurations on its hills, which lead to a major readjustment in site use as the system migrates to the new ‘valleys.’

17.4.1 Generalities

The energy landscape we wish to describe has two types of coordinates; site vertex variables and link edge/link variables between (different) sites (generalizations of latitude and longitude). The energy of the landscape is denoted by its altitude. To be concrete, we think of these landscapes statistically, whereby the likelihood of achieving a particular value of H is given by the expression $\exp(-H/T)$ (the Boltzmann-Gibbs distribution of physics) where we have introduced another parameter T , the temperature in a physical context. The assumption is that the system will evolve from the unlikely ‘peaks’ (high H , so low likelihood) to the more likely ‘valley bottoms’ (low H , so high likelihood). Pursuing the simile further, the long-term evolution of a network can be thought of as a slow buckling of the terrain (comparable to plate tectonics). The short-term volatility, controlled by the ‘temperature’ T , whereby high volatility is ‘hot,’ low volatility ‘cold,’ corresponds to shaking the landscape (earthquakes). The parameters that control the contours of the landscape are measures of site independence or self-sufficiency, as well as constraints on population size, etc. Thus, for example, as populations grow or total trade volume increases, the landscape changes, and the positions of the valleys into which the system wishes to fall changes. This is rather like Broodbank’s increase in the number of links per island as population increases. Volatility here would correspond to short periods of drought, or unexpected changes in local population.

In all models, the ‘Hamiltonian’ energy function H , which characterizes each configuration of the system, separates into four terms:

$$H = -\kappa S - \lambda E + (j P + \mu T). \quad (17.1)$$

In a roughly defined way, H measures the ‘cost’ (in manpower, resources, etc.) of organizing the system of sites and their trading links. The aim is to find the configuration of the network that makes H as small as possible, for fixed values of κ , λ , j and μ .

Earlier, we raised the issue of robust system morphology. This is less a problem than one might think, as a partial resolution is given by the notion of a universality class. By this we mean that, rather than try to prescribe ‘fuzzy’ functions to

accommodate our uncertainty in H , we can hope for a family of ‘crisp’ functions that, provided we ask the right questions, will all give us the ‘same’ answer. The notion of topological congruence, taken from population biology, is most helpful. Functions which can be deformed into one another by stretching and squeezing are topologically congruent.

The individual terms that constitute H are understood as follows:

Sites: S only depends on the properties of the site vertices (usually islands) in isolation. As such, it is a sum of terms, one term for each site, which describes the exploitation of the site as a function of its de-trended population or occupation index (i.e. the fraction v_i of its total resources that have been exploited). Each site i is given a physical location, a fixed characteristic carrying capacity S_i (its effective size) and a variable occupation index v_i to be determined. One possible representation is that the active population at a site is $(S_i v_i)$ with S_i setting the maximum self-sustainable population at a site. Small rocky islands will have small S_i , yet they might have a large population, $v_i \gg 1$, if they play a pivotal role in the global network.

We will denote the total number of sites by N . We have tested our model using the list of 34 known MBA sites shown in Fig. 17.8. Initially, we assume that all sites are equally easy (or difficult) to exploit and have assigned equal relative sizes $S_i = 1$ to all.

At a later stage we could adopt a systematic approach to site location, such as the cultivatable land/population density method used by Broodbank (2000), leading to some loss of simplicity but at no cost to the numerical work.

By itself, v_i takes a minimum at some intermediate value.³ As a simple example of morphology, congruence here means little more than the observation that over-exploitation of resources incurs an increasingly non-linear cost, whereas under-exploitation permits growth.

Edges: E is the edge/exchange/trade term that shows how the sites interact with one another (trade, influence), in a way that depends on both the properties of the sites and the network and on the weight of their interactions. Most simply, it is a sum of terms, one term for every pair of sites that is linked by trade or for other reasons. We associate an edge variable e_{ij} to each link between sites i and j . One interpretation is that e_{ij} represents the trade *from* site i *to* site j and need not equal e_{ji} . We also define an effective distance d_{ij} *from* site i *to* site j , which here is just the distance between the two sites. In later work we will modify d_{ij} to take account of difference between land and sea transportation, prevailing winds and currents and so forth.

Constraints: The final terms (in brackets) enable us to impose constraints on population size (P), and on total trading links (and/or journeys made) in T .

Parameters: The parameters κ, λ, j, μ that control the contours of the landscape are measures of site independence or self-sufficiency, and constraints on population size, etc. Thus, for example, as populations grow or total trade volume increases, the optimal network (lowest energy configuration) changes. All other things being

³ It is not the value of S that is important but its derivative (slope).

equal, increasing λ enhances the importance of inter-site interaction, whereas increasing κ augments the importance of single site behavior. On the other hand, increasing j effectively corresponds to reducing population, and increasing μ reduces exchange.

Transformation Properties: To further constrain H , we demand that it behaves appropriately under special transformations. One such principle is the symmetry of the form of H under the interchange of labels of any two sites. That is, every site is governed by the same type of interactions as any other. This does not mean that every site is identical; we break this homogeneity when we incorporate different resources, S_i , and unequal distances, d_{ij} , between sites.

Finally, we demand that our H is the same whether we have several small sites clustered together in a small region or we treat this cluster as a single site of size equal to the sum of the smaller constituent parts. In this way the precise determination of what was the center of any one site should be unimportant and each site represents local population.

A ‘gravitational’ Hamiltonian: The example we have proposed in our initial proof-of-concept studies that embodies the above is shown in equation 17.2.

$$\begin{aligned}
 H = & -\kappa \sum_i S_i v_i (1 - v_i) - \lambda \sum_{i,j} V(d_{ij}/D) \cdot (S_i v_i) \cdot e_{ij} \cdot (S_j v_j) \\
 & + j \sum_i S_i v_i + \mu \sum_{i,j} S_i v_i e_{ij}
 \end{aligned} \tag{17.2}$$

The sums are over the different sites or over all pairs of sites, labeled by i or j . The first term, proportional to a constant κ , controls the size of sites as if there were no outside contacts. It is the logistic map used for simple models of population dynamics. Sites have negative energy for $0 < v_i < 1$, while the cost is positive for values larger than 1. Note that this term is invariant if we split a site into two by dividing S_i between the two new sites but keep the occupation fraction v_i the same for both new sites – our center-of-mass principle.

The second term allows for interactions, ‘trade.’ It is proportional to the total ‘populations’ at both ends of a link ($S_i v_i$) and to an edge weight variable e_{ij} . This ensures block renormalization, provided we ignore any new edge between the two new sites, since the number of possible edges involved also doubles, so that the total energy remains the same. For such models it is advantageous, in cultural exchange or trade, if both a site and its exchange partner have large resources. We realize that the cultural exchange/transmission that we are considering here is, by no means, simply economic, but, in contemporary economic parlance, we would say that this model embodies the advantages of a large consumer market and producer power.

It is through the interaction term that the effects of distance are included. This is done through a ‘potential’ term that is essentially zero for long range distances and one for short distances. Thus, direct long distance interactions give virtually no benefit and are unlikely to appear in our simulations; they are deemed to carry prohibitively high overheads. We introduce another parameter D which defines the

boundary between short and long distances. D is set to 100 km in all the Monte Carlo simulations shown below, as this is taken to be the distance scale appropriate for sailing in the MBA. By way of contrast, we might imagine that D should be 10 km for a rowing-based EBA simulation. The shape of the potential function used ought not to be too important, but, so far, we have worked only with the form $V(x) = 1/(1+x^{-4})$ which gives the desired behavior. In principle, we also need to introduce a very short distance scale. This is the minimum separation required before we consider two sites to be separate entities. This is needed for our block renormalization analysis to work appropriately. In practice, all our input sites are already deemed to be independent entities so that is not needed for the archaeological data.

17.5 Analytic (Mean-Field) Solutions

Before attempting any numerical modeling with the real island parameters, it is useful to see some of the behavior that might arise, using simple analytic approximations for an idealized network of sites. We make a mean field approximation in which we replace every value of v_i and every value of e_{ij} in $H[v_i, e_{ij}]$ by their average values, v and e , respectively. In doing so, we are removing all volatility and working at zero temperature. We then look for minima of $H(v, e)$, a two-dimensional energy landscape in which, for simplicity, we restrict ourselves to $0 \leq v \leq 1$ and that $0 \leq e \leq 1$. In some cases the lowest values will be at one of the boundaries and indeed the energy landscape will force the system to move to extreme values in one or both parameters.

As we have suggested, increasing λ increases the importance of inter-site interaction, whereas increasing κ increases the importance of single site behavior. If (λ/κ) is relatively small, the latter effect will overwhelm trade effects. In this case, we find that the stable energy minimum has every site close to the population $v = 0.5$ which is the optimal size when only local effects matter. The plot in Fig. 17.9 for small λ does indeed show a valley near this value.

On the other hand, when islands are not self-sufficient and (λ/κ) is relatively large, the latter effect may not be enough to inhibit runaway growth as trade brings benefits that outweigh local overpopulation effects and this is seen in Fig. 17.10.

However, in this situation we have a saddle point, and there are two possible outcomes: the runaway growth or the collapse of the system. That is, it may be better to reduce the population to reduce the penalty of having large populations and suffer the loss of advantageous trade. Iterating this brings us to collapse. Which solution prevails depends on which side of the saddle point leads to the lower valley bottom. In general, this will not be a blanket collapse. There will be a mixture of valleys and cols in this multidimensional landscape, and not all of the latter will be traversed in the direction of local collapse. Nonetheless, this shows the ease with which many sites in the network can either disappear ($v_i = 0$) or cease to communicate ($e_{ij} = 0$).

Roughly, provided λ is large enough, then, as λ increases from zero for fixed κ , there is monotonic growth in average site exploitation, from under-exploitation to

Fig. 17.9 The energy landscape for small (λ/κ) with the vertical axis H and horizontal axes v and e . In this regime, sites appear to be close to their optimal size and edges can have non-trivial values

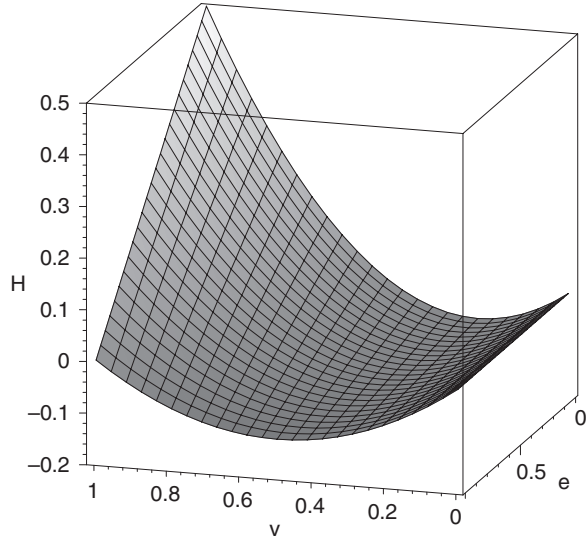
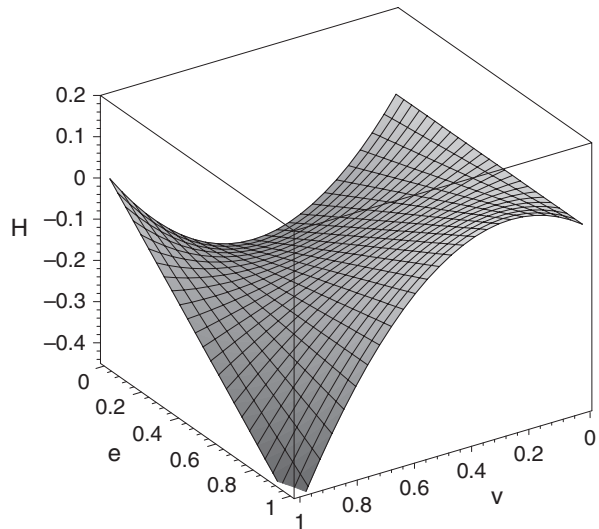


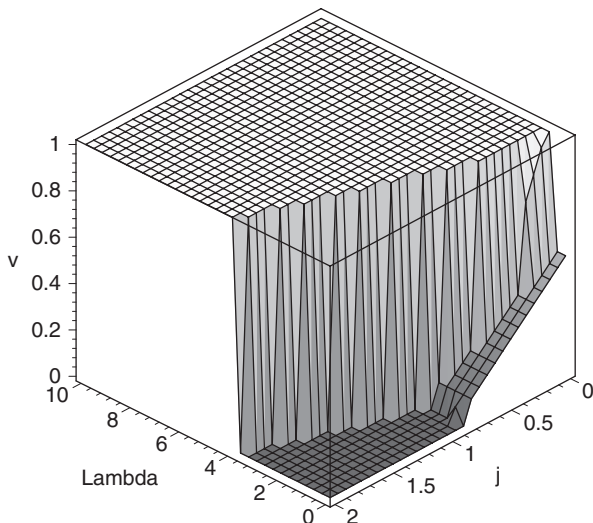
Fig. 17.10 The energy for large (λ/κ) with the vertical axis H and horizontal axes v and e . Now, the network is forced to extreme values



full exploitation. Provided λ is large enough, then, if λ is held fixed and we increase κ , all sites undergo medium exploitation as trading links become unimportant. The major difference occurs when λ (trading strength) decreases for small fixed κ (low self-sufficiency). Then, for only a small reduction in trading strength, exploitation of resources can collapse from full exploitation to no exploitation, which, naively, we might infer as site abandonment. This is shown in Fig. 17.11, in which we see the mean field average of v for varying λ and j , for fixed κ and μ .

This behavior is not specific to the particular form of H given above, but is a consequence of the fact that H has positive eigenvalues which, when they become

Fig. 17.11 The energy landscape for fixed κ and μ . The mean field average of v (vertical axis) is shown against varying λ and j . Vanishing v denotes collapse



negative, lead to instability. Gravitational models guarantee negative eigenvalues because of the way in which large sites interact preferentially with other large sites so that, for a prosperous network, the effect is non-linearly beneficial. This non-linearity has to compete against the non-linear costs of over-exploiting resources. Negative eigenvalues will arise once the system becomes less insular, with less stress on individual island self-sufficiency. In this regard we note the following observation by Broodbank et al. (2005):

For the southern Aegean islands in the late Second and Third Palace periods, an age of intensifying trans-Mediterranean linkage and expanding political units, there may often have been precariously little middle ground to hold between the two poles of (i) high profile connectivity, wealth and population, or (ii) an obscurity and relative poverty in terms of population and access to wealth that did not carry with it even the compensation of safety from external groups.

We note that these rapid collapses are not induced by volatility but correspond to a smooth buckling of the landscape.

17.5.1 Non-Gravity Models

Despite their virtues, models that minimize our ignorance, such as gravity models, are unlikely to be more than roughly correct and may be seriously at fault if applied across the whole network, even if applicable to the more homogeneous regions. As we have stressed, if this is the case, we need more detailed archaeological data than the aggregate data of the 'gravity' models, and the answers will depend on the island meso-scale population distributions.

To take the other extreme, in non-gravitational models, it may be advantageous to connect to bigger sites without any further advantage if one is big oneself. More

simply, it could be that an exchange/trade term at a site only depends on the existence of links to other sites and is insensitive to the resources/population available on the site itself. In contemporary economic parlance, we might term this a ‘supply-side’ model which ignores consumer demand. Conversely, the exchange/trade term at a site might be determined by the sites to which it is linked, insensitive to the resources/population available on those sites. In contemporary parlance, we might term this a ‘demand-side’ model. Because there is not the positive feedback in the virtues of becoming large, it is difficult to see how negative eigenvalues can arise. As a result, global collapse lessens. This is not to say that links or sites do not become abandoned, but they do it smoothly, as a consequence of shifts in the external population and maintenance pressures.

Taken together with the structural differences in the role of intra-island interactions, the existence of these further differences is important in that it shows that no universal network structure exists for island archipelagos. Instead, perhaps contrary to the hopes of some network theorists, the nature of the networks is strongly conditioned by geography and society.

17.6 Numerical Simulations

We are currently in the first stages of applying our models to the Aegean network of Fig. 17.8, using Monte Carlo methods. But before we display the results, it is useful to contrast our approach with PPA as described earlier. Figure 17.12 shows PPA applied to our sites (Table 17.1). In general, PPA emphasizes the closest links and tends to produce tightly connected groups of sites. One way to define these groups is to make sure that, for all the sites in one such group, there is a path from every site to every other site in the *same* group. Importantly, in moving along a path from one site to another, one always moves in the direction given by the arrows.⁴ For instance, in Fig. 17.12, we see that no path can involve a move from Chania (9) to Kastri (12) since the only edge between these two points is in the opposite direction. This link, along with the unidirectional link of Kea (14), ensures that the group of sites centered on the Greek Mainland is considered to be distinct from two other groups, one centered on the Cyclades and another on Crete.

Our optimal Monte Carlo model is not as unambiguous. *A priori*, it is difficult to make sensible estimates for the model parameters, and we have to search for robust ranges where features are visible, much as we have to choose the right scale and coverage when choosing a map for a problem in real life. Figure 17.13 shows a typical result from our Monte Carlo model for a similar density of edges as in the PPA network of Fig. 17.12. As we just noted, there is no direct or indirect link between Crete and the Cyclades in the PPA model if one respects the directionality of the

⁴ These groups are examples of what are called Strongly Connected Components (SCC) in Graph theory. Other definitions of these groups or communities are possible but should always show a similar result for PPA, something which is usually ‘obvious’ visually.

Table 17.1 Key to sites used in the Monte Carlo simulations

Site number	Site name
1.	Knossos
2.	Malia
3.	Phaistos
4.	Kommos
5.	Ayia Triadha
6.	Palaikastro
7.	Zakros
8.	Gournia
9.	Chania
10.	Akrotiri
11.	Phylakopi
12.	Kastri
13.	Naxos
14.	Kea
15.	Karpathos
16.	Rhodes
17.	Kos
18.	Miletus
19.	Iasos
20.	Samos
21.	Petras
22.	Rethymnon
23.	Paroikia
24.	Amorgos
25.	Ios
26.	Aegina
27.	Mycenae
28.	Ayios Stephanos
29.	Lavrion
30.	Kasos
31.	Kalymnos
32.	Myndus
33.	Cesme
34.	Akbuk

PPA links. The shortest possible Cyclades-Crete link would be from Thera (10) to Knossos (1) or Malia (2), which are around 100 km long. However, because each of these sites has several closer neighbors, the long distance Thera-Knossos/Malia links are never present in PPA. On the other hand, while the distance parameter D is 100 km in our Monte Carlo simulation of Fig. 17.13, our model will penalize, but not *a priori* exclude, links between sites separated by distance D or more. In the example show in Fig. 17.13, both Thera-Knossos/Malia links appear. In this case, the regional network benefits of such direct Cycladean-Cretan links seem to compensate for their large physical separation in our model, something the localized rules for network generation of PPA can never describe.

In terms of analysis, the increased complexity of our networks provides several challenges. For instance, the degree of a vertex is no longer a useful measure as edges are likely to carry a non-zero weight. For visualization, we have used a cut-off,

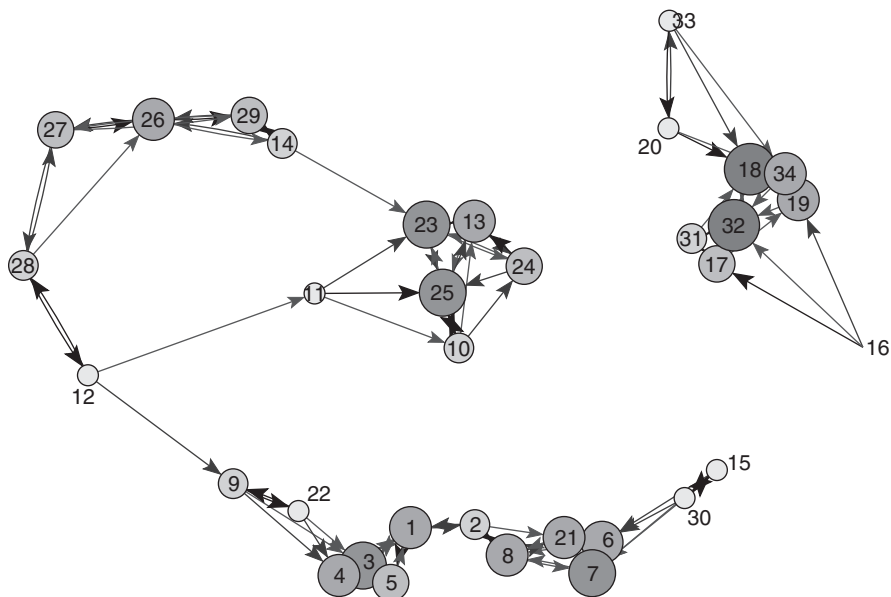


Fig. 17.12 PPA analysis where each site has three outgoing edges to the three nearest sites, the direction indicated by the arrows. The darker and larger the vertex, the larger the in-degree: Miletus (18) and Myndus (32) have largest in-degree of 6, Knossos (1) has in-degree 4, Malia has in-degree 2, Kastri (12) has in-degree one while Rhodes (16) has no incoming connections. However using *betweenness* as a measure we find that Knossos (1), Malia (2) and Kastri (12) are the most central. Note that the PPA does not assign a link between Crete and the Cyclades. Further if the sense of direction given by the arrows is respected then there are four large strongly connected cores: Crete, the Cyclades, the Dodecanese and mainland Greece

and, in our figures, we do not show edges or vertices which are below 10% of the size of the largest in that network. We could use a similar threshold method to map our network onto a simple graph such as that constructed by PPA. However the *raison d'être* of our work is precisely to exploit this as a feature. Thus, we have to introduce new methods for the analysis of the networks of island archaeology.

There are numerous areas where we are planning to make improvements to our models. We may adapt our input distances d_{ij} to reflect actual transport times rather than physical distances. The list and the sizes of the sites can be fitted to archaeological data, both adapting Broodbank's method of assigning sites on the basis of cultivatable area and exploiting modern GIS-based techniques. Within the model, we have variations where we use network distances rather than pure physical distances d_{ij} , both within the Hamiltonian and in the analysis.

Once we have constructed our networks, there are numerous ways to analyze them. We are developing the use of random walkers to rank sites. The Markov chain analysis of Hage and Harary for Oceania archaeology (Hage & Harary, 1991) is of this type but the generic method is used for ranking in many contexts from sports teams (Keener, 1993) to web pages (Brin & Page, 1998). In particular, we are limiting the walks that such agents can make which leads to simple models of

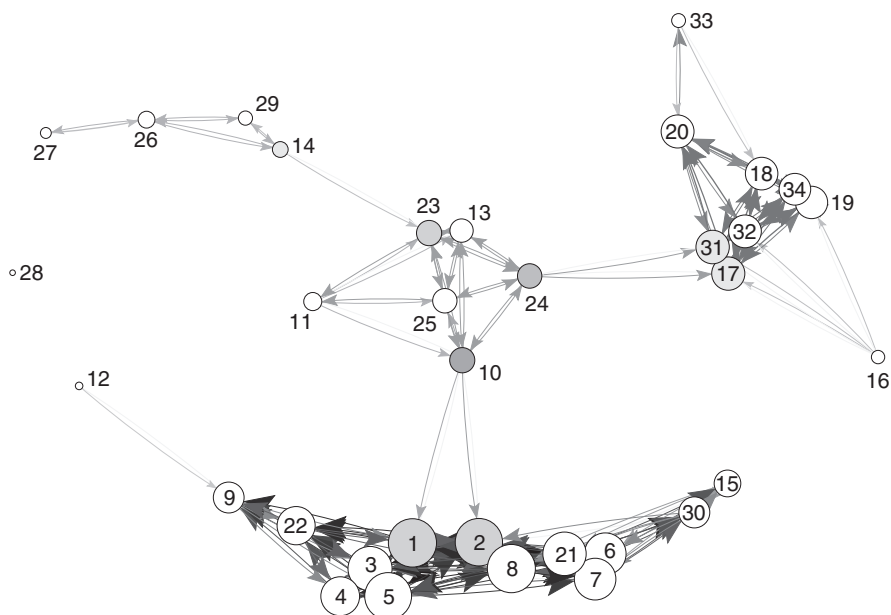
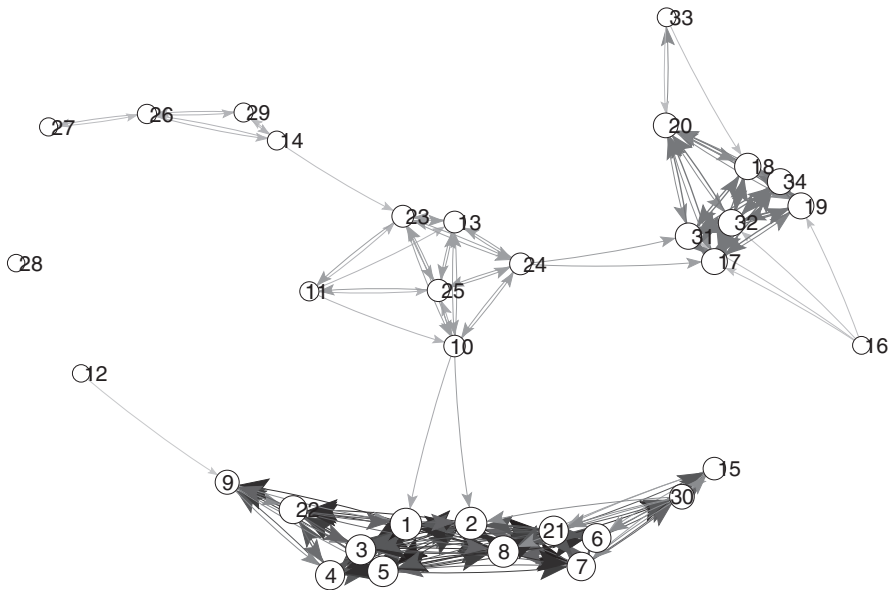


Fig. 17.13 Monte Carlo analysis for $\kappa = 2.0$, $\lambda = 1.0$, $\mu = 0.35$, $j = 0.7$ $D = 100$ km. The size of the vertices is proportional to their strength, the total weight of the in and out going edges. The largest are Gournia (8), Mallia (2) and Knossos (1) followed closely by the rest of central Crete. The Dodecanese are about half the strength and the Cyclades are a third of the strength. The darker the vertex is colored, the larger the *betweenness* and this shows a very different story with sites on the edges of clusters scoring highly. This includes Malia and Knossos but now the Cyclades scores even higher than these indicating their central role if all vertices are treated equally. However all vertices are clearly not equal in our networks and many standard measures of network properties, such as betweenness, are of little use in our work

cultural transmission on our networks, a topic modelled in rather different ways elsewhere (see, for instance, Neiman, 1995; Bentley & Shennan, 2003; Bentley, Hahn, & Shennan, 2004 and references therein).

Finally, we are beginning to study problems of temporal evolution. This can occur in the form of slow ‘adiabatic’ changes, such as population build up, or as quick ‘quenches.’ We can simulate, for example, the change that occurs between Late Minoan IA and IB (c. 1600 BC), when the volcanic eruption of Thera removes that island (the site of Akrotiri in particular) from the system (this is shown in Fig. 17.14, before and after). Whatever the scenario, it is quite possible that the system gets stuck, for a time, in a meta-stable state with the instability only apparent much later. This might be a good model for the transition that occurs after Late Minoan IB (c. 1500 BC), when sites across Crete are destroyed and the balance of power shifts to the Greek mainland (i.e. Mycenae) for the following three centuries. In future work, we hope to move towards the study of longer-term dynamics of this kind.

A



B

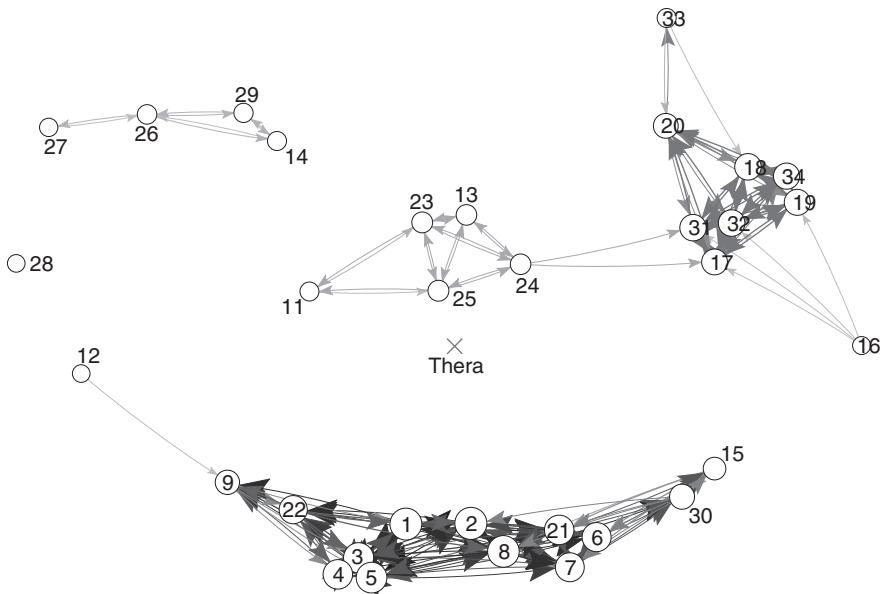


Fig. 17.14 The same values are used, but, in network **B**, Thera has been removed (marked by a cross) while it remains in place in the network **A** (which is the network of Fig. 17.13 displayed in a different manner). The sizes of sites (size of *circles*) are similar, but now the Dodecanese is ranked far higher

17.7 Conclusions

Our assumptions about the way in which the meso-scale of the individual island societies is accommodated within the macro-scale of the overarching MBA network have been an important part of our model. This introduction of what we might call ‘social’ or ‘behavioral’ aspects into the physical system shows that, rather than assuming diffusion or gravitational pull to be intrinsic physical properties of the system, they are only relevant given certain (social) conditions, as exemplified by the ‘gravitational’ MBA Aegean and the very different EBA Cyclades. We can, of course, include other ‘social’ conditions of the network by adjusting, for example, the degree of commitment to local resources or trade, to achieve better fits with perceived or actual past scenarios.

Both the analytical and numerical approaches outlined in the previous sections furnish insights into the articulation of the physical and relational dimensions of regional interaction networks. This articulation is expressed through the gravity models discussed in the analytical section. In terms of the results arising from this investigation, there is clearly still much to do. There are, at this stage, interesting indications that the development in the later MBA of affiliation networks linking certain larger sites might be amenable to further exploration using gravity models. The Monte Carlo analyses run thus far do seem to testify to the importance of the link between north-central Crete and Thera under certain conditions, a pattern that is very clearly seen in the archaeological evidence. Our goal is to explore these conditions much more fully. The same goes for the consistently central role of Knossos, or of other sites in north-central Crete, and the apparent robustness of this pattern to changes in network conditions. There is much we need to do to produce more results and to compare the results with the archaeological data. Our work indicates that these techniques open up numerous possibilities and offer new means of assessing the different kinds of networks that have long preoccupied the social sciences in various manifestations. Post-Watts and Strogatz, other new approaches have also emerged. Rather than use these tools for their own sake, however, we must seek to ensure that our techniques are commensurate with the complexity of the archaeological data and ensure any conclusions are robust against changes in the details of our models.

Acknowledgments We are grateful to participants in various ISCOM workshops over the past three years for their helpful comments on the ideas presented in this paper. On the archaeological side, we warmly thank Cyprian Broodbank and Todd Whitelaw for their insightful comments on earlier drafts, and Andy Bevan for his advice and help on GIS matters. Figures 17.12, 17.13 and 17.14 were created with PAJEK (Program for Large Network Analysis, <http://vlado.fmf.uni-lj.si/pub/networks/pajek/>).

References

Adams, E. (2006). Social strategies and spatial dynamics in Neopalatial Crete: an analysis of the north-central area. *American Journal of Archaeology*, 110(1), 1–36.

- Batty, M. (2005). Network geography: relations, interactions, scaling and spatial processes in GIS. In D. Unwin, & P. Fisher (Eds.), *Re-presenting GIS* (pp. 149–170). Chichester, UK: John Wiley.
- Berg, I. (1999). The southern Aegean system. *Journal of World-Systems Research*, 5(3), 475–484.
- Blake, E. (2002). Spatiality past and present: an interview with Edward Soja, Los Angeles, 12 April 2001. *Journal of Social Archaeology*, 2(2), 139–58.
- Blake, E. (2004). Space, spatiality and archaeology. In L. Meskell, & R. Preucel (Eds.), *A companion to social archaeology* (pp. 230–254). Oxford, UK: Blackwell.
- Bentley, R. A., & Shennan, S.J. (2003). Cultural transmission and stochastic network growth. *American Antiquity*, 68, 459.
- Bentley, R. A., Hahn, M. W., & Shennan, S. J. (2004). Random drift and cultural change. *Proceedings of the Royal Society of London B*, 271, 1443.
- Bender, B. (Ed.). (1993). *Landscape: Politics and perspectives*. Oxford, UK: Berg.
- Brin, S., & Page, L. (1998). The anatomy of a large-scale hypertextual Web search engine. *Computer Networks and ISDN Systems*, 30, 107–117.
- Broodbank, C. (2000). *An Island archaeology of the early cyclades*. Cambridge, UK: Cambridge University Press.
- Broodbank, C., Bennet, J., & Davis, J. L. (2004). Explaining social change: Studies in honour of Colin Renfrew. *McDonald Institute of Archeological Research*, 73–81.
- Broodbank, C., Kiriati, E., & Rutter, J. (2005). From Pharaoh's feet to the slave-women of Pylos? The history and cultural dynamics of Kythera in the Third Palace Period. In A. Dakouri-Hild, & S. E. Sherratt (Eds.), *Autochthon. Papers Presented to O.T.P.K. Dickinson on the Occasion of his Retirement*, BAR Int Series (pp. 70–96). Oxford, UK: Archaeopress.
- Carrington, P. J., Scott, J., & Wasserman, S. (Eds.). (2005). *Models and methods in social network analysis*. Cambridge, UK: Cambridge University Press.
- Chase-Dunn, C., & Hall, T. D. (1997). *Rise and demise: Comparing world-systems*. Boulder, CO: Westview Press.
- Chorley, R. J., & Haggett, P. (Eds.). (1967). *Models in geography*. London, UK: Methuen.
- Clarke, D. L. (1968). *Analytical archaeology*. London, UK: Methuen.
- Clarke, D. L. (1977). Spatial information in archaeology. In D. L. Clarke (Ed.), *Spatial archaeology* (pp. 1–32). London, UK: Academic Press.
- Davis, J. L. (1979). Minos and dextithea: Crete and the cyclades in the later bronze age. In J. L. Davis, & J. F. Cherry (Eds.), *Papers in cycladic prehistory* (pp. 143–157). Los Angeles, CA: Institute of Archaeology.
- Doel, M. (1999). *Poststructuralist geographies: The diabolical art of spatial science*. Edinburgh, UK: Edinburgh University Press.
- Evans, T. S. (2005). Complex networks. *Contemporary Physics*, 45, 455–474.
- Granovetter, M. (1973). The strength of weak ties. *American Journal of Sociology*, 78 1360–1380.
- Granovetter, M. S. (1983). The strength of weak ties: a network theory revisited. *Sociological Theory*, 1, 201–233.
- Hage, P., & Harary, F. (1991). *Exchange in oceania: A graph theoretic analysis*. Oxford, UK: Clarendon Press.
- Hage, P., & Harary, F. (1996). *Island networks: Communication, kinship and classification structures in oceania*. Cambridge, UK: Cambridge University Press.
- Haggett, P. (1965). *Locational analysis in human geography*. London, UK: Arnold.
- Harvey, D. (1973). *Social Justice and the City*. London, UK: Arnold.
- Harvey, D. (1996). *Justice, nature and the geography of difference*. Oxford, UK: Blackwell.
- Hetherington, K. (1997). In place of geometry: the materiality of place. In K. Hetherington, & R. Munro (Eds.), *Ideas of difference: Social spaces and the labour of division* (pp. 183–99). Oxford, UK: Blackwell.
- Hillier, B. (2005). Between social physics and phenomenology. In *Fifth space syntax symposium* (pp. 13–17). June 2005, Delft, The Netherlands.
- Irwin, G. (1983). Chieftainship, kula and trade in Massim prehistory. In J. W. Leach, & E. Leach (Eds.), *The kula: New perspectives on massim exchange* (pp. 29–72). Cambridge, UK: Cambridge University Press.

- Kardulias, P. N. (Ed.). (1999). *Leadership, production and exchange: World-systems theory in practice*. New York, NY: Rowman and Littlefield.
- Keener, J. P. (1993). The Perron-Frobenius theorem and the ranking of football teams. *SIAM Review*, 35(1), 80–93.
- Knapp, A. B., & Ashmore, W. (1999). Archaeological landscapes: constructed, conceptualized, ideational. In W. Ashmore, & A. B. Knapp (Eds.), *Archaeologies of landscape: Contemporary perspectives* (pp. 1–30). Oxford, UK: Blackwell.
- Knappett, C., & Nikolakopoulou, I. (2005). Exchange and affiliation networks in the MBA southern Aegean: Crete, Akrotiri and Miletus. In R. Laffineur, & E. Greco (Eds.), *Emporia: Aegeans in East and West Mediterranean* (pp. 175–184). Liège, Belgium: University of Liège, Aegaeum 25.
- Lefebvre, H. (1991). *The production of space*. Oxford, UK: Blackwell.
- McGuire, R. H. (1996). The limits of world-systems theory for the study of prehistory, In P. N. Peregrine, & G. M. Feinman (Eds.), *Pre-Columbian world systems, monographs in world prehistory* (Vol. 26, pp. 51–64), Madison, WI: Prehistory Press.
- Murdoch, J. (2006). *Post-structuralist geography: A guide to relational space*. London, UK: Sage.
- Neiman, F. D. (1995). Stylistic variation in evolutionary perspective: inferences from decorative diversity and inter-assemblage distance in Illinois Woodland Ceramic assemblages. *American Antiquity*, 60(1), 7–36.
- de Nooy, Mrvar, A., & Batagelj, V. (2005). *Exploratory social network analysis with pajek*. Cambridge, UK: Cambridge University Press.
- Peregrine, P. (1996). Introduction: world-systems theory and archaeology. In P.N. Peregrine, & G. M. Feinman (Eds.), *Pre-Columbian world systems, monographs in world prehistory* (Vol. 26, pp. 1–10). Madison, WI: Prehistory Press.
- Rihll, T. E., & Wilson, A. G. (1991). Modelling settlement structures in Ancient Greece: new approaches to the polis. In J. Rich, & A. Wallace-Hadrill (Eds.), *City and country in the ancient world* (pp. 59–95). London, UK: Routledge.
- Schortman, E. M., & Urban, P. A. (1992). The place of interaction Studies in archaeological thought. In E. Schortman, & P. A. Urban (Eds.), *Resources, power and interregional interaction* (pp. 3–15). New York, NY: Plenum Press.
- Smith, A. T. (2003). *The political landscape: Constellations of authority in early complex polities*. Berkeley, CA: University of California Press.
- Soja, E. (1996). *Thirdspace: Journeys to Los Angeles and other real-and-imagined places*. London, UK: Blackwell.
- Stein, G. J. (1998). World systems theory and alternative modes of interaction in the archaeology of culture contact. In J. G. Cusick (Ed.), *Studies in culture contact: Interaction, culture change and archaeology*. Occasional Paper No. 25, (pp. 220–255), Carbondale, IL: Center for Archaeological Investigations.
- Terrell, J. (1977). *Human biogeography in the solomon Islands*. Chicago, IL: Field Museum of Natural History.
- Thrift, N. (1996). *Spatial formations*. London, UK: Sage.
- Tilley, C. (1994). *A Phenomenology of landscape: Places, paths and monuments*. Oxford, UK: Berg.

Conclusion

David Lane, Denise Pumain and Sander van der Leeuw

Based on the variety of examples that were presented in the different chapters of this book, and the theoretical constructs that emerged during the work of the ISCOM group, we are now confident that some conclusions can be drawn from our attempt at applying a complex system perspective to social systems. In the process, we have identified a few theoretical principles that are essential when introducing the methodologies of complexity in social sciences. We are concluding the volume with a summary of these principles.

A Theory of Complex Social Organisation

We emphasized in our first chapter that when shifting from biology to social sciences, the concept of population thinking, essential to biological evolution theory, has to be replaced by the concept of organisation thinking as the primary foundation in a theory of innovation and social change. Organization thinking requires that no description of a human organisation can separate structure, function and processes. This principle has important consequences for the ways in which the complexity approach can be used in the social sciences. Obviously, it is not very suitable for any attempts at deepening the theoretical concepts and methodological tools of each discipline alone, but on the contrary it incites to an interdisciplinary recognition of social entities and problems requiring converging conceptual approaches to construct a sophisticated theory of social organisation.

Social Organisations are Tangled Hierarchies

The second conclusion concerns the structures of social organisation. Once ‘invented’ (Chapter 2), hierarchical organisation seems to be present in most, if not all, human socio-cultural systems. Different hierarchical levels can be identified

D. Lane (✉)

Department of Social, Cognitive and Quantitative Sciences, University of Modena and Reggio Emilia, Reggio Emilia, Italy

in many social organisations in the guise of institutions with a different ‘sphere’ (financial, religious, medical, military, etc.) and a different ‘span’ of action (family, community, city, province, nation or even wider, international). In many cases these institutions are ‘nested’ at different levels because, at each level, different structures emerge which carry not only quantitatively but also qualitatively different functions. Such organisations often become ‘enslaved’ to more encompassing ones. These in turn absorb the functionalities of the lower level, as part of a process of time-space expansion. Most of this process can be explained as a consequence of the dynamics of interaction between the parts that created the whole, but simple micro-macro processes alone are not enough. On the one hand, the levels are not strictly embedded levels; very often the social hierarchies are tangled hierarchies; this means that relationships between distant levels can play an important role in their evolution. Important interactions can extend over many scales. On the other hand, the interactions are never of a unique kind. Even when they are aiming at building a specific institution to enable a given (specialised) type of transaction (for instance, a seemingly ‘pure’ economical relationship establishing a market) it is necessary to consider other social implications because those will play a non-negligible role in the evolution of that institution. That is why networks envisaged as mere graph representations cannot be considered sufficient to describe such entities. Even processes that reinterpret structures and functionalities in social change are operating through networks. The linkages in these networks never convey a single kind of interaction, so that the meaning of the link is changed through the process that it is involved in. The nodes may be more pervasive, but as the links between them are continuously modified at the scales of time under consideration in the most common descriptions, the nature and functionalities of the nodes change along with them.

Emergence Does not Only Involve Bottom-Up Interactions

When describing the formation and maintenance of social hierarchical organisations, we observe that most of the structures in social systems cannot be created uniquely ‘bottom-up’, in a process in which interactions between elements at one level induce emergent properties that define the structure at the next higher level. One of the fundamental differences of our approach with the ‘population’ approach frequently used in biology is the fact that, in human societies, individuals shape, and are shaped by, the society in a reciprocal interaction. In any dynamic social hierarchy, lower-level behaviours are strongly affected by interactions with elements and processes acting at higher levels and vice-versa. In Chapter 2, we have emphasized the importance of such interaction by stressing the fact that ‘enculturation’ in a collective system of values and norms, leading to the predictability of behavioural expectations, is essential to reduce the potential error rate in social expectations, and thus to optimize social interaction and minimize its costs. But it was not specified which could be the nature of the processes leading to the ‘culture’ involved.

Once explicit social and institutional hierarchies develop, the values of the upper levels of such hierarchies are part of the total culture, and thus exert an important influence on the behaviour of the lower levels of the hierarchies involved. But the

processes occurring between individuals and entities at lower levels also impact on, and constrain, the higher levels.

Translated to our theory of urban systems, we therefore insist that no theory of 'the' city can hold in the social sciences without considering other levels above it, such as the system of cities, or other global systems of relationships. The national or global urban system has a much greater impact on individual cities than the individuals, households and firms living in the city can have by themselves. It is easy nowadays to observe that cities are affected by the same kind of changes all over the world. We have demonstrated that cities co-evolve, in such a way that the speed of their common change (in population, architecture, economy, culture. . .) is much more rapid than the slower changes which can modify their relative positions in terms of size and socio-economic specialisation, or their images (as used now by urban marketing). In that context of highly competitive urban environment, it is very unlikely that a single city can 'invent' or concentrate an innovation to such an extent that it would enable its inhabitants to monitor the further destiny of their city in a chosen direction. The initiatives of the inhabitants are constrained by initiatives taken in other cities, as well as by the aggregate properties of their city (such as the pre-existing size of the city, or its former economic specialisations which shaped its kind of society). These collective properties depend themselves on the upper level of the system of cities, because their values are measured relatively to other cities – when private or public actors search a location for their investments. These 'attributes' which qualify the relative situation of a city within various networks or systems of cities do at any given moment reduce or open the range of its opportunities. The emergence of a city's attributes and its socioeconomic trajectory are by no means resulting from the interactions of the local actors only. Taking multi-level reciprocal interactions into account provides a much more nuanced epistemological position for social sciences than the commonly advocated methodological individualism.

The Place of Scaling Laws in Social Systems

Scaling behaviour can be understood in two ways. Firstly, as in the case of D'Arcy Thompson's (1917) allometric laws, it can establish systematic non-linear relationships between the size of a part and the size of a whole during evolution. This longitudinal view is very often observed by cross-sectional (synchronous) comparison of the sizes of the elements that constitute a whole, as exemplified by Zipf's 'rank size rule' (1949) for urban systems. Such scaling applied to one variable may provide insight in the growth process that produces that kind of statistical distribution. But many studies have demonstrated that this model is under- or even ill-determined as long as the interactions between the elements are not explicit in the growth process (Pumain, 2004; Favaro & Pumain, 2008).

Scaling approaches are perhaps more interesting when applied to variables of different kinds that are measured on the same entities. The major result that emerged from such bivariate comparisons between biological and social systems is that both kinds of systems systematically exhibit increasing efficiency in energy consumption

as they grow in size, but that human social systems differ from biological systems in that with respect to activities linked to innovation, human systems show super-linear scaling, meaning that, proportionately, these activities increase with size. That implies that the creation of something new has a positive feedback effect on urban size, stimulating further growth. But, such creation requires a higher level of investment to appear there. Two complementary explanations have been proposed for this observation. The one developed in Chapter 7 emphasizes the tension between urban infrastructures that scale sublinearly, becoming more efficient as cities grow in size, and other social activities scaling superlinearly, giving rise to accelerating growth, towards a finite time singularity, thus linking inextricably the desired properties of fast economic and technological development to crises of adaptation. In Chapter 8, the explanation given insists on the fact that cities are not isolated but embedded in a system of interacting cities, which are differentiating not only in size but also qualitatively through the process of their co-adaptation to innovations. New activities that are generated by each large innovation cycle locate first in the largest cities (with high potentials of all kinds of capital investment- human, technical cultural or financial-, enabling high returns but at high costs), then spread down the urban hierarchy until they concentrate in cheaper locations in smaller towns. In this interpretation, the historical and spatially and hierarchically interactionist character of the innovation process generates the observed variety of exponents in urban scaling laws.

What are the Constraints in the Evolution of Social Systems?

In general, it seems that any society that attempts to infinitely increase its size in order to ensure its continued existence would blow up. This seems true for non-human as well as human species, whatever the specific conditions that bring about the end. In biology, this thesis underpins the Lotka-Volterra predator-prey dynamic, for example. But such a catastrophic end has also been predicted for human organisations. It is in effect a recurrent theme in many societies and cultures; in our own culture, it underpinned Malthus' (1803) analysis of population and resource growth, for example, but was also highlighted in the early part of the last century by Spengler (1918–1922). Currently, it is on everyone's mind due to our growing awareness that our environment is rapidly degrading. Thus, we have to ask ourselves whether there are in effect limits to growth (Meadows, Meadows, Randers, & Behrens, 1972, Meadows, Randers, & Behrens, 2004)?

For non-human species and groups, there seems indeed a limit to the extent to which the size of groups, or the density of populations, can be scaled up. Humans, however, have thus far been able to overcome many barriers to population, economic and other kinds of growth. Their number has increased from some 60,000 around 140,000 BP to 6,500,000,000 today, and their social organizations from bands of about 500, to cities of 20,000,000 or more. On a number of occasions, their societies have lost momentum and dissipated. But each and every time, human societies have collectively learned (sometimes slowly), or at least modified their behaviour and

managed to grow again. From the archaeologist's perspective, therefore, it is just possible that even the current, very major, limits to growth that are appearing all around us can be overcome.

In that case, one should first of all inquire what might be the limiting constraints? There are numerous candidates, energy, water, space, climate change and political domination among them. Each of these has also been of relevance in the past, to some extent at least. What seems particular to our current predicament, however, is the fact that these phenomena are now interacting globally, and that it is therefore not possible for any one society to either move or disband, or to attach itself to another, more encompassing, system of institutions to spread the risks involved. It seems to us, therefore, that we should posit the problem in a different way, and define as a 'crisis' *the temporary incapacity of a society to process the information necessary to deal with the dynamics in which it is involved* (van der Leeuw, 2006 and 2008). In doing so, we shift the 'cause' of the potential crisis, the 'constraint' that may be exceeded, from the environment to the society itself, arguing that as humans the only way we can transcend the constraint is by accepting the limits of our current ways of doing things, and therefore attempt to overcome them by changing our outlook, our organization, our society and the resources on which it depends. After all, it is our current over-dependency on specific resources and ways of doing that makes us dread their potential shortage so greatly. If our societies could change the ways they were doing things, these problems might decrease in importance.

Interactions Between Demography, Economy, Natural Resources and Ecosystem Services, Local and Global Policies

We are all of course well aware that all the above domains are part of one and the same system, which nowadays encompasses the whole globe. The dynamics in each of these domains, and an almost infinite number of others, interact to the point that we cannot separate them 'cleanly' (see above). Yet, in the last two hundred years, since the establishment of the current academic system of 'disciplines', their study has become more and more isolated, with some important consequences. Let us take the domain of natural resources and ecosystem services as an example. Since the 14th century, in western civilization, the gap between the study of the 'natural environment' and that of society has increased to the point that even such eminently 'human' things such as thinking, feeling, etc., have been 'explained' in terms of complex chemistry (Evernden, 1992). We would argue that that separation ('taking people out of the rest of the world') is to a large extent responsible for many of our current environmental challenges. *Mutatis mutandis*, the same can be said of human demography – it is the emphasis on improving human health, without taking the longer-term consequences for the environment or human demography into account, which has led to the current demographic stress, with its consequences for resources, space, conflict, etc. And the undue emphasis on the economic domain has led to the fact that the current hyper-capitalist financial system is controlling the well-being

of so many. Clearly, re-integration of these domains of study and action is essential if we are to creatively deal with the current ‘crisis’.

Questions of Time

Another set of closely related issues concerns the way the social sciences deal with time. In our society, much more attention has been paid traditionally to looking back towards the past, to history, the search of origins and explanation rather than to looking forward in time, to prediction, creation of scenario’s, etc. Yet humans do dream, expect, anticipate, and direct their actions towards the realization of these expectations. In any effective social science, looking forward in time must therefore have its place. It is one of the great merits of the complex systems approach to have re-emphasized the role of history in the natural sciences by pointing out that there are many situations in which the trajectory of a system cannot be reconstructed from the position of the system at the end of the trajectory. But the emphasis on emergence in the complex systems approach also brings home the fact that we can do better than look back upon history (searching for the origins of the present) by looking forward *with* history, i.e. constructing the present from the past, rather than reconstructing the past from the present. And from there it is but one step (if not always an easy one) to extrapolate towards the future in the form of multiple scenarios. In our opinion, this domain of intellectual activity, which is emergent itself, merits considerable more attention than it is currently getting. Modelling, of course, plays an essential role in developing it – but it should be made crystal clear that this kind of modelling serves to scrutinize the implications of theories and expectations, rather than to ‘re-create the dynamics of the real world’.

Modelling also allows us to deal with another shortcoming of much social science research – the fact that such research only takes a very limited number of temporal scales into account (usually between two and four), whereas of course, in the real world, an almost infinite number of temporal scales interact in any system (though some of these may be more critical to the phenomena observed than others). To see the importance of this, one could look at the long-term of any relationship between a society and its environment. As the interaction between the two develops, humans will first act upon the most frequently observed phenomena. In ‘appropriating’ them (Ingold, 1986), society will to some extent change the dynamics of the environment at different temporal scales, some of them short (and therefore easy to observe), and some of them longer (and therefore unknown for a long time). When the latter later come into play, we call them ‘unintended consequences’. In effect, with time, the risk spectrum changes because people ‘appropriate’ short term, known, risks, but initiate long-term, unknown, ones. In due course, this creates the kind of ‘time-bombs’ that ‘suddenly’ emerge, such as the current greenhouse gas crisis. This is exacerbated by the fact that as this shift occurs, more and more of the society’s attention is drawn to the ‘unintended consequences’ that emerge. The society’s attention is thus drawn away from realizing a long-term vision to dealing with immediate challenges. Little by little, due to the growth in complexity of the system, and the emergence of ever more ‘unintended’ challenges, the range of

temporal scales of which the society is conscious will shift towards the short end of the spectrum. In our opinion, the social sciences should use modelling to investigate this phenomenon and see if they can widen the spectrum of temporal scales they are dealing with, and in particular re-introduce long term perspectives, and study the role of anticipation in societal dynamics.

Innovation as a Threat and as a Way Out

Increasingly, over the last two centuries, our society has focused more and more on innovation and change, to the detriment of history and stability. This may be seen as yet another side effect of the above-mentioned shift in risk spectrum: as more and more challenges become noticeable, the push for change becomes, of course, more important. As mentioned in the introduction, this tendency has now accelerated to a point where the instantiation of several completely new sets of ways to deal with the material, biological and social world threatens to overwhelm society. We are of course referring to the so-called NBIC revolution mentioned in the introduction. Innovation, we must conclude, is a double-edged sword: it allows us to deal with certain challenges, but poses others.

The ambivalent nature of innovation has prompted us to look more closely at the phenomenon and its role in modern society. In doing so, we have – again – profited from the complex systems approach, and in particular its capacity to study emergent phenomena. What, after all, is a more ‘pure’ emergent phenomenon than an invention and its subsequent introduction into society (here termed innovation to distinguish it from invention)? Accordingly, we have developed a ‘generative’ approach that looks at inventions and innovations in context. One way to describe it is to argue that all inventions are born in the context of a ‘mental map’ that is constituted by people, ideas, artefacts, functions, materials and processes. The structure of that map is in itself governed by the ‘culture’ of the society in which the invention takes place. The map constrains the range of inventions that may occur, and enables certain of them to occur with greater probability than others. Without ‘determining’ the kinds of inventions that are likely to be the response to certain challenges, the map nevertheless ‘steers’ the overall inventive capability of the society.

Looking at invention this way, it should be possible – and we have begun to explore this possibility – to become pro-active about invention and innovation, rather than only re-active to innovations once they have emerged. Certainly, this requires the development of more tools and concepts than are currently available. But in our opinion, it would enable our societies to mobilize their talents and energies more directly towards dealing with the challenges that are currently facing us as a result of centuries of pushing the risk-spectrum towards the future, and thereby minimize the disruptions that are likely to be inevitable as ‘the chickens come home to roost’.

References

D’Arcy Thompson, D. W. (1917). *On Growth and Form*. 1992 : Dover reprint of 1942 2nd ed. (1st ed., 1917).

- Evernden, N. (1992). *The social creation of nature*. Baltimore: Johns Hopkins University Press.
- Favaro, J. -M. & Pumain, D. (2008). Urban hierarchies explained by a spatial interaction model including innovation cycles. Université Paris 1, UMR Géographie-cités.
- Ingold, T. (1986). *The appropriation of nature: essays on human ecology and social relations*. Manchester: Manchester University Press
- Malthus, T. R. (1803). *An essay on the principle of population* (1798 1st ed., plus excerpts 1803 2nd ed.), Introduction and assorted commentary on Malthus by Philip Appleman. New York: Norton. www.econlib.org/Library/Malthus/malPlong.html
- Meadows, D. H., Meadows, D. L., Randers, J., & Behrens III, W. W. (1972). *The limits to growth*. New York: Universe Books.
- Meadows, D. H., Meadows, D. L., Randers, J., & Behrens III, W. W. (2004). *The limits to growth: The 30-year update*. London: Chelsea Green Publishing
- Pumain, D. (2004). *Scaling laws and urban systems*. Santa Fe Institute, Working Paper n°04-02-002, 26p.
- Spengler, O. (1918–1922). A. Helps, & H. Werner (Eds.), *The decline of the west*. (C. F. Atkinson, Trans.). New York: Oxford UP, 1991.
- van der Leeuw, S. E. (2006). Crises vécues, crises perçues. In C. Beck, Y. Luginbühl, & T. Muxart (Eds.) *Temps et espaces des crises de l'environnement* (pp. 351–368). Paris: Editions Quae
- van der Leeuw, S. E. (2008). Agency, networks, past and future, In C. Knappett, & L. Malafouris (Eds.), *Material and nonhuman agency*. New York: Springer
- Zipf, G. K. (1949). *Human behaviour and the principle of least-effort*. Cambridge, USA: Addison-Wesley.

Author Index

A

Agar, M., 320
Atlan, H., 89
Arrow, K., 316
Arthur, B., 169, 332

B

Bairoch, P., 203, 212, 256, 352
Bak, P., 177, 364, 376, 447
Batty, M., 177, 331, 433, 438, 452
Berry, B. J. L., 102, 197, 201, 241, 357
Bettencourt, L., 3, 103, 175, 221–235, 438, 448
Bijker, W., 307–308
Bretagnolle, A., 197–219, 240–241, 331, 335

C

Christaller, W., 102, 201, 340, 438, 442, 444

D

Darwin, C., 32–33, 119–126, 127–129, 142
David, P., 290
Dosi, G., 308

F

Ferber, J., 332, 334, 343–344
Florida, R., 171, 228, 256–257
Fontana, W., 386
Friedman, M., 118–119, 144–147

G

Gibrat R., 201, 213, 240–241, 434
Gilbert, N., 362, 389
Gould, S., 129, 414, 415

H

Hägerstrand, T., 239, 313
Hayek, F. Von., 118–119, 122, 130, 134–148
Helbing, D., 103, 224, 225, 433–448
Holland, J., 19, 306, 401

J

Johnson, G., 101, 102

K

Kauffman, S., 364, 387, 420
Keynes J., 135, 145
Knappett, C., 451, 456
Krugman, P., 177
Kühnert, C., 103, 433–438, 446

L

Lakoff, G., 118, 400
Lane, D., 1–40, 43, 85, 256, 263–288,
289–290, 292, 299, 308–309, 311,
317–318, 320, 322, 324, 361, 364, 366,
367, 379, 396, 410, 424, 481–487
Lobo J., 3, 103, 114, 221, 225, 434, 438
Lundvall, B.-A., 313

M

Malinowski, B., 104
Malthus, T., 119–124, 126, 132, 484
Marx, K., 122, 136, 144
Maxfield, R., 11–40, 263–288, 289–290, 292,
299, 308, 311, 317–318, 320, 322, 379,
396, 410
Mayhew, B., 101
Mayr, E., 12, 13, 127

N

Nelson, R., 66, 308–309, 313, 316
Nooteboom, B., 313

P

Pred A., 201, 238–239, 257
Pumain D., 1, 3, 102, 197, 198, 201, 204,
212–213, 215–216, 237–259, 317,
331–358, 434, 438, 448,
481, 483

R

Rosenberg, N., 117, 138, 313, 320
Rossi, F., 6, 289–309, 311–325, 391
Russo, M., 311–325, 391

S

Sanders L., 102, 331–334, 353, 354
Schumpeter, J., 130, 131, 416
Sen, A., 143–144
Serra, R., 324, 361–388, 413, 424
Simon, H., 28, 33, 415, 434
Strogatz, S., 170, 477

T

Tainter, J., 104, 108, 175–176

V

Van der Leeuw, S., 1, 11, 43, 68, 81, 85–114,
234, 317, 389–391, 481–485
Villani, M., 324, 361, 363, 364, 387, 413–432

W

Watts, D., 11, 46, 170
West G., 1, 103, 175, 221–225, 317, 433, 436,
438, 439, 443
White D., 3, 105, 107, 153–190
White, H., 181
Winter, S., 308–309
Wright, H., 162

Z

Zipf, G., 102, 156, 177–178, 201, 212, 213,

Subject Index

A

Agent, 27–28, 37–40, 267, 289–309, 381, 382
Agent based model, 324, 334, 362–363, 387,
389–410, 413–432
Artifact, 26–28, 37–40, 289–309, 366–369,
370–377, 383–385, 413–432, 456

C

Choppers, 90
City, 173–177, 180, 198–200, 202, 203, 208,
209, 211–216, 221–235, 239–242, 247,
248, 255, 336–340, 344–347, 441, 483
Cognitive capacities, 44, 49–50, 52, 62, 63,
65–67, 80–81, 88–89, 392
Competence network, 292, 294, 297, 301, 306,
307, 308, 318, 319, 322–324
Containers, 94–95
Corporation network, 160, 171–172, 298–301,
318–324, 337

D

Diffusion, 28, 31, 135, 138, 156, 168, 171,
174–175, 181–182, 189, 200, 216–217,
219, 238–241, 243, 244, 253, 254, 257,
289, 291, 298, 307, 308, 313, 353, 365, 477
Directedness, 32, 33–36, 40, 286–287, 365,
410
'Discoidal tools', 90
Distributed control, 263, 266–267, 270,
276–277, 281, 282, 286, 290–295,
297–301, 306–308
Distributed growth, 201, 204, 211, 213,
239–241, 334

E

Emergence, 4, 25, 26–27, 29, 86, 98, 100–103,
106, 110–112, 124, 130, 162, 164, 165,
205, 209, 238, 239, 241–243, 249, 254,
279, 298, 322–323, 332, 334, 335, 342,

356, 371, 383, 386–387, 414, 426, 428,
431–432, 454, 456, 482–483, 486–487
Evaluation, 30, 47, 52, 59–60, 308, 311–312,
315, 317, 318, 322, 324, 347, 393–394,
396, 401, 418–424, 427–429, 432, 446
Exaptation, 39, 40, 129–130, 413–416, 423,
425–432
Exaptive bootstrapping, 13, 37–40, 364

F

Functionality, 13, 17, 20–22, 26, 31, 32,
35, 37–40, 44–47, 49, 55, 59–60, 62,
71–72, 76–82, 112, 140–141, 200, 204,
257, 264–269, 271, 281, 286, 290,
299, 306–308, 322, 365, 393, 396,
398, 401–403, 410, 413, 415, 418–420,
422–428, 430–431

G

Generative potential, 286, 287, 318, 319, 322,
365, 379
Generative relationships, 263–309, 318, 319,
321–322, 324, 365, 410

H

Hierarchy, 25, 28, 29, 33, 46, 47, 86, 102, 128,
154–168, 172–177, 183, 188–189, 201,
204, 205, 209, 215–217, 238–239, 244,
249–250, 253–256, 274, 317, 332, 352,
357, 444–445, 482, 484
Hominins, 44, 49, 58, 60–68, 81
Hunter-gatherers, 51, 60, 61, 64–65, 72, 75,
103

I

Imitation, 47–48, 59, 62–63, 81, 338, 367–368,
456
Information, 1–3, 5, 33, 40, 46, 48, 52–54,
81, 85–86, 99–102, 104–105, 107–113,

- 137–138, 155–158, 167–168, 173, 182, 183, 188, 200, 216, 219, 222, 224, 234–235, 237, 239, 256–257, 259, 264–265, 269, 271, 285–286, 290, 293, 297–299, 302–307, 314, 316, 320, 323–325, 334, 338–339, 341, 357, 362–363, 389–410, 413–414, 418, 421–427, 430–431, 435, 438, 442–446, 448, 459, 485
- Innovation, 16–17, 20–21, 26, 43–82, 85–89, 97–100, 109–114, 153–190, 216–219, 230, 233, 237–259, 298–301, 311–325, 353–354, 361–388, 407–410, 421, 487
- Innovation cycle, 157, 188, 237–259, 331, 338, 341, 484
- ‘Innovation innovation’, 43–82, 85, 109
- Interaction, 26–27, 55, 72–76, 201–202, 307–309, 340–346, 401–402, 482–483
- K**
- Kinship, 44, 58–59, 62, 65, 67, 71–72, 75–76, 85, 87, 96, 105, 162, 164–165, 189
- L**
- Landscape, 96–97, 201, 451, 464, 466–467, 469–471
- Local governance, 52, 97, 139, 145, 163, 171, 199, 265, 300, 305, 321, 338, 339, 386, 436, 437, 468, 487
- Local system, 5, 322
- M**
- Market, 5, 32, 38, 102, 118–119, 121–124, 129–131, 133–138, 141, 200, 209, 223, 242, 254, 316–318, 342, 346, 369, 428, 437, 448, 465, 468, 482, 483
- Market system, 263–288, 290–291
- Multi-agents system, 332, 334, 337
- N**
- Narrative structure, 308, 322
- Network, 53, 105–107, 109, 153–190, 266–267, 298–301, 318–324, 337, 345–346, 433–448, 452, 457–459
- O**
- Organisation, 3–6, 481–482, 484
- P**
- Policy, 108, 112, 127, 145, 147, 188, 190, 222, 270, 289, 311–325
- Primates, 47–48, 50, 63–68, 71, 77, 82, 89
- Problem solving, 86, 97–100, 105, 111, 175–176, 189–190, 442
- R**
- Reciprocity, 27–28, 364–365
- Recursive reasoning, 44, 67–70, 80
- S**
- Scaffolding structure, 111, 129, 256–257, 291, 292, 294, 295, 298–299, 307, 308, 318, 323
- Scaling law, 5, 6, 33, 173, 222, 223, 226, 230–231, 234–235, 254, 258–259, 433–448, 483–484
- Sensitivity, 342, 348, 351, 353
- Simulation, 6, 112, 163, 324–325, 331–335, 338, 341–349, 351–353, 355, 362–364, 367–368, 373, 375–376, 383–384, 394, 403–405, 428–431, 468–469, 472–473
- Societal evolution, 197, 331
- Societal theory, 4, 44, 76, 81, 97, 110, 137, 331, 335, 351, 487
- Spatial interaction, 199, 201, 340, 341
- Stone tools, 68, 88–89, 91, 93–94, 113, 391
- Sustainability, 235, 256–259, 316, 331
- System of cities, 198, 201–203, 239, 244, 254–255, 259, 331–332, 335–336, 338–339, 342, 344, 433, 483
- T**
- Territory, 97, 107, 108, 198, 202, 204, 205, 210–212, 219, 239, 256, 259, 281, 300, 323–324, 335, 340, 345, 356–357
- U**
- Urban function, 198, 201, 224, 244, 253, 332, 336–338, 340–348, 350, 353, 358
- Urban growth, 5, 104, 175, 199, 205, 210–215, 233, 239–242, 258, 331, 333, 342–343, 435, 439, 448
- V**
- Validation methods, 155, 332, 348–351