# Multi-scale Quantum Models for Biocatalysis

## Modern Techniques and Applications

Edited by
Darrin M. York
and Tai-Sung Lee

Series Editor
Jerzy Leszczynski

$e^-$

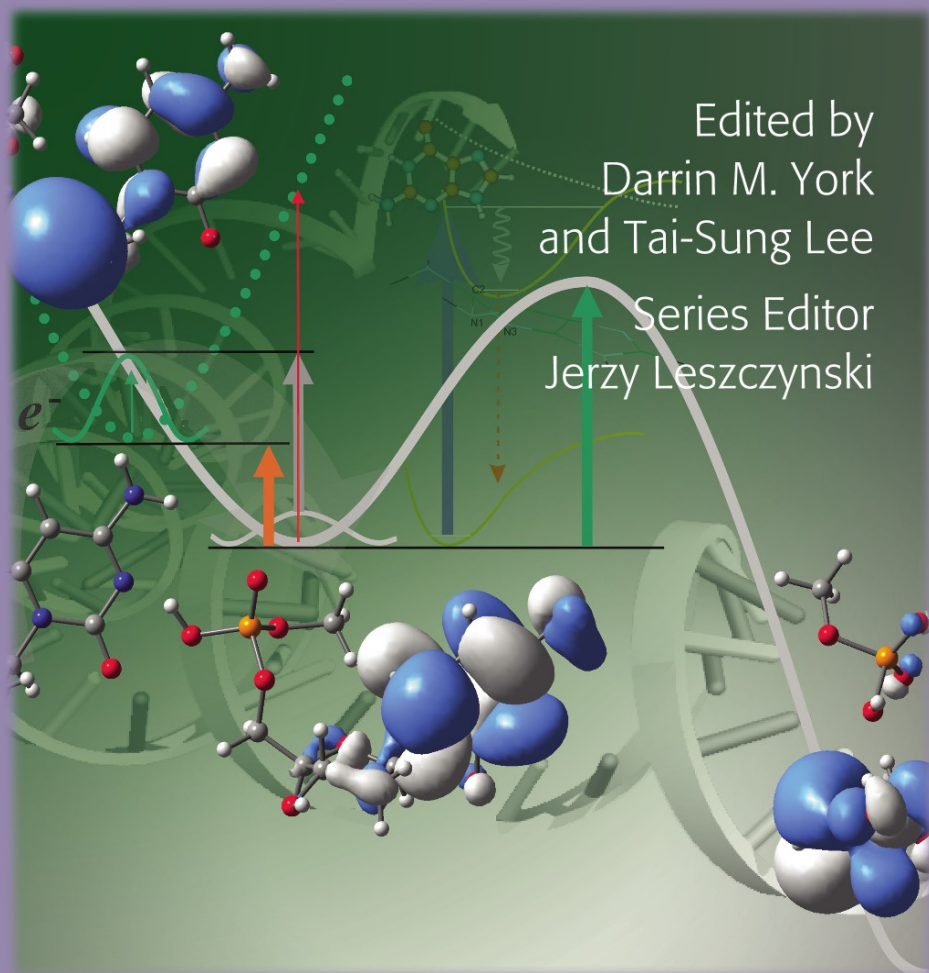Multi-scale Quantum Models for Biocatalysis

# CHALLENGES AND ADVANCES IN COMPUTATIONAL CHEMISTRY AND PHYSICS

Volume 7

# Multi-scale Quantum Models for Biocatalysis

## Modern Techniques and Applications

*Edited by*

Darrin M. York
*University of Minnesota, Minneapolis, MN, USA*

Tai-Sung Lee
*University of Minnesota, Minneapolis, MN, USA*

Springer

*Editors*
Prof. Darrin M. York
University of Minnesota
Dept. Chemistry
207 Pleasant St.
Minneapolis MN 55455
USA
york@umn.edu

Dr. Tai-Sung Lee
University of Minnesota
Dept. Chemistry
207 Pleasant St.
Minneapolis MN 55455
USA
leex2750@umn.edu

*Cover design:* Boekhorst Design BV

Printed on acid-free paper

Springer is part of Springer Science+Business Media (www.springer.com)

# CONTENTS

v

**Part III   Biocatalysis Applications**

# PREFACE

Predictive computational modeling of biological processes is one of the most significant and important present day challenges faced by the computational chemistry and biology community. The complexity of these problems demands the use of methods that are able to accurately model intra and intermolecular forces, adequately sample relevant configurational space, and establish realistic solvated environments. At one end of the spectrum, catalytic processes may require a high-level quantum chemical description, whereas at the other end of the spectrum large-scale conformational rearrangements or molecular recognition events may require very long simulation times or specialized methods for configurational sampling. These computational modeling challenges are amplified when such processes are coupled. In order to address these problems, so-called "multiscale" simulation models are required. Here, "multi-scale" implies the integration of a hierarchy of methods that span a broad range of spatial and temporal domains and work in concert to provide insight into complex problems. This book aims to provide a collection of state-of- the-art multiscale quantum simulation techniques recently developed and applied to tackle fundamental biocatalysis problems.

With 14 contribution chapters from the experts in the field, this book has three sections that group together different aspects of multiscale quantum simulations. The first section consists of four chapters that describe strategies for "multiscale" quantum models. In Chapter 1, Seabra, Swails, and Roitberg present an overview of combined quantum-classical calculations. In Chapter 2, Lundberg and Morokuma describe the generalized "ONIOM" method and its use in enzyme reactions. In Chapter 3, Cisneros and Yang examine methods for free energy simulation of enzyme catalysis with *ab initio* QM/MM methods. Chapter 4, contributed from Gao and co-workers, deals with coupled electronic structure and internuclear quantum contributions to biochemical reactions.

Section 2 mainly focuses on the current efforts to improve the accuracy of quantum calculations using simplified empirical model forms. McNamara and Hillier, in Chapter 5, summary their work on improving the description of the interactions in biological systems via their optimized semiempirical molecular models. Piquemal and co-workers present recent advances in the classical molecular methods, aiming at better reproduction of high-level quantum descriptions of the electrostatic interactions in Chapter 6. In Chatper 7, Cui and Elstner describe a different semiempir-

ical approach from the "traditional MNDO" methods, the Self-Consistent-Charge Density-Functional-Tight-Binding (SCC-DFTB) method. Chapter 8 describes a method, developed by Gordon and co-workers to obtain the approximate intermolecular potentials by merging potentials from molecular segments. Chapter 9, by Lopes et al., formulates the explicit inclusion of electronic polarizability in molecular modeling and dynamics studies.

The last section, Section 3, consists of 5 chapters focused on the applications of important biological systems. In Chapter 10, Khandogin describes modeling protonation equilibria using constant pH simulation. Quantum Mechanical studies of the photophysics of DNA and RNA bases are presented by Kistler and Matsika in Chapter 11. In Chapter 12, Zhang applies a pseudo-bond QM/MM approach to study histone modifying enzymes. In Chapter 13, Merz and co-workers rationalize the experimentally observed substrate selectivity and the product regioselectivity in Orf2-catalyzed prenylation through multiscale quantum calculations. Finally, Chapter 14 is the contribution from Lee, York and co-workers that describe a multiscale simulation strategy aimed at unraveling the mechanisms of ribozyme catalysis.

While all contributions are independent, they collectively paint a broad picture of current developments in multiscale quantum model methods at the frontier of the field. As such, this book will serve as an important reference to the scientific community. Finally, we are especially grateful to our authors who contributed excellent chapters to this volume.

<div align="right">

Tai-Sung Lee

Darrin M. York

</div>

Minneapolis, MN

# Part I
# Overview of Methodologies

# CHAPTER 1

# MIXED QUANTUM-CLASSICAL CALCULATIONS IN BIOLOGICAL SYSTEMS

GUSTAVO M. SEABRA, JASON SWAILS, AND ADRIAN E. ROITBERG

*Quantum Theory Project and Department of Chemistry, University of Florida, Gainesville, FL 32611-8435, USA, e-mail: roitberg@qtp.ufl.edu*

**Abstract:**     The development and applications of hybrid quantum mechanical/molecular mechanics calculations to the computational study of enzymatic reactions is a quickly growing field. The present chapter describes some of our group's efforts in this area over the last ten years. Increases in computational power coupled to methods for increased sampling will surely make this type of techniques a commonplace tool in the chemistry set

**Keywords:**     QM/MM methods, Jarzynski approximation, Enzyme mechanisms

## 1.1.     INTRODUCTION

Due to the large size of biological molecules, certain approximations need to be made when using computational tools for their study. It is common to use approximate classical molecular mechanics (MM) Hamiltonians coupled to molecular dynamics (MD) or Monte Carlo (MC) sampling techniques for the computational studies of such systems. These MM Hamiltonians apply parameterized force fields to describe molecular properties, greatly reducing the computational complexity of the calculation. For example, Eq. (1-1) shows the functional form for the potential energy, $U(\mathbf{R})$, as a function of the position of all atoms ($\mathbf{R}$) used by the MD program Amber [1, 2]. The terms on the right hand side represent bond, angles, dihedrals, van der Waals and electrostatic potential energies, respectively:

$$U(\mathbf{R}) = \sum_{bonds} K_r(r - r_{eq})^2 + \sum_{angles} K_\theta(\theta - \theta_{eq})^2 + \sum_{dihedrals} \frac{V_n}{2}\left[1 + \cos(n\phi - \gamma)\right]$$

$$+ \sum_{A<B}^{atoms}\left(\frac{A_{AB}}{R_{AB}^{12}} - \frac{B_{AB}}{R_{AB}^6}\right) + \sum_{A<B}^{atoms}\frac{Q_A Q_B}{\varepsilon R_{AB}}. \tag{1-1}$$

In order to run a MM calculation, the parameters $K_r$, $r_{eq}$, $K_\theta$, $\theta_{eq}$, $V_n$, $\gamma$, $A_{AB}$, $B_{AB}$, and the charges $Q_A$ must be adjusted for different molecules or residues, usually to reproduce some experimental property or quantum calculation. The term "force field" refers to a specific set of those parameters, derived to work together. Besides the program, the name Amber is also used to refer to the particular family of force fields according to Eq. (1-1), although in this case the proper reference to the force field should also include the specific version details (e.g. Amber ff99SB [3], ff03 [4], etc.).

Notwithstanding their approximated nature, these methods have provided invaluable insight to the understanding of a range of biologically relevant processes such as ligand binding [5, 6], enzyme reaction mechanisms [7], protein folding [8], refolding [9], and denaturation [10], supporting the analysis of complex experimental data and structures.

Despite the continuous effort in the development of new and more reliable force fields [3], processes essential for the study of enzymatic catalysis remain beyond the reach of classical mechanics, such as bond breaking and forming and charge fluctuations as a function of geometry [11], or describing parts of the potential energy surface far from equilibrium, such as transition states [12]. Part of the difficulties encountered by MM simulations comes from the quadratic nature of the two first terms in Eq. (1-1), which pre-empts bond breaking. Also, effects due to electronic rearrangement are neglected by fixing the value of the atomic charges. This last hurdle can be partially overcome by the use of polarizable force fields [13–15], however at an increased computational cost. In some cases, although computationally expensive, it is possible to treat a model system purely by quantum mechanics (QM) methods [16, 17], but the effect of the environment is usually either neglected or simulated by a continuum dielectric approximation.

In this contribution, we describe work from our group in the development and application of alternatives that allow the explicit inclusion of environment effects while treating the most relevant part of the system with full quantum mechanics. The first methodology, dubbed MD/QM, was used for the study of the electronic spectrum of prephenate dianion in solution [18] and later coupled to the Effective Fragment Potential (EFP) [19] to the study of the Claisen rearrangement reaction from chorismate to prephenate catalyzed by the chorismate mutase (CM) enzyme [20].

A second approach is based on the methodology first explored in the seminal work by Warshel and Levitt as early as 1976 [21], and is the use of hybrid quantum mechanics/molecular mechanics (QM/MM) calculations whereby a subsection of the system is treated by QM methods, the remainder (environment) is treated by standard molecular mechanics (MM) methods, and a coupling potential is used to connect the two regions [22]. This methodology will then be exemplified with work developed in this group in recent years [23–26].

## 1.2. MD/QM

The MD/QM methodology [18] is likely the simplest approach for explicit consideration of quantum effects, and is related to the combination of classical Monte Carlo sampling with quantum mechanics used previously by Coutinho et al. [27] for the treatment of solvent effects in electronic spectra, but with the variation that the MD/QM method applies QM calculations to frames extracted from a classical MD trajectory according to their relative weights.

In the MD/QM technique each tool is used separately, in an attempt to exploit their particular strengths. Classical molecular dynamics as a very fast sampling technique is first used for efficient sampling of the conformational space for the molecule of interest. A cluster analysis of the MD trajectory is then used to identify the main conformers (clusters). Finally QM calculations, which provide a more accurate (albeit more computationally expensive) representation of the system, can be applied to just a small number of snapshots carefully extracted from each representative cluster from the MD-generated trajectory.

### 1.2.1. Chorismate Mutase

The Shikimate pathway is responsible for biosynthesis of aromatic amino acids in bacteria, fungi and plants [28], and the absence of this pathway in mammals makes it an interesting target for designing novel antibiotics, fungicides and herbicides. After the production of chorismate the pathway branches and, via specific internal pathways, the chorismate intermediate is converted to the three aromatic amino acids, in addition to a number of other aromatic compounds [29]. The enzyme chorismate mutase (CM) is a key enzyme responsible for the Claisen rearrangement of chorismate to prephenate (Scheme 1-1), the first step in the branch that ultimately leads to production of tyrosine and phenylalanine.

Besides the obvious biological interest, chorismate mutase is important for being a rare example of an enzyme that catalyses a pericyclic reaction (the Claisen rearrangement), which also occurs in solution without the enzyme, providing a unique



*Scheme 1-1.* Conversion of chorismate to prephenate

opportunity to compare the catalyzed and uncatalyzed reactions, and has been the
focus of numerous studies, both experimental and theoretical [18, 20, 23, 26, 28–41].

### 1.2.2.        The Electronic Spectrum of the Prephenate Dianion

Although low energy structures for prephenate have been reported before [40], these
have been optimized using gas-phase quantum mechanics, and are not compatible
with the structure determined for the prephenate inside the active site of CM [41].
The first calculation of the electronic spectrum of prephenate inside the active site of
the enzyme was done by our group [18]. Using the MD/QM method described, we
were also able to obtain an electronic spectrum for prephenate in solution.

In that work, following the MD/QM methodology described, a 3 ns long (fully
classical) molecular dynamics run of prephenate in a box of TIP3P waters was used
to obtain a trajectory. Cluster analysis of the MD trajectory revealed two main con-
formers, differing mostly in the ring puckering and on the distance between the OH
and the $COO^-$ attached to the ring, as shown in Figure 1-1. The two conformers are
shown Figure 1-2, and were called "loose" (a), where the OH and $COO^-$ are inde-
pendently solvated and corresponding to a high value for the angle in Figure 1-1, and
"tight" (b), H bond between OH and $COO^-$, and corresponding to low angle value.



*Figure 1-1.* (**a**) Value of the dihedral angle, H(20)-C(6)-C(5)-H(21), plotted versus time for the duration
of the dynamics run. (**b**) Histogram of the dihedral angles found in part (**a**). Adapted from Ref. [18]

*Figure 1-2.* The two dominant conformers of prephenate in solution. Conformer (**a**) has both the OH and $COO^-$ groups solvated by the environment. Conformer (**b**) has a strong H-bond between the OH and the $COO^-$ groups. Adapted from Ref. [18]

As can be seen from the histogram in Figure 1-1(b), the "loose" conformation is preferred over the tight one, a result only possible with inclusion of solvent effects. *Ab-initio* calculations of those conformers show that, without the inclusion of solvent effects, the "tight" conformer is preferred by 7.4 kcal/mol, while the inclusion of solvent effects (with polarizable continuum model, PCM) shifts the preference towards the "loose" conformer, which becomes more stable than the "tight" one by 0.1 kcal/mol.

To identify the chromophore, the structures in the ground and first excited state were first partially optimized at the CASSCF level. Comparison of the two structures allowed the identification of the chromophore as the $CO-COO^-$ moiety. To calculate the electronic spectrum of prephenate in solution, 6 structures were then chosen from each cluster in the MD trajectory, being 2 where the $CO-COO^-$ dihedral angle corresponds to the average value and four within one standard deviation from the average (2 at each side). To guarantee better statistics, the pairs were chosen at different time frames. The electronic spectrum for each of the 12 conformers was then calculated at the CASSCF/CEP-31G level and averaged, and the intensities were obtained by weighting the results according to the relative abundance of each conformer during the classical runs.

Figure 1-3 shows a comparison of the calculated and experimental spectra. There are two high-energy transitions predicted by theory, around 190 and 240 nm. The first one is clearly visible in the experimental spectrum at the same wavelength, and roughly the same broadening. The second one was predicted to have intensity of 0.01 relative to the first peak, and is probably masked by the broad 200 nm band. Both are $\pi \rightarrow \pi^*$ transitions localized in the C=C bond in the ring and are likely to be present in most molecules found in the aromatic acid pathways, which hinders

*Figure 1-3.* Comparison between experimental and theoretically derived spectra for prephenate anion in solution. The vertical lines correspond to the theoretical spectrum for 12 conformers (3 lines for each) with intensities computed as described in the main text. The experimental spectrum is presented as a dark line (with the highest energy intensity also normalized to 1). The inset shows the near-UV absorption in greater detail. Adapted from Ref. [18]

its applicability for specific detection of prephenate. However, the theory also predicts a transition around 330 nm unique to the $CO—COO^-$ chromophore, which is restricted to prephenate and phenyl pyruvate, in the reaction catalyzed by CM (Figure 1-3, inset). The agreement between experiment and theory, including transition energies, overall shape and broadening and relative intensity of the band, is excellent, and was made possible only by the inclusion of solvent effects.

Since CASSCF calculations are necessary to obtain accurate excitation energies the examination of a large number of snapshots would be intractable. The MD/QM methodology instead allows the use of a few, carefully chosen representative structures for which to carry quantum simulations.

## 1.3.    THE EFFECTIVE FRAGMENT POTENTIAL

Another method that has been applied by our group to the study of enzymatic reactions is the Effective Fragment Potential (EFP) method [19]. The EFP method (developed at Mark Gordon's group at Iowa State University) allows the explicit inclusion of environment effects in quantum chemical calculations. The solvent, which may consist of discrete solvent molecules, protein fragments or other material, is treated explicitly using a model potential that incorporates electrostatics, polarization, and exchange repulsion effects. The solute, which can include some

number of solvent molecules as well, is treated in a fully ab initio manner, using an appropriate level of electronic structure theory.

The EFP method was used to optimize and simulate the electronic spectrum of the prephenate anion inside the chorismate mutase active site [18]. The optimization was done with prephenate treated at the RHF level and the first shell of protein residues represented by EFPs, and starting from a crystal structure of CM obtained from *B. subtilis* [41]. The optimized structure agrees well with the observed structure in the active site, still resembling the structure in Figure 1-2a but slightly more compacted towards the tight conformation, which can be expected given the substantial shielding of the repulsions between the carbonyl and carboxylate bonds by the enzyme environment. Using EFP allows the explicit inclusion of the environment in the CASSCF calculations for the spectrum. The prephenate bound to CM was predicted to absorb in the near-UV around 340 nm.

EFPs were also used as the QM part of the MD/QM method described above to study the reaction path of the CM-catalyzed rearrangement [20]. In this study, the CM active site was divided in a "chemically active" and "spectator" regions. The chemically active regions, composed of the substrate and the protein residues directly involved in the chemistry of the reaction, were optimized at the RHF/4-31G SBK level and energies recalculated with MP2 at the same geometry, while the spectator region was represented by EFPs and kept rigid. The starting structures for the quantum optimizations were obtained from the experimental X-ray structures, and also from molecular dynamics, as in the MD/QM method described above. The use of MD snapshots allowed the consideration of water molecules that do not appear in the X-ray structure because their residence time is small compared to the time scale of the experiment, but turned out to be important for the reaction. The number of water molecules in the first solvation shell of the reactant chorismate inside the CM active site was around 7–8 from the MD calculation.

We were the first group to point to an important effect of Glu-78 in the mechanisms of CM.

## 1.4.    QM/MM

The second approach described here for inclusion of environment effects is the use of hybrid quantum mechanics/molecular mechanics methods (QM/MM). In a QM/MM calculation [21, 22], the system is partitioned in two regions: A QM region, typically consisting of a relatively small number of atoms relevant for the specific process being studied, and a MM region with all the remaining atoms.

The total Hamiltonian ($\hat{H}$) for such a system is written as:

$$\hat{H} = \hat{H}^{QM} + \hat{H}^{MM} + \hat{H}^{QM/MM}, \tag{1-2}$$

where $\hat{H}^{QM}$ and $\hat{H}^{MM}$ are the Hamiltonians for the QM and MM parts of the system, and are calculated using either the QM method chosen or the usual force field

equations, respectively. The remaining term, $\hat{H}^{QM/MM}$, describes the interaction between the QM and MM parts:

$$\hat{H}^{QM/MM} = \hat{H}_{vdW}^{QM/MM} + \hat{H}_{elect}^{QM/MM} + \hat{H}_{bonds}^{QM/MM}. \qquad (1\text{-}3)$$

The first term on the right hand side of Eq. (1-3) stands for the van der Waals interaction between quantum and classical atoms

$$E_{vdW}^{QM/MM} = \sum_{\alpha}^{QM} \sum_{A}^{MM} \left[ \frac{A_{\alpha A}}{R_{\alpha A}^{12}} - \frac{B_{\alpha A}}{R_{\alpha A}^{6}} \right] \qquad (1\text{-}4)$$

and is calculated using the standard 12-6 Lennard-Jones equation and parameters derived from the force field in use for both the QM and MM atoms. It has been shown that the use of the MM parameters in this interaction does not introduce significant errors in the calculation [42].

The second term on the right hand-side of Eq. (1-3) accounts for the electrostatic interaction between classical and quantum zones, and will depend on the specifics of the QM implementation.

The final term in Eq. (1-3) becomes necessary if there are covalent bonds across the boundaries of the QM and MM subsystems. The treatment of such covalent bonds across boundaries is still the topic of active research and will not be discussed here [43–46]. The method of choice in the work reviewed here is the *link atom* approach, originally introduced by Singh and Kollman [47], which has found widespread use in QM/MM calculations with a number of variations being developed later including those by Bersuker et al. [48] and Morokuma et al. [49]. From all the methods to treat frontier covalent bonds, the link atom approach is the simplest to implement, and has been shown to give satisfactory results if used carefully [50]. In this treatment, a "link atom", which is usually – but not necessarily – a hydrogen atom, is placed along the bond between the QM and MM atoms at a suitable distance from the QM atom (e.g. ~1Å for H-link atom) to fill its valence, and which is treated as a regular QM atom by the QM calculation. The forces exerted on the link atom, as well as its charge, must later be scaled and redistributed between the QM and MM regions according to some rules that depend on the specific implementation.

This QM/MM approach was used in the implementation of an interface between the MM program Amber and the DFT QM program SIESTA (Spanish Initiative for Electronic Simulation of Thousands of Atoms) [23, 51]. The SIESTA program can use flexible basis sets consisting of linear combination of finite atomic orbitals defined in a real space grid, has been optimized to yield order-N scaling for large systems, and is extremely faster than conventional Gaussian- or plane-wave-based schemes for medium sized molecules. The nuclei and inner electrons are represented by norm-conserving pseudopotentials.

In this implementation, the QM and MM parts are combined according to the Hamiltonian described by Eq. (1-2), where the $\hat{H}^{QM}$ term is the Kohn-Sham

Hamiltonian which, in a nonlocal pseudopotential approximation, can be written as

$$\hat{H}^{QM} = T + \sum_\alpha V_\alpha^{nl} + \sum_\alpha V_\alpha^{na}(\mathbf{r}) + \delta V^H(\mathbf{r}) + V^{XC}(\mathbf{r}), \qquad (1\text{-}5)$$

where $\alpha$ is the QM atom index, $T$ is the kinetic energy operator, $V^{nl}$ is the nonlocal part of the pseudopotential, $V^{na}$ is a "neutral atom" potential that combines the local part of the pseudopotential and an isolated atom contribution, $\delta V^H$ is the Hartree potential associated with a small density fluctuation from a neutral atom reference $[\delta\rho = \rho(\mathbf{r}) - \rho^{na}(\mathbf{r})]$, and $V^{XC}$ is the exchange correlation potential.

The environment affects the charge density of the QM zone in a self-consistent way by the addition of a point charge potential ($V^{MM}$)to the Hartree potential as in:

$$\delta V_{QM-MM}^H(\mathbf{r}) = \delta V^H(\mathbf{r}) + \sum_{i=1}^{MM} V_i^{MM}(\mathbf{r}), \qquad (1\text{-}6)$$

$$V_i^{MM}(\mathbf{r}) = \begin{cases} \dfrac{q_i}{|r_i - \mathbf{r}_\alpha|}, & \text{for } |r_i - \mathbf{r}_\alpha| > R_c \\[2mm] \dfrac{q_i}{R_c}, & \text{for } |r_i - \mathbf{r}_\alpha| \leq R_c \end{cases} \qquad (1\text{-}7)$$

Here, $R_c$ is a cutoff distance which is made necessary to prevent the potential to become too steep in certain points of the grid and generate instabilities in the numerical integration, and is generally between 0.2 and 0.3Å.

This interface between Amber and SIESTA (developed by the Estrin Group in Buenos Aires, Argentina) was used to investigate the conversion of chorismate to prephenate (Scheme 1-1) [23, 26]. The reaction coordinate was taken as the combination of the distances describing the breaking and forming of bonds, $\xi = d_1 - d_2$, as depicted in Scheme 1-1. In the first study [23] the reaction path was investigated by restrained energy minimization, where a quadratic restraint term given by

$$V_R = \frac{1}{2}k(\xi - \xi_0)^2 \qquad (1\text{-}8)$$

where $k$ is an adjustable force constant and $\xi_0$ is the value of the reaction coordinate for a particular configuration is added to the potential energy. A value of 200 kcal/mol Å$^{-2}$ was used for $k$. To construct the path, an unrestricted QM/MM minimization was performed for the reactant or product, to generate an initial configuration for the reaction path, which is then mapped out by adding $V_R$ to the potential energy, with $\xi$ varying from $-2.0$ to $2.0$ Å, and performing energy minimizations at each step. The restraint energy can be subtracted to obtain the actual energy at each configuration.

This method was used to obtain potential energy profiles for the reaction in vacuum, in aqueous solution and in the enzyme environment [23]. For the enzymatic reaction, two different choices of the quantum system were considered: one where

*Scheme 1-1.* Transition state for the conversion of chorismate into prephenate. Also indicated are the Glu78 and Arg90 residues from chorismate mutase

only the substrate was considered quantum and the other where both the substrate and the charged side chains glu78 and arg90 where also included in the QM zone (Scheme 1-1), allowing the study of possible charge transfer and polarization effects in these neighbor charged residues. In all cases, the QM zone was treated with the PBE functional [52, 53] and DZVP basis sets with a pseudoatomic orbital energy of 30 meV and a grid cutoff of 150 Ry. For the classical region, the force field from Wang et al. [54] and the TIP3P model for water [55] were used.

The profiles obtained are shown in Figure 1-4. The computed activation energy in vacuum was 32.4 kcal/mol, compared to 13.8 kcal/mol in solution and 5.3 and 4.3 kcal/mol in enzyme, for the smaller and larger QM subsystems, respectively. The experimental activation energies in water and in enzyme are 20.7 [56] and 12.7 kcal/mol [57]. Although the computed activation energies are higher than the experimental ones, the catalytic effect of the enzyme, calculated as the difference between the activation energies in solution and in the enzyme environment, estimated as 7.5–8.5 kcal/mol from the calculations, is very similar to the experimental value of 8.0 kcal/mol. Analysis of the individual energy contributions revealed that the major contribution to the enzymatic activity is the stabilization of the transition state due to more favorable electrostatic interactions between the transition state and the enzyme and a minor steric compression (hence destabilization) of the substrate, which are not present in the aqueous environment [23].

One weakness of this treatment, however, is that it neglects entropic contributions. Entropic contributions were considered in the free energy profiles (FEP) calculated earlier using umbrella sampling [58] and Monte Carlo Free energy Perturbation [59], both using a QM/MM scheme and the AM1 Hamiltonian for the QM part. Our group used the same SIESTA DFT-based QM/MM method described above

*Figure 1-4.* Energy profiles for the reaction of chorismate to prephenate. **(a)** Profile in vacuum for the forward (*squares*) and reverse (*filled circles*) reactions. **(b)** Profiles for forward reaction in water (*filled circles*), and in the enzyme with only the substrate in the QM zone (*squares*) and with substrate plus chorismate mutase side chains glu78 and arg90 in the QM zone (*diamonds*)

to calculate the free energy potential for the chorismate to prephenate reaction in the enzyme environment using Multiple Steering Molecular Dynamics (MSMD [60–62]) simulations [26]. In MSMD, a time-dependent external potential is added to the time-independent Hamiltonian, $H_0(\mathbf{r})$. Choosing this perturbation as a harmonic potential with force constant $k$, and whose center $\lambda_0$ moves at a constant velocity $v$ along a chosen reaction path represented by $\lambda(r)$, the final Hamiltonian can be written as

$$H(\mathbf{r}, \lambda) = H_0(\mathbf{r}) + \frac{1}{2}k\,[\lambda(r) - \lambda_0 - vt]^2 \,. \qquad (1\text{-}9)$$

Assuming an infinite number of realizations of the process, equilibrium properties can be obtained from the non-equilibrium dynamics by the Jarzynski equality,

$$e^{-\beta \Delta G(\lambda)} = \left\langle e^{-\beta W(\lambda)} \right\rangle, \qquad (1\text{-}10)$$

where $\Delta G(\lambda)$ and $W(\lambda)$ are the free energy change and external work performed onto the system as it evolves along the reaction path, respectively, and the brackets indicate an ensemble average. In practice, a finite number of independent simulations



*Figure 1-5.* Free energy profile for the reaction from chorismate (RC $\approx$ 1.75) to prephenate (RC $\approx$ −1.75), obtained using MSMD and Jarzynski's equality and pulling speeds of 2.0 Å/ps (*red*) and 1.0 Å/ps (*green*), and using umbrella sampling (*blue*)

using a slow enough pulling speed is performed, starting from initial structures selected from a previously equilibrated ensemble.

In the MSMD calculations of Ref. [26], 20 initial structures were obtained from the final 2 ns of a pure MM molecular dynamics run of the system at the reaction coordinate corresponding to the reactant. Each structure was then equilibrated for 0.5 ps at the QM/MM level and 300 K with a time step of 0.5 fs, where the substrate was treated as QM and the rest of the system as MM. The reaction coordinate was then changed from 1.8 Å (chorismate) to $-1.8$ Å (prephenate) using a force constant of 200 kcal/(mol Å) and a constant pulling speed of 2.0 Å/ps for 15 of the structures, and 1.0 Å/ps for the remaining 5. The same process was repeated for the reverse reaction, for a total of 40 realizations of the process, and the results combined to build the profiles shown in Figure 1-5. For comparison, Figure 1-5 also shows the results from an umbrella sampling simulation with a total of 12 windows of 5 ps each, centered at snapshots taken from the constrained energy minimizations described above (from Ref. [23]). There's virtually no difference between the profiles calculated using the different methods. Although the $\Delta G^{\ddagger}$ values obtained ($\sim$8 kcal/mol) are lower than the experimental value ($\sim$15 kcal/mol) a negative entropic effect is calculated, in agreement with the experiment ($-9.1$ eu).

## 1.5.    SEMI-EMPIRICAL QM/MM IN AMBER

Starting with release 9.0, the Amber package also includes a new, redesigned native support for semi-empirical QM/MM Hamiltonians in its MD module, SANDER. Our group has been involved in the implementation of the Density Functional Tight Binding (DFTB), method, an approximation of the Density Functional Theory (DFT) based on a second order expansion of the DFT Hamiltonian [63] and its Self-Consistent-Charge variation (SCC-DFTB [64, 65]) into Amber's QM/MM [24]. The SCC-DFTB method was recently shown to reproduce MP2/cc-pVTZ geometries with accuracy similar to other semi-empirical methods such AM1 and PM3 [66].

We have recently used the Amber/SCC-DFTB method for investigating the mechanism of action of an enzyme from *Tripanosoma cruzi*, the causative agent of Chagas's Disease [67]. Expressed on the parasite's surface, the enzyme *trans*-sialidase (TcTS) transfers sialic acid (SA) from host glycoconjugates to *T. cruzi*'s surface mucins [68–70], thus providing *T. cruzi* with the means to evade the host's initial immune response. Strong experimental evidence supports this enzyme's value to the parasite's survival [71–75], and its absence in mammals makes it a promising target to treat Chagas' disease.

The suggested mechanism of TcTS is shown in Figure 1-6. A Tyr342/Glu230 pair acts as a unique nucleophile in this mechanism. Tyr342 attacks the anomeric C atom of the sialic acid (SA). Glu230 helps Tyr342 by accepting its phenolic proton. At the same time the C—O bond between the donor sugar and the SA weakens, resulting in a conformational change for SA. The acid forms a planar structure around the anomeric C, corresponding to the transition state geometry. Debate in the literature

*Figure 1-6.* The proposed mechanism of action of the *Tripanosoma cruzi*'s enzyme *trans*-sialidase (TcTS)

concerns whether the mechanism involves an oxocarbenium ion collapsing into a covalent intermediate or being stabilized via electrostatic interaction with Tyr342 [76, 77].

The coordinates involved are depicted in Scheme 1-2. Two-dimensional Potential Energy Surface (PES) scans of the $r_1/r_2$ and $r_3/r_4$ surfaces were performed where the $r$ values were restrained to values in a grid, and the geometry of the remaining atoms optimized, were followed by 10 ns windows of umbrella sampling



*Scheme 1-2.* Quantum zone used for QM.MM calculations of the reaction catalyzed by *trans*-sialidase. The relevant coordinates are shown in *red*

*Figure 1-7.* Free energy surface for the first step of the *trans*-sialidase mechanism as obtained from QM/MM calculations with SCC-DFTB

calculation starting from each of the minimized structures. The Free Energy Surface (FES) obtained by the Weighted Histogram Analysis Method (WHAM) [78–80] from the umbrella calculations for the $r_3/r_4$ scan is shown in Figure 1-7. The results support the 2-step mechanism, where the first step is a proton transfer from Asp59 to the oxygen connecting the sialic acid to the sugar chain, followed by the breaking of the glycosidic bond to form a planar carbocation (Figure 1-7), and the second step is the carbocation attack by the nucleophile pair formed by Tyr342/Glu230.

## 1.6.    CONCLUSIONS

This chapter reviewed some of our group's contributions to the development and application of QM/MM methods specifically as applied to enzymatic reactions, including the use of sequential MD/QM methods, the use of effective fragment potentials for reaction mechanisms, the development of the new QM/MM interface in Amber, as well as the implementation and optimization of the SCC-DFTB method in the Amber program. This last implementation allows the application of advanced MD and sampling techniques available in Amber to QM/MM problems, as exemplified by the potential and free energy surface surfaces for the reaction catalyzed by the *Tripanosoma cruzi* enzyme *trans*-sialidase shown here.

## ACKNOWLEDGMENTS

## REFERENCES

1.  Case DA, Cheatham TE III, Darden T, Gohlke H, Luo R, Merz KM Jr, Onufriev A, Simmerling C, Wang B, Woods RJ (2005) J Comput Chem 26:1668–1688
2.  Case DA, Darden TA, Cheatham TE I, Simmerling CL, Wang J, Duke RE, Luo R, Merz KM, Pearlman DA, Crowley M, Walker RC, Zhang W, Wang B, Hayik S, Roitberg A, Seabra G, Wong KF, Paesani F, Wu X, Brozell S, Tsui V, Gohlke H, Yang L, Tan C, Mongan J, Hornak V, Cui G, Beroza P, Mathews DH, Schafmeister C, Ross WS, Kollman PA (2006) Amber 9, University of California, San Francisco
3.  Hornak V, Abel R, Okur A, Strockbine B, Roitberg A, Simmerling C (2006) Proteins: Struct, Funct, Bioinform 65:712–725
4.  Duan Y, Wu C, Chowdhury S, Lee MC, Xiong G, Zhang W, Yang R, Cieplak P, Luo R, Lee T, Caldwell J, Wang J, Kollman P (2003) J Comput Chem 24:1999–2012
5.  Michel J, Verdonk ML, Essex JW (2006) J Med Chem 49:7427–7439
6.  Hornak V, Okur A, Rizzo RC, Simmerling C (2006) Proc Natl Acad Sci USA 103:915–920
7.  Basner JE, Schwartz SD (2005) J Am Chem Soc 127:13822–13831
8.  Simmerling C, Strockbine B, Roitberg AE (2002) J Am Chem Soc 124:11258–11259
9.  Affentranger R, Tavernelli I, et al. (2006) J Chem Theory Comput. 2:217–228
10. Caflisch A, Karplus M (1994) Proc Natl Acad Sci USA 91:1746–1750
11. Jono R, Shimizu K, Terada T (2006) Chem Phys Lett 432:306–312
12. Hermans J (2007) J Phys Chem A [This volume]
13. Pengyu Ren JWP (2002) J Comput Chem 23:1497–1506
14. Ren P, Ponder JW (2004) J Phys Chem B 108:13427–13437
15. Kaminsky J, Jensen F (2007) J Chem Theory Comput 3:1774–1788
16. Borowski T, Georgiev V, Siegbahn PEM (2005) J Am Chem Soc 127:17303–17314
17. Marothy SAD, Blomberg MRA, Siegbahn PEM (2007) J Comput Chem 28:528–539
18. Roitberg AE, Worthington SE, Holden MJ, Mayhew MP, Krauss M (2000) J Am Chem Soc 122:7312–7316
19. Day PN, Jensen JH, Gordon MS, Webb SP, Stevens WJ, Krauss M, Garmer D, Basch H, Cohen D (1996) J Chem Phys 105:1968–1986
20. Worthington SE, Roitberg AE, Krauss M (2001) J Phys Chem B 105:7087–7095
21. Warshel A, Levitt M (1976) J Mol Biol 103:227–249
22. Field MJ, Bash PA, Karplus M (1990) J Comput Chem 11:700–733
23. Crespo A, Scherlis DA, Marti MA, Ordejon P, Roitberg AE, Estrin DA (2003) J Phys Chem B 107:13728–13736
24. Seabra GM, Walker RC, Elstner M, Case DA, Roitberg AE (2007) J Phys Chem B 111:5655–5664
25. Crespo A, Marti MA, Roitberg AE, Amzel LM, Estrin DA (2006) J Am Chem Soc 128:12817–12828
26. Crespo A, Marti MA, Estrin DA, Roitberg AE (2005) J Am Chem Soc 127:6940–6941
27. Coutinho K, Saavedra N, Serrano A, Canuto S (2001) J Mol Struct THEOCHEM 539:171–179
28. Dewick PM (1998) Nat Prod Rep 15:17–58

29. Kloosterman H, Hessels GI, Vrijbloed JW, Euverink GJ, Dijkhuizen L (2003) Microbiology 149:3321–3330
30. III HLW, Hodocek M, Gilbert ATB, Gill PMW, III HFS, Brooks BR (2007) J Comput Chem 9999:NA
31. Ruiz-Pernia JJ, Silla E, Tunon I, Marti S (2006) J Phys Chem B 110:17663–17670
32. Zhang XD, Zhang XH, Bruice TC (2005) Biochemistry 44:10443–10448
33. Ranaghan KE, Ridder L, Szefczyk B, Sokalski WA, Hermann JC, Mulholland AJ (2004) Org Biomol Chem 2:968–980
34. Marti S, Andres J, Moliner V, Silla E, Tunon I, Bertran J (2004) J Am Chem Soc 126:311–319
35. Szefczyk B, Mulholland AJ, Ranaghan KE, Sokalski WA (2004) J Am Chem Soc 126:16148–16159
36. Vamvaca K, Vogeli B, Kast P, Pervushin K, Hilvert D (2004) PNAS 101:12860–12864
37. Marti S, Moliner V, Tunon I, Williams IH (2003) Org Biomol Chem 1:483–487
38. Worthington SE, Roitberg AE, Krauss M (2003) Int J Quantum Chem 94:287–292
39. Guo H, Cui Q, Lipscomb WN, Karplus M (2001) Proc Natl Acad Sci USA 98:9032–9037
40. Kast P, Tewari YB, Wiest O, Hilvert D, Houk KN, Goldberg RN (1997) J Phys Chem B 101: 10976–10982
41. Chook YM, Gray JV, Ke H, Lipscomb WN (1994) J Mol Biol 240:476–500
42. Riccardi D, Li G, Cui Q (2004) J Phys Chem B 108:6467–6478
43. Konig PH, Hoffmann M, Frauenheim T, Cui Q (2005) J Phys Chem B 109:9082–9095
44. Reuter N, Dejaegere A, Maigret B, Karplus M (2000) J Phys Chem A 104:1720–1735
45. Zhang YK, Lee TS, Yang WT (1999) J Chem Phys 110:46–54
46. Gao JL, Amara P, Alhambra C, Field MJ (1998) J Phys Chem A 102:4714–4721
47. Singh UC, Kollman PA (1986) J Comput Chem 7:718–730
48. Bersuker IB, Leong MK, Boggs JE, Pearlman RS (1997) Int J Quantum Chem 63:1051–1063
49. Maseras F, Morokuma K (1995) J Comput Chem 16:1170–1179
50. Field MJ, Albe M, Bret C, Proust-De Martin F, Thomas A (2000) J Comput Chem 21:1088–1100
51. Jos, Soler M, Emilio A, Julian DG, Alberto G, Javier J, Pablo O, Daniel S, nchez P (2002) J Phys: Condens Matter 2745
52. Perdew JP, Burke K, Ernzerhof M (1997) Phys Rev Lett 78:1396
53. Perdew JP, Burke K, Ernzerhof M (1996) Phys Rev Lett 77:3865
54. Wang J, Cieplak P, Kollman PA (2000) J Comput Chem 21:1049–1074
55. Jorgensen WL (1981) J Am Chem Soc 103:335–340
56. Andrews PR, Smith GD, Young IG (1973) Biochemistry 12:3492–3498
57. Galopin CC, Zhang S, Wilson DB, Ganem B (1996) Tetrahedron Lett 37:8675–8678
58. Marti S, Andres J, Moliner V, Silla E, Tunon I, Bertran J, Field MJ (2001) J Am Chem Soc 123: 1709–1712
59. Guimaraes CRW, Repasky MP, Chandrasekhar J, Tirado-Rives J, Jorgensen WL (2003) J Am Chem Soc 125:6892–6899
60. Jarzynski C (1997) Phys Rev Lett 78:2690–2693
61. Roitberg AE (2005) Ann Rep Comp Chem 1:103–111
62. Xiong H, Crespo A, Marti M, Estrin D, Roitberg A (2006) Theor Chem Acc: Theory, CompMod (Theor ChimActa) 116:338–346
63. Frauenheim T, Porezag D, Elstner M, Jungnickel G, Elsner J, Haugk M, Sieck A, Seifert G (1998) Mat Res Soc Symp Proc 491:91–104
64. Elstner M, Frauenheim T, Kaxiras E, Seifert G, Suhai S (2000) Computer Simulation of Materials at Atomic Level 357–376
65. Elstner M, Porezag D, Jungnickel G, Elsner J, Haugk M, Frauenheim T, Suhai S, Seifert G (1998) Phys Rev B: Condens Matter 58:7260–7268

66. Sattelmeyer KW, Tirado-Rives J, et al. (2006) J Phys Chem A 110:13551–13559
67. Chagas C (1909) Mem Inst Oswaldo Cruz 1:159–218
68. Ferguson MAJ, Murray P, Rutherford H, McConville MJ (1993) Biochem J 291:51–55
69. Frasch ACC (2000) Parasitol Today 16:282–286
70. Schenkman S, Eichinger D, Pereira MEA, Nussenzweig V (1994) Annu Rev Microbiol 48:499–523
71. Schenkman S, Jiang M-S, Hart GW, Nussenzweig V (1991) Cell 65:1117–1125
72. Costa F, Franchin G, Pereira-Chioccola VL, Ribeirao M, Schenkman S, Rodrigues MM (1998) Vaccine 16:768–774
73. Pereira-Chioccola VL, Costa F, Ribeirao M, Soares IS, Arena F, Schenkman S, Rodrigues MM (1999) Parasite Immunol 21:103–110
74. Villalta F, Smith CM, Burns JM, Chaudhuri G, Lima MF (1996) Ann N Y Acad Sci 797:242–245
75. Schenkman RPF, Vandekerckhove F, Schenkman S (1993) Infection and Immunity 61:898–902
76. Horenstein B, Yang J, Schenkman S (2002) Abstracts of Papers, 224th ACS National Meeting, Boston, MA, United States, August 18–22, 2002 BIOL-110
77. Watts AG, Damager I, Amaya ML, Buschiazzo A, Alzari P, Frasch AC, Withers SG (2003) J Am Chem Soc 125:7532–7533
78. Kumar S, Rosenberg JM, Bouzida D, Swendsen RH, Kollman PA (1995) J Comput Chem 16: 1339–1350
79. Roux B (1995) Comput Phys Commun 91:275–282
80. Kumar S, Rosenberg JM, Bouzida D, Swendsen RH, Kollman PA (1992) J Comput Chem 13: 1011–1021

# CHAPTER 2

# THE ONIOM METHOD AND ITS APPLICATIONS TO ENZYMATIC REACTIONS

MARCUS LUNDBERG AND KEIJI MOROKUMA

*Fukui Institute for Fundamental Chemistry, Kyoto University, Takano Nishihiraki-cho 34-4, Kyoto 606-8103, Japan, email: morokuma@euch4e.chem.emory.edu*

**Abstract:** ONIOM is a flexible hybrid scheme that can combine the most suitable computational methods for a given system without previous parameterization. The reason for its flexibility is that all calculations are performed on complete molecular systems, and the total energy is obtained from an extrapolation scheme. Most commonly used is the combination of a quantum mechanics and a molecular mechanics method (ONIOM QM:MM), and we describe applications of this method to several enzymatic systems, e.g., glutathione peroxidase and methylmalonyl-CoA mutase. The role of the protein is highlighted by comparing models with and without explicit inclusion of the protein matrix. We also outline future directions for the application of ONIOM to enzymes. One of the major deficiencies of QM/MM models in general, including ONIOM QM:MM, is the poor description of electrostatic interactions between the QM and the MM region. An attractive alternative to QM:MM is to take advantage of the multi-layer capability of ONIOM and design three-layer QM:QM':MM models. In this scheme QM' is a relatively fast molecular orbital method that can describe charge transfer and mutual polarization between the reactive region and the protein surroundings

**Keywords:** ONIOM, QM/MM protein environmental effects, Bacteriorhodopsin, Methane, Monooxygenase, Isopenicillin N synthase, Glutathione peroxidase, Methylmalonyl-CoA mutase, PLP-dependent b-lyase

## 2.1. INTRODUCTION

Despite advent of theoretical methods and techniques and faster computers, no single theoretical method seems to be capable of reliable computational studies of reactivities of biocatalysts. Ab initio quantum mechanical (QM) methods may be accurate but are still too expensive to apply to large systems like biocatalysts. Semi-empirical quantum methods are not as accurate but are faster, but may not be fast enough for long time simulation of large molecular systems. Molecular mechanics (MM) force field methods are not usually capable of dealing with bond-breaking and formation

pertinent with biocatalysis. Therefore hybrid methods combining different levels of theoretical methods for different parts of large molecular systems have to play a role.

We have been involved in developing the ONIOM method, a versatile hybrid scheme allowing multiple theoretical methods to be combined in multiple layers. It allows combining more than one quantum mechanical (QM) methods as well as an MM method. We have implemented the ONIOM method into a popular electronic structure package Gaussian. We have been applying the ONIOM method to a variety of problems, including thermochemistry, homogenous catalysis, solution chemistry and carbon nanotube chemistry.

In the present article, we will at first briefly overview the ONIOM methodology, with an illustration of a benchmark test of a three-layer ONIOM(QM:QM:MM) method, which we consider a method of future. Then we will review our recent studies of biocatalysis in which we used the ONIOM(QM:MM) method to examine the effects of protein environments on the mechanisms of enzymatic reactions, with an emphasis on metalloenzymes.

## 2.2.    OVERVIEW OF THE ONIOM METHOD

### 2.2.1.    ONIOM Method

ONIOM is a multi-layered hybrid method based on an extrapolation assumption, combining different levels of methods for different parts of a system [1–10]. Taking two-layer ONIOM as an example, the ONIOM(high level:low level) energy, $E_{real}^{ONIOM}$, is an approximation to the energy at the high level for the real system, $E_{real}^{high}$, referred to as the target, and is given by:

$$E_{real}^{ONIOM} = E_{model}^{high} + E_{real}^{low} - E_{model}^{low} \qquad (2\text{-}1)$$

where *high* and *low* refer to the high and low level theoretical methods, respectively, while *model* and *real* refer to the model and real systems, respectively. The model system is a part cut from the real system. The model system is mended by *link atoms* to satisfy the valencies if covalent bonds are cut. With the term "energy", we typically refer to a relative energy, such as the binding energy or the barrier height.

The accuracy of the ONIOM method is defined as the error of ONIOM relative to the target calculation:

$$Err_{real}^{ONIOM} = E_{real}^{ONIOM} - E_{real}^{high} \qquad (2\text{-}2)$$

$Err_{real}^{ONIOM}$ depends critically on the partitioning of the model system and the reliability of the low-level method used in ONIOM. The accuracy of any ONIOM combination can be tested using the S-values, which are defined as:

$$S_{\text{real-model}}^{\text{low}} = E_{\text{real}}^{\text{low}} - E_{\text{model}}^{\text{low}}$$
$$S_{\text{real-model}}^{\text{high}} = E_{\text{real}}^{\text{high}} - E_{\text{model}}^{\text{high}}$$

(2-3)

With these definitions of the S-values, the ONIOM energy can be written as:

$$E_{\text{real}}^{\text{ONIOM}} = E_{\text{model}}^{\text{high}} + S_{\text{real-model}}^{\text{low}}$$

(2-4)

Therefore:

$$\text{Err}_{\text{real}}^{\text{ONIOM}} = S_{\text{real-model}}^{\text{high}} - S_{\text{real-model}}^{\text{low}}$$
$$= \Delta S_{\text{real-model}}^{\text{high-low}}$$

(2-5)

Thus, the error of the ONIOM approximation is zero if $S_{\text{real-model}}^{\text{high}} = S_{\text{real-model}}^{\text{low}}$, namely, if the S value (the energy difference between the real and the model) at the low level is equal to the S value at the high level. This represents a situation where the effect of the surroundings on the reactive region is equal in the two methods, a condition that can be satisfied even if the two methods give different energies for the reaction itself. The S-value is therefore a useful tool in the calibration of hybrid methods.

ONIOM can combine two MO levels like ONIOM(QM:QM), which is a unique feature that is not available to QM/MM methods. However, the most popular combination is ONIOM(QM:MM), combining QM with MM. This is essentially equivalent to generic QM/MM. However, there are some subtle cancellation or double-counting differences for the case where a covalent bond is cut. For this, we refer to a detailed discussion published elsewhere [8]. QM:MM or QM/MM applications have typically been used without appropriate accuracy or S-value tests, as the benchmark full QM calculation for the real system is often impossible. In Section 2.2.2, we will examine one such test in detail.

It is very crucial to make an appropriate cut of the QM region in a QM:MM or QM/MM approach. In general, the larger the QM region, the more reliable the results, although the cost increases substantially. Our detailed examination also showed that the potential surface can be discontinuous when there is bond breaking and forming in the QM region closer than three bonds away from the MM region [8].

There are two ways of handling the interaction between the QM region and MM region; one way is to calculate electrostatic QM–MM interaction with the MM method (sometimes called mechanical embedding, or ME) and the other is to include the QM–MM interaction in the QM Hamiltonian (called electronic embedding or EE). The major difference is that in the ME scheme the QM wave function is the same in the gas phase and the electrostatic interaction is included classically, while in the EE scheme the QM wave function is polarized by the MM charges. The EE scheme is substantially more expensive than ME scheme, as the SCF iteration needs to be performed until self-consistency is achieved for QM electron distribution. Although the polarization effects are called important, as we will show later,

for structure and reactivities of metalloenzymes, we ourselves have not found a case where ME and EE make decisive differences, at least in cases where the QM model is appropriately selected.

To the authors' opinion, (quite arbitrary) charge scaling at the border of the QM and MM region is the most unpleasant feature of the QM/MM or QM:MM scheme. The MM charges in the vicinity of QM region gives excessive electrostatic interaction between QM and MM region. In the Amber MM force field, the electrostatic interaction between the atom pairs separated by one and two bonds (first and second neighbor atoms) is scaled to zero, that between the atom pairs separated by three bonds (third neighbor atoms) is scaled by 1/1.2 and the others are unscaled. This scaling is often used for ME calculations. Similar scaling is used for the QM–MM interaction Hamiltonian in EE calculations to avoid excessive electrostatic interaction and polarization, but for non-bonded interactions the full MM charge is still used. For a given computational cost, it can therefore be beneficial to use a larger QM region in ONIOM-ME, and thereby move the MM charges away from the reactive region, than to use an ONIOM-EE system where the MM charges are very close to the reactive region.

The 3-layer combination, ONIOM(QM:QM:MM), shown in Figure 2-1, is also a unique combination not available in the generic QM/MM approach and we recommend this method strongly, as the border between QM and MM regions is pretty far away form the "active" part of the system and the effects of scaling is



*Figure 2-1.* Representation of the three-layer ONIOM method (Reprinted with permission from Morokuma et al. [11]. Copyright © 2006 American Chemical Society.)

minimal; in addition, the middle layer QM method, even if a semi-empirical method is used, is fully polarizable. Because of lack of appropriate low-level QM methods, especially for transition metals, this method has not yet been tested or used very extensively. In a QM:MM calculation the bonded terms of the model system in general cancel out in the ONIOM extrapolation scheme (Eq. 2-1) since they appear in both $E_{real}^{low}$ and $E_{model}^{low}$. The results therefore are not very sensitive to the quality of the MM parameters, and ONIOM QM:MM can therefore be applied to any type of system. The situation is different when a semi-empirical QM method is used as mid or low layer. In a molecular orbital description, the effect of the surroundings depends on the orbital description of the reactive region. An accurate environmental effect thus requires that these orbitals are at least qualitatively correct, also for the low-level method.

In the next section, we show an example of test calculations of the three-layer ONIOM method.

## 2.2.2.    Benchmark Test of Three-Layered ONIOM Method

To illustrate the potential of the three-layer ONIOM method, we show results from a systematic comparison of three- and two-layer ONIOM methods with full QM benchmark calculations [11]. The system studied is a zwitterionic peptide, $NH_3^+$—$CH^nBu$—$CO$—$NH$—$CH_2$—$CO$—$NH$—$CH^nBu$—$COO^-$, and the partition scheme illustrated in Figure 2-2 is used. In this partition, both model and mid



model (dark) = $(H_3N^+)$CHL-CO-NHL , L is link hydrogen.
mid (dark+light) = $(H_3N^+)$CHBu-CO-NH-$CH_2$-CO-NHL
real (all) = $(H_3N^+)$CHBu-CO-NH-$CH_2$-CO-NH-CHBu($COO^-$)

*Figure 2-2.* The three-layer partition (B3LYP:AM1:Amber) used in the recent test (Adapted from Morokuma et al. [11]. Reprinted with permission. Copyright © 2006 American Chemical Society.)

systems have a charge of $+1$ while the real system is neutral. Three theoretical levels adopted are high-level quantum HQ: B3LYP/6-31G*, low-level quantum LQ: AM1 and MM: Amber 96. In Table 2-1 and 2-2, the combination HQ:HQ:HQ is nothing but the full HQ calculation and serves as the target calculation which ONIOM combinations are trying to approximate. Table 2-1 gives the root-mean-square (RMS) deviations of all the bond distances, bond angles and dihedral angles of the model system (excluding link hydrogens) and the real system, respectively, from those of the target calculation, One notices that both the pure Amber MM method, MM:MM:MM, and the pure semiempirical AM1 method, LQ:LQ:LQ, give very large deviations in optimized geometries of the entire peptide as well as those of the model system; among the two, Amber seems to do a little better than AM1. Among various ONIOM combinations, HQ:HQ:MM (equivalent to 2-layer HQ:MM with a large QM region) gives the smallest error. All the ONIOM combinations using HQ in the model system give rather small deviations in the model part of the geometry, which is expected but not necessarily automatic. An interesting finding is that this is true even for the geometry of the entire peptide, the real system; namely, the

*Table 2-1.* RMS errors of ONIOM optimized geometries (relative to the target calculation HQ:HQ:HQ) of the $NH_3^+ — C^nBuH — CO — NH — CH_2 — CO — NH — CH^nBu — COO^-$ system (Reprinted with permission from Morokuma et al. [11]. Copyright © 2006 American Chemical Society.)

| ONIOM combination[a] | Estimated Cost[b] | Atoms in the model system only | | | All atoms in the real system | | |
|---|---|---|---|---|---|---|---|
| | | 9 Bond lengths (Å) | 13 Bond angles (deg) | 10 Dihedral angles (deg) | 47 Bond lengths (Å) | 87 Bond angles (deg) | 98 Dihedral angles (deg) |
| HQ:HQ:HQ | 10000 | – | – | – | – | – | – |
| HQ:HQ:LQ | 200 | 0.009 | 1.34 | 9.04 | 0.014 | 1.27 | 6.32 |
| HQ:HQ:MM | 100 | 0.017 | 1.20 | 2.87 | 0.010 | 1.24 | 5.38 |
| *HQ:HQ:MM (EE)* | *500* | *0.017* | *1.19* | *2.85* | *0.010* | *1.24* | *5.38* |
| HQ:LQ:LQ | 110 | 0.012 | 1.62 | 14.27 | 0.022 | 1.70 | 8.04 |
| HQ:LQ:MM | 21 | 0.019 | 1.51 | 4.62 | 0.018 | 1.62 | 6.10 |
| HQ:MM:MM | 10 | 0.018 | 1.54 | 2.28 | 0.013 | 1.55 | 6.49 |
| *HQ:MM:MM(EE)* | *50* | *0.015* | *1.38* | *2.97* | *0.012* | *1.69* | *8.05* |
| LQ:LQ:LQ | 100 | 0.035 | 3.74 | 28.40 | 0.025 | 2.29 | 15.10 |
| LQ:LQ:MM | 11 | 0.033 | 3.77 | 23.87 | 0.022 | 2.22 | 11.62 |
| LQ:MM:MM | 2 | 0.032 | 3.87 | 17.41 | 0.018 | 2.20 | 9.33 |
| MM:MM:MM | 1 | 0.031 | 3.87 | 27.85 | 0.017 | 2.06 | 12.10 |

[a]HQ=B3LYP/6-31G*, LQ=AM1, MM=Amber. For the combinations involving the MM method, B3LYP/6-31G* RESP charges are adopted in the Amber calculation, and mechanical embedding (ME) is used unless specified as electronic embedding (EE).
[b]Very rough estimate of relative cost for a very large system, based on assumed cost: MM=$(10^{-3},$ $10^{-2}, 1)$, LQ=$(1, 10, 10^2)$, HQ=$(10, 10^2, 10^4)$ for (model,mid,real) systems for ME, respectively. For EE the time for HQ and LQ calculation was multiplied by a factor of 5 to reflect the charge-iteration process.

Table 2-2. 1-, 2- and 3-layered ONIOM calculations for the deprotonation energy (in kcal/mol) of $NH_3^+ - C^nBuH - CO - NH - CH_2 - CO - NH - CH^nBu - COO^-$ system using the optimized geometries by the respective methods (Reprinted with permission from Morokuma et al. [11]. Copyright © 2006 American Chemical Society.)

| Combination[a] | E(high/model) | E(med/mid) | E(med/model) | E(low/real) | E(low/mid) | E (ONIOM) | S(med/mid-model) | S(low/real-mid) |
|---|---|---|---|---|---|---|---|---|
| HQ:HQ:HQ | 241.70 | 245.51 | 241.70 | 322.09 | 245.51 | 322.09 | 3.81 | 76.59 |
| HQ:HQ:LQ | 244.04 | 248.18 | 244.04 | 244.96 | 175.48 | 317.66 | 4.15 | 69.48 |
| HQ:HQ:MM | 243.25 | 243.11 | 238.68 | 107.13 | 68.50 | 281.75 | 4.43 | 38.63 |
| *HQ:HQ:MM (EE)* | | *244.47* | | *69.56* | *107.14* | *282.05* | | *37.58* |
| HQ:LQ:LQ | 243.61 | 174.55 | 172.43 | 244.03 | 174.55 | 315.21 | 2.11 | 69.48 |
| HQ:LQ:MM | 239.82 | 170.29 | 168.31 | 104.93 | 67.25 | 279.47 | 1.98 | 37.68 |
| HQ:MM:MM | 238.50 | 67.46 | 70.29 | 106.68 | 67.46 | 274.89 | −2.83 | 39.22 |
| *HQ:MM:MM(EE)* | *258.13* | | *84.27* | *109.22* | | *283.09* | | |
| LQ:LQ:LQ | 158.24 | 160.56 | 158.24 | 238.08 | 160.56 | 238.08 | 2.32 | 77.52 |
| LQ:LQ:MM | 158.69 | 160.83 | 158.69 | 108.44 | 66.72 | 202.55 | 2.14 | 41.72 |
| LQ:MM:MM | 158.99 | 68.03 | 69.33 | 110.17 | 68.03 | 199.82 | −1.30 | 42.13 |

[a]See footnote a of Table 2-1.

correct geometry of the model system seems to dictate the errors in the rest of the system.

It is noted that when the MM is adopted already in the middle system (with small HQ region), i.e. in HQ:MM:MM, the electronic embedding (EE) gives a substantially better geometry than the mechanical embedding (MM), in particular in the geometry of the model system. However, when the MM is used only in the real system (with large HQ region), i.e. HQ:HQ:MM, the differences between EE and ME are negligible, because, as suggested early, the problematic boundary between QM ad MM layers is now on the outside peripheral of the middle system and is located far away from the model system, and the polarization of the QM layer due to the MM charges becomes less important.

Energies required to deprotonate the NH3+—COO— zwitterion to NH2—COO— are shown in Table 2-2. This benchmark is extremely sensitive to electrostatic effects, because it includes a change in the total charge of the system, and not necessarily representative of enzymatic reactions. The deprotonation energy of 322.09 kcal/mol at the pure B3LYP/6-31G* level is the target result which ONIOM approximations are trying to reproduce. At first we pay attention to the combinations without MM. The most expensive HQ:HQ:LQ combination, the 2-layer HQ:LQ method with large HQ region, as expected gives the smallest error of only −4.4 kcal/mol, or only 1.3% underestimation of the deprotonation energy. If one can afford a large HQ region (HQ:HQ) as well as LQ for the entire system, obviously this is an excellent approximation. The next level of approximation, HQ:LQ:LQ gives a little larger error of −6.9 kcal/mol, with smaller cost. The LQ:LQ:LQ or pure semiempirical AM1 calculation is not worth considering as this method is unable to describe the deprotonation reaction even qualitatively, with an absolute error of over 80 kcal/mol. These combinations are all quantum calculations and are likely to remain to be too expensive (See rough estimated cost in Table 2-1) in the near future as tools for exploring potential energy surfaces of reactions of most very large (>thousands of atoms) biological systems.

In most real biomolecular calculations, one will have to use MM as the lowest level method for at least a part of the very large system. Of the methods that use MM, HQ:HQ:MM with the RESP charges in the mechanical embedding (ME) has an error in deprotonation energy of −40 kcal/mol, followed by HQ:LQ:MM of −43 kcal/mol and then HQ:MM:MM of −47 kcal/mol. The HQ:HQ:MM contains a large HQ region and is expensive. The three-layered HQ:LQ:MM method, which is inexpensive because the mid layer is calculated by the inexpensive LQ method, lost only 2.3 kcal/mol over the more expensive HQ:HQ:MM. Compared to the target calculation, the major error appears when the negatively charged carboxylate is moved from a QM to an MM region. This division is not representative for a normal QM:MM or QM/MM application.

HQ:MM:MM, the ONIOM(QM:MM) with small QM region, with ME gives an error larger by 7 kcal/mol than QM:QM:MM with large QM region. The electronic embedding in QM:MM:MM reduces the error by 9 kcal/mol from the corresponding ME but is expensive because it has to iterate QM calculations to converge the

polarized charges. Thus one can conclude clearly for this example that HQ:LQ:MM is an excellent approximation to the impractical HQ:HQ:MM method and is the method of choice which improves over QM:MM:MM or the standard QM/MM with very little additional cost of semi-empirical calculation for the middle system. Again the combinations that use AM1 for highest level, LQ:LQ:MM and LQ:MM:MM, are in error over 120 kcal/mol and are not worth considering.

One notes that the results of QM:MM:MM depend sensitively on the choice of the charges used in the Amber calculation. The use of the Mulliken charges, for instance, in QM:MM:MM increases the error from −47 kcal/mol with RESP charges to −86 kcal/mol. This implies that the results will also depend sensitively on how to arbitrarily "scale" the near-border charges for the QM–MM interaction, because the problematic QM–MM boundary is very close to the reaction center.

The performance of different methods can be evaluated more systematically by examining the S-values for the deprotonation energy. Looking at the S-value between the middle and model systems, S(med/mid-model) in Table 2-2, one sees that the target S value is 3.8–4.4 kcal/mol for the HQ level, derived from HQ:HQ:HQ, HQ:HQ:LQ and HQ:HQ:MM calculations. The S-value for LQ is in the range of 2.0–2.3 kcal/mol from a variety of combinations involving LQ for the middle and model system. This implies that the LQ in this middle system introduce an error of 2–4 kcal/mol in the deprotonation energy, suggestion that AM1 in the middle system is a good choice of the method. The S-value for MM is −2.8 to −1.3 kcal/mol with the RESP charges and −36 to −38 kcal/mol with the Mulliken charges. MM with RESP is not bad at least in this region, but is very sensitive to the choice of charges.

Now we switch our attention to the S-value between the real and middle systems, S (low/real-mid). The target S-value for the deprotonation energy is 77.6 kcal/mol for the HQ level. The S-value for LQ is in the range of 69–78 kcal/mol from HQ:HQ:LQ to LQ:LQ:LQ. LQ is not bad for the real system but is too costly. The S-value for MM is 38–42 kcal mol/mol. Here we find a large source of error in deprotonation energy by using the Amber method for the real system.[1] The present results clearly show that MM, even used as the low level method in the outermost region of a large system, can introduce a substantial error in the energetics.

To summarize, we have systematically tested all possible three- and two-layer ONIOM combinations of high-level QM (HQ=B3LYP/6-31G*), low-level QM (LQ=AM1), and MM (Amber) for the deprotonation energy and structure of a test molecule, an ionic form of a peptide. We find the errors introduced in the ONIOM approximation, in comparison with the target HQ (or HQ:HQ:HQ) calculation, generally increases in the order:

---

[1] There is no well-defined way of calculating deprotonation energy with Amber. However, this does not matter as it totally cancels out in ONIOM by taking the difference between the real and model system.

HQ:HQ:HQ (target) < HQ:HQ:LQ < HQ:LQ:LQ < HQ:HQ:MM
< HQ:LQ:MM << HQ:MM:MM << LQ:LQ:LQ < LQ:LQ:MM < LQ:MM:MM

For realistic systems, the HQ calculation for the middle system and the LQ calculation for the real system can be expensive; we have no choice but using MM in the real system. The AM1 as the highest level (semiempirical QM/MM) has too large an error to be useful, while also a QM–MM boundary close to the region of the action, in ONIOM(QM:MM:MM), produces large errors.

Although its capability has not been fully exploited, we recommend the three-layer ONIOM(HQ:LQ:MM) method as the best trade-off between accuracy and computational cost. It treats the inner-most active center (small model) with a high level quantum mechanical (HQ) method. The active center plus nearby environment (middle system) is handled with a low level quantum mechanical (LQ) method, which provides a proper quantum mechanical description of the exchange as well as charge–charge interaction, can polarize the wave function of the active center, and allows charge-transfer between the active center and the environment. The real system is handled with a molecular mechanics (MM) method. The problematic boundary between the MM layer and the QM layer is sufficiently distant from the active center where the bond breaking and forming takes place, and the intrinsically arbitrary choice of charge assignment and scaling does not affect the outcome of the calculations.

## 2.3.    APPLICATIONS TO ENZYMATIC REACTIONS

### 2.3.1.    Active-Site and Protein Models

The understanding of the catalytic function of enzymes is a prime objective in biomolecular science. In the last decade, significant developments in computational approaches have made quantum chemistry a powerful tool for the study of enzymatic mechanisms. In all applications of quantum chemistry to proteins, a key concept is the active site, i.e. a local region where the chemical reactivity takes place. The concept of the active site makes it possible to scale down large enzymatic systems to models small enough to be handled by accurate quantum chemistry methods.

The following two theoretical approaches have been widely employed in the investigation of enzyme-catalyzed reactions: (1) the "active-site QM-only" approach, and (2) the hybrid Quantum Mechanics/Molecular Mechanics (QM/MM) approach. In the "active-site QM-only" approach, the enzyme's active site is modeled using from tens to hundred atoms, and treated with high-level ab-initio, DFT or semi-empirical methods. The neglected electronic effects of the surrounding protein can be approximately incorporated using a homogenous dielectric medium. These models, in the present contribution called "active-site" or "active-site QM-only" models have been particularly useful for metal-containing enzymes because the electronic and geometric structure of their active sites are dominated by the metal and its first coordination sphere [12–18].

However, the active site is only a conceptual tool and the assignment of the active-site atoms is more or less arbitrary. It is not possible to know beforehand which residues and protein interactions that will turn out to be important for the studied reaction. Hybrid QM/MM methods have been used to extend the "active site only" models by incorporating larger parts of the protein matrix in studies of enzymatic reactions [19–22]. The problem to select active-site residues appears both for active-site and QM/MM models, but in the latter, explicit effects of the surrounding protein (i.e. atoms outside the active-site selection) can at least be approximately evaluated. As this and several other contributions in this volume show, this is in many cases highly desirable.

The present chapter reviews applications in biocatalysis of the ONIOM method. The focus is on studies performed in our research group, in most cases using the two-layer ONIOM(QM:MM) approach as implemented in Gaussian [23]. The studied systems include: methane monooxygenase (MMO), ribonucleotide reductase (RNR) [24, 25], isopenicillin N synthase (IPNS) [26], mammalian Glutathione peroxidase (GPx) [27, 28], $B_{12}$-dependent methylmalonyl-CoA mutase [29] and PLP-dependent β-lyase [30]. These systems will be described in more detail in the following sections. ONIOM applications to enzymatic systems performed by other research groups will be only briefly described.

In the ONIOM(QM:MM) scheme as described in Section 2.2, the protein is divided into two subsystems. The QM region (or "model system") contains the active-site selection and is treated by quantum mechanics (here most commonly the density functional B3LYP [31–34]). The MM region (referred to as the "real system") is treated with an empirical force field (here most commonly Amber 96 [35]). The real system contains the surrounding protein (or selected parts of it) and some solvent molecules. To analyze the effects of the protein on the catalytic reactions, we have in general compared the results from ONIOM QM:MM models with "active-site QM-only" calculations. Such comparisons make it possible to isolate catalytic effects originating from e.g. the metal center itself from effects of the surrounding protein matrix.

Studying the full protein system in the same way as an "active-site QM-only" model introduces a variety of challenges. Since the conformational space of the protein is very extensive, care must be taken so that changes in protein structure are directly coupled to the reaction coordinate, and not to arbitrary conformational changes during the optimization. A conservative solution is to use similar structures for reactant and product, but even this is technically very difficult to achieve and often requires repeated iterations between reactant and product. Improved ONIOM optimization algorithms [9, 36] decrease the number of bad steps during the geometry optimization and improve the possibility of staying in the same local MM minima during a reaction step.

However, in some cases the reaction coordinate actually extends from the initial active-site selection into the protein, and the "same-configuration" solution is not adequate. One example appears in the study of Methylmalonyl-CoA mutase described below. Another drawback of a static optimization scheme is that it

neglects possible protein effects on free energies. To include these effects would require extensive configurational sampling, which has not been accomplished with the presently described DFT QM:MM potentials.

When comparing different computational approaches to enzyme systems, several different factors have to be considered, e.g., differences in high-level (QM) method, QM/MM implementation, optimization method, model selection etc. This makes it very difficult to compare different QM/MM calculations on the same system. Even comparisons with an active-site model are not straightforward. It can be argued that adding a larger part of the system into calculaton always should make the calculation more accurate. At the same time, introducing more variables to the calculation also increases the risk of artificial effects.

Computational methods are normally evaluated by benchmarks where relative energies are a dominating factor. This is not possible for applications to enzymatic systems because relatively little energetic data is available. Experimental studies can provide information about possible intermediates and reaction paths, but not about relative energies. Spin states of metal centers in proteins are in many cases available from spectroscopic data, but again with little information on relative energies. Instead, the main information comes indirectly from turnover rates (which can be converted to reaction barriers using transition state theory). The problem is that an accurate determination of a transition state in a protein system is still extremely difficult. With so many variables, cancellation of errors may lead to reasonable barriers even if the true origin of the catalytic effect is not correctly described.

In contrast to energetics, there exists a large amount of structural data, either from X-ray crystallography or NMR. For metal centers, accurate interatomic distances and local structures are also available from EXAFS. Computationally, geometries are also considerably easier to determine than transition state barriers and relatively independent of the choice of QM method (within reasonable choices). These facts make comparisons between experimental and optimized geometries rather straightforward and reasonable criteria. In a majority of our studies, we therefore start by analyzing how the QM:MM description affects the geometries of the active site.

In the end, what is unique about computational methods is their ability to describe transition states and intermediates. This is why the calculation of reaction mechanisms has achieved such a prominent position in quantum biochemistry. We will therefore spend a considerable amount of time to describe when improved active-site geometries can be expected to give important beneficial effects on reaction energies. In addition, we will try to describe how the non-bonded interactions between active site and surrounding protein affect relative energies.

The next section contains the most relevant findings from ONIOM applications to enzymatic systems performed in our group. This is followed by a discussion of the important protein effects and how this information can be used to improve the modeling of enzymatic reactions.

Whether the protein effect is considered important depends heavily on the type of process that is being studied. If the target is to understand the color-tuning effects of

different proteins involved in vision, a 5-kcal/mol protein effect constitutes more or less the entire problem. However, for transition metal enzymes the limited accuracy of the computational method does not allow for an accurate determination of barrier heights and relative energies. Instead the target is usually to discriminate between alternative reaction mechanisms. In mechanistic studies, an effect of 5 kcal/mol when going from an active-site to a QM:MM model may not change the proposed mechanism and must therefore be considered comparably less important.

## 2.3.2. Enzymatic Systems Studied with ONIOM

### 2.3.2.1. Bacteriorhodopsin

The protonated Schiff base of retinal (PSBR, see Figure 2-3) is the chromophore in a large number of light-sensitive proteins. Well-known examples are rhodopsin, which can be found in the human retina and bacteriorhodopsin, which is a light-driven proton pump in the photosynthetic system of some archae bacteria. These biological systems present interesting model challenges because the excited state of the chromophore requires advanced methods (e.g. CASSCF or CASPT2), at the same time as the color-tuning effect of the large surrounding protein must be taken into account. This modeling complexity makes PSBR an ideal target for multi-scale methods like ONIOM.

In our first ONIOM study, we showed the advantages of combining two molecular orbital (MO) methods in calculations of the chromophore itself. Compared to a full CASSCF treatment of a scaled chromophore (PSBN in Figure 2-3), a two-layer ONIOM (CASSCF:CIS) calculation where only parts of the conjugated system (PSBN8 in Figure 2-3) is included in the model system, reproduces the



*Figure 2-3.* Protonated Schiff-base of retinal (PSBR) and computational models used in ONIOM QM:QM calculations (*left*). Electrostatic effects of the surrounding protein on excitation energies in bacteriorhodopsin evaluated using TD-B3LYP:Amber (*right*). (Adapted from Vreven and Morokuma [37] (Copyright © American Institute of Physics) and Vreven et al. [38]. Reprinted with permission.)

excitation energy of the $S_0 \rightarrow S_1$ transition at one tenth of the computational cost [37]. The ONIOM partition gives reasonable results even though the excitation in the conjugated $\pi$ system in part extends beyond the model into the real system. Other alternatives for the low-level calculation are TD-HF or TD-B3LYP.

In a second study, the environmental effect of the protein was taken into account in a (TD-B3LYP:Amber) calculation [38]. In this investigation, the whole chromophore was treated by QM while the surrounding protein was treated by MM. Compared to gas phase, large effects on the excitation energy ($\sim$6 kcal/mol) could be observed inside the protein. The largest effect came from changes in chromophore geometry, but the electrostatic effects of the protein are also important. For a given geometry, the surrounding amino acids cause a blue-shift of the emission energy, as seen in Figure 2-3. The discrepancy between the experimental and the calculated emission energy (8 kcal/mol) can partly be attributed by the lack of proper polarization of the QM:MM interface. In ONIOM, this polarization effect can be included with the use of a three-layer calculation, one attractive alternative being CASSCF:CIS:Amber.

### 2.3.2.2.    Non-heme Di-Iron Enzymes: Methane Monooxygenase and Ribonucleotide Reductase

Metalloenzymes with non-heme di-iron centers in which the two irons are bridged by an oxide (or a hydroxide) and carboxylate ligands (glutamate or aspartate) constitute an important class of enzymes. Two of these enzymes, methane monooxygenase (MMO) and ribonucleotide reductase (RNR) have very similar di-iron active sites, located in the subunits MMOH and R2 respectively. Despite their structural similarity, these metal centers catalyze very different chemical reactions. We have studied the enzymatic mechanisms of these enzymes to understand what determines their catalytic activity [24, 25, 39–41].

Methane monooxygenase is a classic monooxygenase in which two reducing equivalents from NAD(P)H are utilized to split the O—O bond of $O_2$. Later, one oxygen atom is reduced to water while the second oxygen atom is incorporated into the substrate to yield methanol [42–45].

Ribonucleotide reductase is responsible for the conversion of the four biological ribonucleotides (RNA) into their corresponding deoxy forms (DNA). Although RNR is not an oxygenase during its primary catalyzed reaction (the conversion of ribonucleotides), it activates oxygen to generate a stable tyrosyl radical that is essential to the overall mechanism [46–49]. The common link between the chemistry of MMO and RNR is the activation of $O_2$ by the di-iron active site.

The resting states of the of the two enzymes are the oxidized forms $MMOH_{ox}$ and $R2_{met}$ while the reduced forms $MMOH_{red}$ and $R2_{red}$ show the highest activity towards oxygen. X-ray structures of all these states are available [50–52]. The purpose of our initial ONIOM study [24] was to understand how the protein affected the geometry of the dinuclear iron sites in MMO and RNR. We therefore optimized the structures of these enzymes using: (1) an active-site model (B3LYP/lanl2dz), (2) ONIOM2 (two-layer B3LYP:Amber) and (3) ONIOM3 (three-layer B3LYP:

*Figure 2-4.* Comparison of optimized and X-ray structures for the active site of RNR. The X-ray structure of R2$_{met}$ is superimposed on the optimized structures from "active-site QM-only" (*left*) and ONIOM2 (*middle*) models. The plot shows the quality of the optimizations evaluated as the root-mean-square deviations (in Å) compared to the X-ray structures of RNR and MMO (*right*). (Adapted from Torrent et al. [24]. Reprinted with permission. Copyright © 2002 Wiley Periodicals, Inc.)

HF/STO-3G: Amber – for R2$_{met}$ only). In all models the QM part consists of the two Fe centers and the first shell ligands of four formates, two imidazoles, and a few oxo, hydroxo, and/or aquo groups (see Figure 2-4).

When one superposes the B3LYP-optimized "active-site QM-only" structures on the experimental structures, as shown in Figure 2-4 for R2$_{met}$, there are noticeable differences in the geometry of the active site. In particular, the imidazole rings (representing histidine ligands) rotate away from their experimental positions. The difference between the QM-only optimized structure and the X-ray structure has been quantized as the root-mean-square (RMS) and maximum deviations. The results are 0.62 Å and 1.01 Å respectively. The ONIOM2 (B3LYP:Amber) model contains the active site in a QM description and a large part of the four α-helical fragments that surround the active site in an MM description. The most notable difference compared to the to the "active-site QM-only" result is that the new ONIOM2 optimized structures agree much closer with experiment. The RMS deviation is reduced from 0.62 to 0.34 Å when the surrounding protein is added (corresponding value for MMOH$_{ox}$ is from 1.02 to 0.51 Å). Figure 2-4 shows this very clearly; now all the active-site atoms are in closer vicinity of their experimental position. The local structure of the active site of these metalloenzymes is thus in part controlled by the protein environment.

We also performed optimization for R2$_{met}$ using the three-layer ONIOM3 (B3LYP: HF/STO-3G: Amber). In addition to the atoms shown in Figure 2-4, an additional 45 side-chain and backbone atoms were treated at the Hartree-Fock/STO-3G level. The resultant RMS and maximum deviations are 0.23 and 0.36 Å, respectively, compared to 0.34 and 0.52 Å for QM:MM. This indicates that the electronic effects of the protein residues, evaluated only classically in the QM:MM (B3LYP:Amber) treatment, can be further improved with the use of the ONIOM3 QM:QM:MM method.

As briefly mentioned above, the reduced form of MMO reacts with oxygen to initiate substrate oxygenation. To further analyze the protein effects on this reaction, the dioxygen-binding step was treated with two-layer ONIOM (B3LYP:Amber) [25]. The overall setup was similar to the one used for evaluating active-site geometries.

In these ONIOM calculations we optimized positions of all atoms of amino acids that had at least one atom within 7 Å from either iron center.

As expected from the results above, including the protein has significant effects on the geometry of the $O_2$-bound state. Looking at individual bond distances, one important effect is that the distance between the iron ($Fe^{2+}$) ions shortens upon including the protein, from 3.88 Å in the active-site model to 3.46 Å in the protein. In the crystal structures of $MMOH_{red}$, the Fe—Fe distance ranges from 3.26 Å to 3.31 Å for different protomers. Thus, the Fe—Fe distance calculated by ONIOM is in better agreement with experimental findings compared to active-site model calculations. Possibly more important for $O_2$ binding is that the Fe—$O_2$ and O—O distances change significantly. In the active-site model the Fe—$O_2$ distance is very long, 3.22 Å, while in the protein environment it is much shorter, 2.16 Å (see Figure 2-5). Further, the O—O distance in dioxygen increases from 1.26 Å to 1.34 Å upon insertion of the QM model into the protein environment, indicating a higher degree of O—O activation in the protein.

These geometrical observations are supported by the population analysis. Without the protein the dioxygen remains essentially un-activated, with a net charge close to zero and with spin corresponding to a pure triplet state. In the protein, the spin densities drop to 0.63–0.69 $e$ while the charge becomes $-0.4$ $e$. These changes indicate a partial electron transfer from the iron cluster, an important step in $O_2$ activation. The protein effect can be partly explained by the change in hydrogen-bond environment between active-site and ONIOM model. In the protein, the dioxygen molecule accepts hydrogen bonds from two water ligands (O(7)H$\cdots$O(14) and O(8)H$\cdots$O(13)) and these hydrogen bonds facilitate the negative charge transfer to the dioxygen molecule. In the active-site model (left in Figure 2-5), the ligands are oriented differently and these important hydrogen bonds never form. These results suggest that at the early stages of the methane monooxygenase catalytic cycle the protein effectively



*Figure 2-5.* Geometries of the $O_2$-bound state optimized using the active-site model (*left*) and an ONIOM model (*right*). Note the large differences in geometry of the two calculations, especially the hydrogen bonds donated to $O_2$ in the ONIOM model (marked in *grey*) (Adapted from Hoffman et al. [25]. Reprinted with permission. Copyright © 2004 Wiley Periodicals, Inc.)

facilitates the activation of the dioxygen molecule, and thus, plays an important role in the catalytic reaction.

### 2.3.2.3. Isopenicillin N Synthase

Oxygen activation is a central theme in biochemistry and is performed by a wide range of different iron and copper enzymes. In addition to our studies of the dinuclear non-heme iron enzymes MMO and RNR, we also studied oxygen activation in the mononuclear non-heme iron enzyme isopenicillin N synthase (IPNS). This enzyme uses $O_2$ to transform its substrate ACV to the penicillin precursor isopenicillin N [53], a key step in the synthesis of the important β-lactam antibiotics penicillins and cephalosporins [54, 55].

As in MMO and other oxygen-activated enzymes, binding of $O_2$ to the metal center is a critical step. Active-site model calculations of mononuclear non-heme iron enzymes using the B3LYP functional usually estimate $O_2$ binding to be endoergic by 10 kcal/mol or more [56]. However, such a large endoergicity does not seem consistent with the observed reactivity. In other studies, extensions of the QM treatment by including the protein in a QM/MM model gave significant effects on the non-heme Fe—$O_2$ interactions [57, 58] one example being $O_2$ binding in MMO described above [25]. We therefore modeled the $O_2$ binding step using active-site and ONIOM (B3LYP:Amber) models and compared calculated geometries and binding energies [26]. The binding energy is obtained by comparing the $O_2$-bound state (product) with enzyme and $O_2$ as two separate systems (reactant).

The active-site model (and the ONIOM "model" system) includes Fe, one aspartate and two histidine ligands, a water ligand and selected parts of the substrate (see Figure 2-6). The 2-histidine-1-carboxylate ligand theme is shared by several other non-heme iron enzymes [59]. For the protein system, we used two different



*Figure 2-6.* Real and model system used in ONIOM calculations of the reaction mechanism in isopenicillin N synthase

X-ray structures, one representing a state prior to $O_2$ binding (five-coordinate Fe) and another X-ray structure with NO bound as an analogue for $O_2$ (six-coordinate Fe) [60]. In both cases the "real" system consists of the full protein.

The most important result of the B3LYP active-site calculations is that even in the most stable state, $O_2$ remains unbound by 10.7 kcal/mol, mainly due to large entropy effects. Including the protein environment from the X-ray structure where iron is six-coordinated stabilizes dioxygen binding by 8–10 kcal/mol. In methane monooxygenase, there were large differences in the geometry of the $O_2$-bound state, with a significantly shorter Fe—$O_2$ bond in the ONIOM model. However, in IPNS no such change in geometry occurs. When the protein environment is added to the active-site model, the Fe—$O_2$ distance changes by $\leq 0.02$ Å and the O—O distance changes by $\leq 0.01$ Å. The electronic structure also looks very similar. Instead, the change in binding energy apparently comes from a different description of the reactant structure (where $O_2$ is not yet bound). As in MMO, the histidine ligands are prone to rotate away from their original positions when using "active-site QM-only" models. For the coordinatively unsaturated iron center, these artificial geometry changes become much more pronounced especially for His270 (and the substrate) (see Figure 2-7). The larger structural relaxation in the five-coordinate structure leads to an artificial stabilization of the reactant and an underestimation of the binding energy. This exaggerated flexibility of the active-site model appears even though the link atom hosts have their Cartesian coordinates frozen to their relative positions in the X-ray structure. In the ONIOM calculations both states are geometrically well described and the error is corrected. For $O_2$ binding in IPNS, the geometric effect on the binding energy is up to 6 kcal/mol, depending on binding mode of dioxygen.

To check whether including the protein simply stabilized the state most similar to the X-ray structure (i.e. if the static optimization approach was too rigid), we performed identical $O_2$ binding calculations using an X-ray structure representing the five-coordinate reactant state [60]. The results for the binding energy turned out to be almost identical (within 1 kcal/mol). It seems that independent of starting structure,



*Figure 2-7*. Origins of the increased $O_2$ binding energy in IPNS when the protein is included in an ONIOM model. (**A**) A comparison of the optimized geometries from an active-site model (*silver*) and an ONIOM protein model (*dark grey*), show that the artificial structural relaxation of the active-site model is more pronounced for the reactant state than for the product state. (**B**) Contributions to $O_2$ binding from the surrounding protein, evaluated only at the MM level (Adapted from Lundberg and Morokuma [26]. Reprinted with permission. Copyright © 2007 American Chemical Society.)

the protein environment provides the same stabilization of $O_2$. This suggests that for reactions with limited changes in the geometry of the active site, the static optimization approach method is flexible enough.

Further stabilization of the $O_2$-bound state is provided by interactions with the surrounding protein, effects evaluated at the MM level. These effects are between 4.5 and 5.0 kcal/mol and the largest contributions come from van der Waals interactions (see Figure 2-7). This is similar to the findings for $O_2$-binding in hemerythrin, where van der Waals interactions favored $O_2$ binding by as much as 6 kcal/mol [57]. Since $O_2$ does not have any interactions with the protein in the reactant, any interactions in the product state directly affects the binding energy. The active-site models neglect these long-range interactions and therefore underestimate the binding energy.

We are also studying the full substrate reaction using the same active-site and ONIOM models as described above. We started by modeling the reaction mechanism without protein and later followed the same reaction path with the surrounding protein included. In short, the reaction mechanism proposed from the active-site model is valid also when the protein has been added. Transition state barriers do not change significantly (except for the initial step that includes $O_2$ binding), but relative energies of some intermediates change significantly. One important effect of the protein is a better description of the hydrogen bond network (as in MMO). Another effect comes from strong non-bonded interactions between the protein and a water molecule that is released during the catalytic reaction. These effects will be discussed in more detail in Section 2.3.3.

### 2.3.2.4.    Glutathione Peroxidase

Glutathione peroxidases (GPx) constitute a family of selenoproteins, which demonstrates a strong anti-oxidant activity and protect cells against oxidative damage [61]. These enzymes use glutathione to reduce reactive oxygen species like hydrogen peroxide and organic peroxides. One of the peculiar features of the enzyme active site is the occurrence of a selenocysteine residue. Biochemical [61], kinetic [62] and crystallographic [63] studies indicate that this residue directly participates in the catalytic process.

Although experimental studies provide significant amounts of information regarding the structure and the catalytic activity of these enzymes, several issues concerning the structure (presence of water in the active site) and the catalytic mechanism remained unresolved. Based on the complete X-ray structure of human plasma GPx (2.9 Å resolution) [64], we performed active-site and ONIOM QM:MM calculations of structure and reaction mechanism [27, 28, 65].

The X-ray structure shows that enzyme is a tetramer, with two asymmetric units containing two dimer. Each dimer has two selenocysteine residues. From this structure we extracted a full monomer (see Figure 2-8) to use as our "real" system in the ONIOM calculations. The active-site selection includes the selenocysteine residue because it is suggested to play a critical role in the catalytic cycle. Furthermore, the residues Tyr48, Gly50, Leu51, Gln83, and Trp157 that form a part of the cage around

*Figure 2-8.* (**A**) X-ray structure of GPx dimer, (**B**) monomer used as the "real" system in ONIOM calculations, and (**C**) structure of the active site including two water molecules (Adapted from Prabhakar et al. [27]. Reprinted with permission. Copyright © 2004 American Chemical Society.)

the selenocysteine residue are also included in the active-site selection. The initially constructed system consisted of 3113 atoms with 86 atoms in the QM region. As discussed below, this system was later extended by adding two water molecules to the active site (see Figure 2-8).

More critical for the modeling of the reaction mechanism is how the geometry of the active site is reproduced. The most obvious differences between active-site and ONIOM models is that with the protein environment included, Gln83 and Trp157 largely retain their positions from the X-ray structure. The effect of the protein environment on the structure of the active site is also reflected in the RMS deviations between the optimized and the X-ray structures, which are 1.48 Å and 0.97 Å for the "active site only" and ONIOM(B3LYP/6-31G(d):Amber)-ME calculations, respectively. Using the cheaper ONIOM(HF/STO-3G:Amber)-ME method, the RMS deviation of the active-site atoms increases significantly (to 1.22 Å). Slightly surprising, the treatment of the QM–MM interactions using electronic embedding scheme, ONIOM(B3LYP/6-31G(d):Amber)-EE, also gives a larger RMS deviation (1.17 Å).

Comparisons between optimized and X-ray structures were once again made by calculating root-mean-square (RMS) deviations. When comparing all heavy atoms in the protein, the total RMS deviation is approximately 1.7 Å, irrespective of method for the model system or the ONIOM implementation (mechanical, ONIOM-ME, or electronic embedding, ONIOM-EE). The largest deviations occur for residues in the vicinity of the second monomer. Therefore, adding the second monomer to the model should improve the calculated geometries.

One reason for the relatively large RMS deviations, compared to the active sites of MMO and RNR, is that the active-site residues are not coordinated to the selenium (see Figure 2-8). The lack of a structural anchor leads to a relatively unstable active-site geometry. An alternative formulation is that the presence of a metal center with strong ligand interactions is one reason the active-site model works comparatively well for many metal enzymes.

Another important reason for the significant deviation between calculated and X-ray structures can be the low resolution (2.9 Å) of the X-ray structure. It is well known that X-ray structures may miss disordered water molecules inside the enzyme. The X-ray structure of the bovine erythrocyte GPx has a significantly higher resolution (2.0 Å) and that structure contains two water molecules in the active site [63]. Unfortunately that X-ray structure is not complete. In order to test for the presence of water molecules at the active site of the mammalian GPx, calculations were performed with two additional water molecules at the active site. This reduced the RMS deviation to from 0.97 to 0.79Å for ONIOM(B3LYP/6-31G(d):Amber)-ME and suggests the presence of water molecules also in the active site of mammalian GPx. In our investigation of the reaction mechanism, these water molecules turns out to be critical.

The next step of the investigation was to model the mechanism of hydrogen peroxide reduction by two molecules of glutathione ($H_2O_2$+ 2GSH → GS-SG + 2$H_2O$). Initially, we used an active-site model to invest different reaction pathways [65]. We then added the surrounding protein in an ONIOM calculation to analyze the effects of the protein, not only on the structure but, more importantly, also on relative energies [28]. The total catalytic reaction in GPx can be divided into three elementary reactions. The first reaction of the catalytic process, (E—SeH) + $H_2O_2$ → (E—SeOH) + $H_2O$, proceeds via a stepwise pathway with an overall barrier of 17.1 kcal/mol. This is in good agreement with the experimental barrier of 14.9 kcal/mol [66]. During the reaction, the Gln83 residue plays a key role as a proton acceptor, which is consistent with experiments [63]. The second elementary reaction, (E—SeOH) + GSH → (E—Se—SG) + $H_2O$, proceeds with a barrier of 17.9 kcal/mol. The third, and last, reaction, (E—Se—SG) + GSH → (E—SeH) + GS—SG, is initiated by the coordination of the second glutathione molecule. The calculations suggest that the amide backbone of the Gly50 residue directly participates in this reaction. The two water molecules are absolutely vital because they act as proton shuttles between the second glutathione molecule and the selenocysteine residue. This reaction proceeds with the barrier of 21.5 kcal/mol, and is suggested to be a rate-determining step of the entire GPx catalyzed reaction (See Figure 2-9).

To understand the effect of the protein on this modeled reaction mechanism, we selected the first reaction step, $H_2O_2$ reduction by a glutathione molecule for further investigations using the ONIOM (QM:MM) method [28]. The computational setup was similar to the structural study, but the effects of the additional water molecules were added from the active-site model. It is assumed that the reaction coordinate is the same as in the active-site study and no additional reaction pathways were investigated. An important point of the present ONIOM study is the full optimization of QM:MM transition states using the novel ONIOM algorithms [9].

The initial step of the reaction is coordination of the hydrogen peroxide. In the presence of the surrounding protein, Trp157 forms hydrogen bonds with $H_2O_2$, whereas in the "active-site QM-only" model Trp157 is hydrogen bonded to Gln83. Upon binding, $H_2O_2$ forms strong hydrogen bonds with Cso49 (selenocysteine), Gln83 and Gly50 residues. The binding energy is 9.2 kcal/mol. The corresponding

*Figure 2-9.* Reaction scheme for the complete catalytic cycle in glutathione peroxidase (*left*). Numbers represent calculated reaction barriers using the active-site model. The detailed potential energy diagram for the first elementary reaction, (E-SeH) + H$_2$O$_2$ → (E-SeOH) + H$_2$O, calculated using both the active-site (*dashed line*) and ONIOM model (*grey line*) is shown to the right (Adapted from Prabhakar et al. [28, 65]. Reprinted with permission. Copyright © 2005, 2006 American Chemical Society.)

binding energy in the active-site model is 6.3 kcal/mol (see Figure 2-9). In the next stage of the reaction the selenenic acid (E-Se-OH) is formed. The "active-site QM-only" study suggest that this reaction proceeds by a two-step mechanism: first formation of a selenolate anion (E—Se$^-$), and then O—O bond cleavage. Formation of the selenolate anion (E—Se$^-$) occurs via proton transfer from the Se through an oxygen atom of hydrogen peroxide to the neighboring Gln83. The computed barrier for the creation of the selenolate anion is 16.4 kcal/mol (12.8 kcal/mol in the active-site model). When comparing the transition state geometries, no major geometrical changes can be observed. In the protein, the Gln83H$^+$—O$^1$ and the Se—O$^1$H bond distances (see Figure 2-10 for labels) are longer by 0.05 and 0.03 Å respectively.

    In the second step of this reaction the O—O bond of H$_2$O$_2$ is cleaved. During this process, one hydroxyl fragment (O$^1$H) is transferred to the selenolate anion (R—Se$^-$) to form selenenic acid (R—SeO$^1$H), while simultaneously the second



*Figure 2-10.* Fully optimized ONIOM QM:MM transition states for the enzymatic reaction in glutathione peroxidase. Labels (TS-II-III and TS-III-IV) correspond to labels in the potential energy diagram in Fig. 3-7. Numbers show important bond distances in Å. (Adapted from Prabhakar et al. [28]. Reprinted with permission. Copyright © 2006 American Chemical Society.)

hydroxyl fragment ($O^2H$) accepts the previously transferred proton from Gln83 to form a water molecule. The calculated barrier for this process is 6.0 kcal/mol (7.6 kcal/mol in the active-site model). Since this step follows the 12.0 kcal/mol endothermic selenolate anion formation step (see Figure 2-9), the overall barrier for the formation of selenenic acid (E—Se—OH) becomes 18.0 kcal/mol (17.3 in the active-site model). This is in reasonable agreement with the experimentally measured barrier of 14.9 kcal/mol.

For the present reaction, the presence of surrounding protein only marginally affects the barrier (it increases by 0.7 kcal/mol). A possible reason for the small protein effects could be that in the present model, the active site is not deeply buried inside the enzyme; instead it is located on the interface of two monomers. Still, addition of the protein environment had effects on the active-site geometry. The reason this does not affect the total barrier height is that when comparing transition state and reactant, the protein effect appears to be relatively constant.

### 2.3.2.5.    Methylmalonyl-CoA Mutase

Methylmalonyl-CoA mutase (MCM) catalyzes a radical-based transformation of methylmalonyl-CoA (MCA) to succinyl-CoA. The cofactor adenosylcobalamin (AdoCbl) serves as a radical reservoir that generates the 5′-deoxyadenosine radical (dAdo•) via homolysis of the Co—C5′ bond [67]. The mechanisms by which the enzyme stabilizes the homolysis products and achieve an observed $\sim 10^{12}$-fold rate acceleration are yet not fully understood. Co—C bond homolysis is directly kinetically coupled to the proceeding hydrogen atom transfer step and the products of the bond homolysis step have therefore not been experimentally characterized.

Theoretical studies that has investigated the homolysis step in different enzymatic systems [68–70] reveal that small models comprising only the corrin ring and two ligands are insufficient and that inclusion of more amino acids are essential to stabilize the radical intermediates. Recently, a QM/MM study of the initial phase of the glutamate mutase-catalyzed reaction found a large electrostatic stabilization by the surrounding protein [70]. In our study of MCM we employed the ONIOM QM:MM approach to reveal the role of the protein in the rupture of the Co—C5' bond [29].

Four models of the MCM enzyme were used in the present study, one "active-site QM-only" model and three different ONIOM QM:MM models. The first ONIOM model "3req" represents the unreactive, open conformation of MCM (PDB ID: 3REQ) [71]. It does not contain the substrate molecule and the cobalt-carbon bond is intact. The "4req" model is based on an X-ray structure of the closed and reactive conformation of the enzyme (PDB ID: 4REQ) [71]. Results obtained with the 4req and 3req models should give deeper insight into the effect of the conformational changes that occur during substrate binding. The final ONIOM model, "4req-mca", is based on 4req but with the substrate removed from the active site. In both 4req models, the cobalt-carbon bond is broken in the X-ray structure. Comparisons of the 4req and 4req-mca models should highlight the role of the substrate in the homolysis step.

The ONIOM protein system contains the substrate, methylmalonyl-CoA, bound to the active site, the cofactor (AdoCbl) and all amino acids within a 15-Å radius from the cobalt atom. The active-site selection contains a truncated AdoCbl and the imidazole ring of its lower ligand. The QM part was calculated using the BP86 functional [31, 72] because it gives better agreement with experimental Co—C bond energies [73, 74]. This a different choice of functional compared to the other studies in the present review.

Initially relaxed potential energy scans along the cobalt-carbon bond were carried out. The "active-site QM-only" model (squares in Figure 2-11) shows a smooth energy increase that reaches about 26 kcal/mol at a Co—C5′ distance of ∼3.2 Å, Similar results are obtained for the unreactive conformation of the enzyme (3req model, triangles in Figure 2-11). For the closed form of the protein (4req-mca) a stabilization of the forming radical pair can be observed (circles in Figure 2-11). This is in agreement with earlier observations [68, 70]. These results indicate that the conformational switch following substrate binding lowers the barrier for homolysis by >10 kcal/mol.

The most interesting model is of course the 4req model with the substrate, which corresponds to the reactive state of the enzyme. Upon reaching a maximum on the potential energy surface, the 4req model shows a dramatic drop in energy, indicating discontinuity in the energy profile. This change in energy is paralleled by dramatic changes in the conformational features of the dAdo moiety (see Figure 2-12). This indicates that the cobalt-carbon bond alone becomes inadequate as the sole reaction coordinate in the region of greatest interest.



*Figure 2-11.* ONIOM protein model (*left*) with QM atoms shown as spheres and MM atoms as sticks (substrate MCA atoms are shown as tubes). The graph to the *right* shows potential energy profiles obtained by relaxed scans along the Co—C5′ bond in MCM for different computational models (see text for details) (Adapted from Kwiecien et al. [29]. Reprinted with permission. Copyright © 2006 American Chemical Society.)

*Figure 2-12.* (**A**) Overlay of the initially optimized 4req structure (*white*) with the TS structure (*dark grey*). The displayed atoms include the QM part and fragments of the MCA substrate (treated by MM). (**B**) Energy profile and stationary points for homolysis of the Co−C5′ bond in MCM calculated using the 4req model (Adapted from Kwiecien et al. [29]. Reprinted with permission. Copyright © 2006 American Chemical Society.)

The search for the true transition state was started from the points on the 4req potential energy surface that encompass the maximum energy point and the discontinuity point. We succeeded in finding a transition state for the homolysis step with one imaginary frequency that corresponds to Co—C5′ bond breaking. It is worth noting that this is the first time the convergence of the transition state of AdoCbl homolysis has been achieved despite numerous previous attempts. The intrinsic reaction coordinate (IRC) calculations ascertained that the product (labeled P in Figure 2-12) is the homolysis radical pair. On the reactant side on the IRC path we reached a stationary point, which we refer to as the intermediate (I), which differs from the initial structure (S) in the dAdo conformation. In the intermediate (I) the conformational change has caused slight Co—C5′ bond elongation, which is connected to a destabilization of the AdoCbl. The subsequent homolysis from (I) is characterized by a very small energy barrier of ~3 kcal/mol and is exothermic by ~7.5 kcal/mol as illustrated in the energy profile in Figure 2-12. Thus overall, the homolysis step has a barrier of ~10 kcal/mol and is endothermic by ~2.5 kcal/mol. In the case of MCM, inclusion of the protein environment around the active site is critical for calculating the energy profile of cobalt-carbon bond homolysis.

### 2.3.2.6.    *PLP-Dependent β-Lyase*

The final detailed example concerns the NifS CsdB protein. This enzyme belongs to a pyridoxal 5′–phosphate (PLP)-dependent family of enzymes [75, 77]. These enzymes react with L-cysteine and L-selenocysteine to generate L-alanine and

respective elemental forms of S or Se [76]. They are also suggested to participate in the formation of Fe-S clusters and play an important role in nitrogen fixation [77, 78]. The X-ray structure of the NifS CsdB protein from *E. coli* has been determined at 2.0 Å resolution [79] and we performed ONIOM QM:MM calculations to study the structure of the active site [30].

The dimeric structure of the NifS CsdB protein includes two distinct PLP-containing active sites per dimer, but the two active sites are spatially separated and apparently autonomous. We therefore used only the monomer and surrounding water molecules in the calculations. The entire model consists of 7992 atoms with 84 atoms from the active-site selected as the "model" part and treated quantum mechanically.

The full optimization of the "active-site QM-only" model gives extremely large RMS deviation (6.27 Å and 4.13 Å at the AM1 and B3LYP/6-31G(d) level respectively). Keeping the backbones of some amino acids "frozen" in their positions in the X-ray structure provides much improved RMS deviations of 2.45 Å and 1.95 Å for AM1 and B3LYP/6-31G(d) respectively. This clearly shows that keeping the backbones of key amino acids frozen retain some steric effects of the surrounding protein. However, the comparison of these results with the RMS deviations of 1.83 Å and 1.49 Å, from the optimization of the entire monomer at the ONIOM(AM1:Amber) and ONIOM(B3LYP/6-31G(d):Amber) levels once again demonstrates the importance of the protein surrounding on the structure of the enzyme active site.

### 2.3.2.7.    *Other ONIOM Studies*

In addition to the studies studied by use and presented above, we will briefly present some examples from other research groups. Due to the flexibility of the ONIOM scheme, it can be used to address many different types of problems in biocatalysis. This is by no means an extensive listing, but it illustrates main types of ONIOM applications to date.

As shown by the calculations of bacteriorhodopsin (Section 2.3.2.1), ONIOM is an excellent tool for excited-state reactions in biology. The important rhodopsin system has been studied both by TD-B3LYP:Amber [80] and CASSCF:Amber [81]. Another example of the combination of CASSCF with Amber for the surrounding protein can be found for the yellow protein [82].

A common way to benefit from the ability to combine different molecular orbital methods in ONIOM is to combine a DFT or ab-initio description of the reactive region with a semi-empirical treatment of the immediate protein environment, including up to 1000 atoms. Due to the requirement for reliable semi-empirical parameters, as discussed in Section 2.2.1, this approach has primarily been used for non-metal or Zn-enzymes. Examples include human stromelysin-1 [83], carboxypeptidase [84], ribonucleotide reductase (substrate reaction) [85], farnesyl transferase [86] and cytosine deaminase [87]. Combining two ab-initio methods of different accuracy is not common in biocatalysis applications, and one example from is an ONIOM (MP2:HF) study of catechol O-methyltransferase [88].

Some of the above mentioned studies also use two-layer ONIOM QM:MM approaches to include the full protein in an MM description. Other examples of QM:MM calculations of metal enzymes include heme oxygenase [89], nitrate reductase [90] and peptide deformylase [91]. Finally, we note that the ONIOM (HF:Amber) potential energy surface has been directly used in a molecular dynamics study (ONIOM/MD) of cytidine deaminase [92].

### 2.3.3. Important Effects of the Protein Environment – Insights from Applications

#### 2.3.3.1. *QM:MM Models Improves Calculated Structures*

In all the studied systems addition of the surrounding protein in an ONIOM model clearly improves the calculated active-site geometries. This is clearly illustrated in Figure 2-13, which shows the root-mean-square deviation between calculated and experimental structures for four of the studied enzymes.

There seems to be two major reasons for this improvement in calculated geometries; improved hydrogen-bond networks and a better description of metal coordination. Apart from these two effects, we did not find any *major* changes in active-site geometry that could be attributed to the surrounding protein.

Selection of an active-site model almost always leads to truncations of the hydrogen-bond network. Upon optimization of the active-site structure, this may lead to the formation of artificial hydrogen bonds that disrupt the structure. Freezing selected coordinates in the active-site model can prevent some of these hydrogen bonds to form. Another remedy could be to include more residues around the metal center, but larger QM models are much more expensive and there will probably still be truncated hydrogen bonds, although further away from the reaction center.



*Figure 2-13*. RMS deviations between X-ray structure and optimized geometries for active-site and ONIOM (B3LYP:Amber) models

Instead, including the protein environment in a QM:MM description gives a correct and stable hydrogen-bond pattern at a low computational cost.

The second major reason QM:MM models provide better geometries is that "active-site QM-only" models are too flexible and allows exaggerated rotations of metal ligands, especially histidines. The best illustrations of this phenomenon are given in Figure 2-4 (methane monooxygenase) and Figure 2-7 (isopenicillin N synthase). The rotation of the histidine ligands might be due to the lack of hydrogen-bonding partners to $N_\delta$-H (assuming $N_\varepsilon$ ligates to the metal).

However, in many systems even the QM:MM models show large deviations from the experimental active-site geometry (see Figure 2-13). In some cases this is due to an incomplete protein model (e.g. including one monomer of a dimeric protein as in glutathione peroxidase). Another reason can be remaining errors in the X-ray structures (i.e. the two missing water molecules in glutathione peroxidase). However, from a computational perspective it is interesting that use of an additional (QM) buffer layer between QM and MM atoms, as in the three-layer (B3LYP:HF:Amber) calculations of ribonucleotide reductase, further improves the calculated geometries. This suggests the usefulness in the three-layer ONIOM method to improve the interactions between the active site and the surrounding atoms.

### 2.3.3.2.    *Correlations Between Structures and Energies*

In many cases the most interesting results of a computational study are the relative energies of transition states and intermediates because they determine the reaction mechanism. In this section we will try to outline when improved active-site geometries can be expected to have important effects on relative energies.

After a correct description of the electronic state, the most important factor is probably a correct description of the hydrogen bonds. In the present paragraph we provide two illustrations of when the protein's effect on the hydrogen-bond network are important for the enzymatic reaction. The first example is a rather obvious formation of an artificial hydrogen bond in isopenicillin N synthase. The iron-bound dioxygen species (discussed above) abstracts a hydrogen atom from the substrate and forms an iron-bound peroxide. In the active-site model, this peroxide is stabilized by a hydrogen bond to the carboxylate group (see Figure 2-14) and this interaction is formed even though the coordinates of the carbon atom of the carboxylate are kept frozen. When the protein is included, the carboxylate instead stays hydrogen bonded to other residues and do not interact with the peroxide. The effect on the relative energy of the peroxide intermediate is 5 kcal/mol.

The second illustration related to hydrogen bonds is less obvious, but appeared for the $O_2$ binding step in methane monooxygenase (Section 2.3.2.2). In the active-site model, the dioxygen molecule does not receive any hydrogen bonds (see Figure 2-5). In the protein model, the hydrogen-bond network is different and two water ligands are able to donate hydrogen bonds to the dioxygen. These hydrogen bonds helps to polarize $O_2$ and are thus at least partly responsible for the difference in electronic structure between the two models. However, looking only at the active-site model,

*Figure 2-14.* Illustration of the different hydrogen bonding patterns for an iron-bound peroxide in IPNS using an active-site model (*left*) and an ONIOM QM:MM model (*right*)

there are no specific interactions that can be assigned as "artificial", and it is therefore not obvious that the hydrogen-bond network should be incorrectly described.

A rotation of a metal ligand out of its X-ray position does not seem to affect the electronic structure. For many reactions the energetic effect of the exaggerated movements in active-site models is small, because the associated relaxation energy remains constant during the reaction. However, in reactions where the local coordination environment change, like a transition from a five to a six-coordinate structure, the energetic effect can be significant. In isopenicillin N synthase the rotation of one of the histidine ligands (and the substrate) lead to an artificial stabilization of the coordinatively unsaturated reactant structure by up to 6 kcal/mol. Adding larger parts of the protein restricts these movements and should improve the modeling of all reactions all reactions that describe changes in metal coordination, e.g. substrate binding, product release, and ligand exchange reactions.

### 2.3.3.3. Energetic Effects of Non-bonded Interactions Between Protein and Active Site

Including the protein in a QM/MM description makes it possible to evaluate the effects of non-bonded interactions between active site and surrounding protein. As mentioned above, in the ONIOM implementation electrostatic effects of the protein environment can either be evaluated classically (mechanical embedding, ONIOM-ME) or semi-classically by including the surrounding charges into the QM Hamiltonian (electronic embedding, ONIOM-EE). In the present investigations, the most common choice has been to use the computationally cheaper ONIOM-ME method. This makes it possible to use a larger QM part, and thus better evaluate local electronic effects, possibly at the expense of a less accurate evaluation of the long-range electrostatic effects.

Long-range electrostatic effects are sometimes proposed to be important for enzymatic activity. We have therefore compared results for mechanical and electronic embedding in some of the present projects. For $O_2$ binding in isopenicillin N synthase, the changes when applying ONIOM-EE instead of ONIOM-ME are relatively small. Key distances (i.e. Fe—O and O—O distances) change by up to 0.03 Å, spin and charge populations change by up to 0.1 $e$, while the binding energy changes by only 0.4—0.5 kcal/mol. For methane monooxygenase, the $O_2$ binding energy was not calculated, but changes in electronic structure (i.e. spin and charge populations) are similar to what was observed for IPNS. The major geometric effect in MMO concerns the Fe—Fe bond distance that goes from 3.46 Å in ONIOM-ME to 3.53 in ONIOM-EE. Compared to the results from X-ray (3.26–3.31 Å) electronic embedding actually performs worse, although both implementations show significant deviations. This may indicate remaining shortcomings in either model or experiment that prevent a fair comparison. Still, both ONIOM calculations perform significantly better than the active-site model (Fe—Fe distance of 3.88 Å). A comparison of the RMS deviations compared to X-ray geometries in the active site of glutathione peroxidase also gives larger deviations for ONIOM-EE (1.17 Å) compared to ONIOM-ME (0.97 Å). In other systems, we have not noted any major change in geometries between the two different ONIOM implementations.

When looking at the reaction mechanisms of glutathione peroxidase and isopenicillin N synthase, we did not find any reaction step where the transition state is significantly stabilized by long-range electrostatic interactions (i.e. electrostatic interactions outside the active-site model). However, it is should be added that most transition states have been calculated using ONIOM-ME.

Still, there are instances where we see large non-bonded environmental effects on relative energies. A clear example is provided by the potential energy profile for glutathione peroxidase (Figure 2-9). The formation of a water molecule from hydrogen peroxide is 13 kcal/mol more exothermic in the protein model compared to the active-site model. Similar results are obtained for water formation in isopenicillin N synthase. When a polar molecule is released from the active site and moves to the periphery of the QM part, it forms strong interactions (hydrogen bonds) with residues only included in the QM:MM description. Since these interactions are neglected in an "active-site QM-only" model, this leads to large protein effects on relative energies. In these steps we also find large differences between ONIOM-ME and ONIOM-EE. However, when studying reaction mechanisms, these effects are usually not very important because the reactions tend to be highly exothermal (as in GPx). The relative energy of the intermediate does therefore not affect the barrier of the next step.

Finally we point out that binding processes (as $O_2$ binding in isopenicillin N synthase) are generally much more affected by non-bonded interactions than other types of reactions. The reason is that when a new molecule is added to the model, all interactions with the surrounding protein affect the binding energy. For $O_2$ binding in IPNS the total effect of the non-bonded interactions were 4–5 kcal/mol (mainly van der Waals interactions, see Figure 2-7). During the following reaction these non-bonded interactions do not change significantly and thus have a very limited effect on the calculated transition state barriers.

### 2.3.3.4.    *More Complete Reaction Coordinates in Larger Models*

In the present investigations, there was one enzymatic system (methylmalonyl-CoA mutase) where the extension of the active-site model in an ONIOM QM:MM description significantly decreased the reaction barrier. Homolytic cleavage of the Co—C bond is stabilized by a rotation of the ribose ring in the dAdo cofactor, but this stabilization could only be observed when modeling the reactive closed configuration of the enzyme (see Figures 2-11 and 2-12). The reason the active-site model fails to describe the reaction is that it does not contain all the important reaction coordinates. When the same QM model is used in the ONIOM QM:MM description, the reaction coordinate extends into the MM region. For the present reaction, the relative energies are still reasonably described because the changes in the MM part are limited to a few rotations and torsions. Thus the QM:MM method is a cost-effective way to accurately describe the methylmalonyl-CoA-mutase system.

In a broader perspective, any unknown reaction coordinate can potentially include a significant number of atoms. For example, it cannot be ruled out that some enzymatic reactions are coupled to large changes in protein structure. With static optimization techniques such large changes are unlikely to be discovered, unless they are directly included in the transition state search. In addition, the relative energies of different protein conformations should then be accurately balanced against the relatively small energy differences between transition state and reactant.

Another example of more complicated reaction coordinates are local changes in the active-site environment, including but not limited to, water molecules leaving or entering the active site. Water molecules are the most flexible species in a protein structure and play many important roles. They may act as bridges in proton transfer reactions, stabilize charge distributions appearing during a reaction or constitute an integral part of important hydrogen-bond networks. There are several examples of the importance of water molecules in enzymatic reactions, and we provide just two examples. In glutathione peroxidase the presence of additional active-site water molecules is required both for stabilization of transition states and to correctly describe proton transfer pathways [65]. In methane monooxygenase, a single water molecule (coupled to a change in coordination mode of one of the bidentate carboxylate ligands) affects relative energies of an $Fe(IV)$—$(O^{2-})_2$—$Fe(IV)$ intermediate by 17.5 kcal/mol [39].

## 2.4.    CONCLUSIONS AND PERSPECTIVE

Despite the availability of fast computers and efficient codes for accurate quantum chemistry calculations, it is not likely in the near future that we will be able to study chemical reactions in proteins taking all the proteins atoms into quantum mechanical calculations. Hybrid methods in which different parts of large molecular systems are treated by different theoretical levels of methods are likely to play a key role in such studies for the coming decade or more. The ONIOM method we have developed is a versatile hybrid method that allows combining different quantum mechanical methods as well as molecular mechanics method in multiple layers, some features of

which are not available in generic QM/MM approaches. Concerned with unpleasant connection problems between QM and MM regions in ONIOM QM:MM or QM/MM approaches, we consider a three-layered ONIOM(QM:QM:MM) method as the method of the future, in which the QM-MM connection is kept far away from the region where the reaction takes place and the intermediate QM region (treated typically by a semiempirical QM method) is fully polarizable and allows charge transfer.

We have presented several examples of enzymatic reactions we have investigated using the ONIOM method, in particular studies of metalloenzymes with the ONIOM QM:MM method. When exploring mechanisms of metalloenzymes, the demand for accuracy is relatively low, and active-site models often provide good guesses for the reaction coordinate. Still, our results clearly show that including the surrounding protein leads to more accurate geometries, and from that observation, we make the logical conclusion that improved geometries also leads to more accurate energies. By comparing different modeling approaches, we can identify the most important protein effects: better descriptions of hydrogen bond patterns, a more complete description of the reaction coordinate, correct metal-ligand coordination, and non-bonded interactions with the surrounding protein. Depending on the purpose of the study, and the type of reaction, these effects can be anything from useful corrections to highly critical.

As we discussed, what is missing in our present studies is dynamics of the protein to cover wider range of geometrical changes and also to include statistical averaging of the complex motion of protein environment. At present, we are developing dynamical methods that can be used to better describe complicated geometry changes that might occur during the reaction, while still maintaining the flexibility and accuracy provided by the ONIOM scheme. Such a dynamical implementation may also help to evaluate protein contributions to free energies.

## ACKNOWLEDGMENTS

## REFERENCES

1. Svensson M, Humbel S, Froese RDJ, Matsubara T, Sieber S, Morokuma K (1996a) J Phys Chem 100:19357–19363
2. Svensson M, Humbel S, Morokuma K (1996b) J Chem Phys 105:3654–3661
3. Humbel S, Sieber S, Morokuma K (1996) J Chem Phys 105:1959–1967
4. Dapprich S, Komáromi I, Byun KS, Morokuma K, Frisch MJ (1999) J Mol Struct (THEOCHEM) 461–462:1–21

5. Vreven T, Morokuma K (2000a) J Comput Chem 21:1419–1432
6. Morokuma K, Musaev DG, Vreven T, Basch H, Torrent M, Khoroshun DV (2001) IBM J Res Dev 45:367–395
7. Morokuma K (2002) Philos Trans R Soc London, Ser A 360:1149–1164
8. Vreven T, Byun KS, Komáromi I, Dapprich S, Montgomery JA Jr, Morokuma K, Frisch MJ (2006a) J Chem Theor Comp 2:815–826
9. Vreven T, Frisch MJ, Kudin KN, Schlegel HB, Morokuma K (2006b) Mol Phys 104:701–714
10. Vreven T, Morokuma K (2006c) In: Spellmeyer D (ed) Annual Report Computational Chemistry, Vol. 2, Elsevier, pp 35–52
11. Morokuma K, Wang Q, Vreven T (2006) J Chem Theor Comp 2:1317–1324
12. Basch H, Mogi K, Musaev DG, Morokuma K (1999) J Am Chem Soc 121:7249–7256
13. Musaev DG, Basch H, Morokuma K (2002) J Am Chem Soc 124:4135–4148
14. Baik MH, Newcomb M, Friesner RA, Lippard SJ (2003) Chem Rev 103:2385–2420
15. Siegbahn PEM, Blomberg MRA (2004) Chem Rev 100:421–437
16. de Visser SP, Kumar D, Cohen S, Shacham R, Shaik S (2004) J Am Chem Soc 126:8362–8363
17. Noodleman L, Lovell LT, Han W, Li J, Himo F (2004) Chem Rev 104:459–508
18. Shaik S, Kumar D, de Visser SP, Altun A, Thiel W (2005) Chem Rev 105:2279–2328
19. Gao J (1996) Rev Comput Chem 7:119–185
20. Schoeneboom JC, Cohen S, Lin H, Shaik S, Thiel W (2004) J Am Chem Soc 126:4017–4034
21. Friesner RA, Guallar V (2005) Annu Rev Phys Chem 56:389–427
22. Klähn M, Braun-Sand S, Rosta E, Warshel A (2005) J Phys Chem B 109:15645–15650
23. Frisch MJ, Trucks GW, Schlegel HB, Scuseria GE, Robb MA, Cheeseman JR, Montgomery Jr JA, Vreven T, Kudin KN, Burant JC, Millam JM, Iyengar SS, Tomasi J, Barone V, Mennucci B, Cossi M, Scalmani G, Rega N, Petersson GA, Nakatsuji H, Hada M, Ehara M, Toyota K, Fukuda R, Hasegawa J, Ishida M, Nakajima T, Honda Y, Kitao O, Nakai H, Klene M, Li X, Knox JE, Hratchian HP, Cross JB, Bakken V, Adamo C, Jaramillo J, Gomperts R, Stratmann RE, Yazyev O, Austin AJ, Cammi R, Pomelli C, Ochterski JW, Ayala PY, Morokuma K, Voth GA, Salvador P, Dannenberg JJ, Zakrzewski VG, Dapprich S, Daniels AD, Strain MC, Farkas O, Malick DK, Rabuck AD, Raghavachari K, Foresman JB, Ortiz JV, Cui Q, Baboul AG, Clifford S, Cioslowski J, Stefanov BB, Liu G, Liashenko A, Piskorz P, Komaromi I, Martin RL, Fox DJ, Keith T, Al-Laham MA, Peng CY, Nanayakkara A, Challacombe M, Gill PMW, Johnson B, Chen W, Wong MW, Gonzalez C, Pople JA (2004) Gaussian 03, Revision C.02, Gaussian Inc Wallingford CT
24. Torrent M, Musaev DG, Basch H, Morokuma K (2002a) J Comput Chem 23:59–76
25. Hoffmann M, Khavrutskii IV, Musaev DG, Morokuma K (2004) Int J Quant Chem 99:972–980
26. Lundberg M, Morokuma K (2007) J Phys Chem B 111: 9380–9389
27. Prabhakar R, Musaev DG, Khavrutskii IV, Morokuma K (2004) J Phys Chem B 108:12643–12645
28. Prabhakar R, Vreven T, Frisch MJ, Morokuma K, Musaev DG (2006) J Phys Chem B 110: 13608–13613
29. Kwiecien RA, Khavrutskii IV, Musaev DG, Morokuma K, Banerjee R, Paneth P (2006) J Am Chem Soc 128:1287–1292
30. Prabhakar R, Morokuma K, Musaev DG (2005a) J Comp Chem 26:443–446
31. Becke AD (1988) Phys Rev A 38:3098–3100
32. Becke AD (1993) J Chem Phys 98:1372–1377
33. Becke AD (1993) J Chem Phys 98:5648–5652
34. Lee CT, Yang WT, Parr RG (1988) Phys Rev B 37:785–789
35. Cornell WD, Cieplak P, Bayly CI, Gould IR, Merz KM Jr, Ferguson DM, Spellmeyer DC, Fox T, Caldwell JW, Kollman PA (1995) J Am Chem Soc 117: 5179–5197 (1995)

36. Vreven T, Morokuma K, Farkas Ö, Schlegel HB, Frisch MJ (2003b) J Comp Chem 24:760–769
37. Vreven T, Morokuma K (2000b) J Chem Phys 113:2969–2975
38. Vreven T, Morokuma K (2003a) Theor Chem Acc 109:125–132
39. Torrent M, Musaev DG, Basch H, Morokuma K (2001a) J Phys Chem B 105:8452
40. Torrent M, Mogi K, Basch H, Musaev DG, Morokuma K (2001b) J Phys Chem B 105:8616
41. Torrent M, Vreven T, Musaev DG, Morokuma K, Farkas Ö, Schlegel HB (2002b) J Am Chem Soc 124:192–193
42. Wilkins PC, Dalton H (1994) Biochem Soc Trans 22:700–704
43. Liu KE, Lippard SJ (1995) J Adv Inorg Chem 42:263–289
44. Lee SY, Lipscomb JD (1999) Biochemistry 38:4423–4432
45. Dunietz BD, Beachy MD, Cao YX, Whittington DA, Lippard SJ, Friesner RA (2000) J Am Chem Soc 122:2828–2839
46. Stubbe J (1990) J Biol Chem 265:5329–5332
47. Nordlund P, Sjöberg BM, Eklund H (1990) Nature 345:593–598
48. Reichard P (1993) Science 260:1773–1777
49. Mulliez E, Fontecave M (1999) Coord Chem Rev 775:185–186
50. Nordlund P, Eklund H (1993) J Mol Biol 232:123–164
51. Logan DT, Su XD, Åberg A, Rengström K, Hajdu J, Eklund H, Nordlund P (1996) Structure 4: 1053–1064
52. Whittington DA, Lippard SJ (2001) J Am Chem Soc 123:827–838
53. Baldwin JE, Bradley M (1990) Chem Rev 90:1079–1088
54. Schenk WA (2000) Angew Chem Int Ed 39:3409–3411
55. Andersson I, Terwisscha van Scheltinga AC, Valegård K (2001) Cell Mol Life Sci 58:1897–1906
56. Bassan A, Borowski T, Siegbahn PEM (2004) Dalton Trans 20:3153–3162
57. Wirstam M, Lippard SJ, Friesner RA (2003) J Am Chem Soc 125:3980–3987
58. Nemukhin AV, Grigorenko BL, Topol IA, Burt SK (2006) Int J Quant Chem 106:2184–2190
59. Koehntop KD, Emerson, JP, Que L Jr (2005) J Biol Inor Chem 10:87–93
60. Roach PL, Clifton IJ, Hensgens CMH, Shibata N, Schofield CJ, Baldwin JE (1997) Nature 387:827–830
61. Flohé L (1989) In: Dolphin D, Avramovic O, Poulson R (eds) Glutathione John Wiley & Sons, NewYork pp 644–731
62. Flohé L, Loschen G, Gunzler WA, Eichele E (1972) Hoppe Seyler's Z Physiol Chem 353:987–999
63. Epp O, Ladenstein R, Wendel A (1983) Eur J Biochem 133:51–69
64. Ren B, Huang W, Åkesson B, Ladenstein R (1997) J Mol Biol 268:869–885
65. Prabhakar R, Vreven T, Morokuma K, Musaev DG (2005b) Biochemistry 44:11864–11871
66. Roy G, Nethaji M, Mugesh G (2004) J Am Chem Soc 126:2712–2713
67. Banerjee R (2003) Chem Rev 103:2081–2081
68. Dölker N, Maseras F, Siegbahn PEM (2004) Chem Phys Lett 386:174–178
69. Freindorf M, Kozlowski PM (2004) J Am Chem Soc 126:1928–1929
70. Jensen KP, Ryde U (2005) J Am Chem Soc 127:9117–9128
71. Mancia F, Evans PR (1998) Structure 6:711–720
72. Perdew JP (1986) Phys Rev B 33:8822
73. Jensen KP, Ryde U (2003) J Phys Chem A 107:7539–7545
74. Kuta J, Patchkovskii S, Zgierski MZ, Kozlowski PM (2006) J Comput Chem 27:1429–1437
75. Fujii T, Maeda M, Mihara H, Kurihara T, Esaki N, Hata Y (2000) Biochemistry 39:1263–1273
76. Mihara H, Fujii T, Kato S, Kurihara T, Hata Y, Esaki NJ (2002) Biochemistry 131:679–685
77. Zheng L, White RH, Cash VL, Dean DR (1994) Biochemistry 33:4714–4720

78. Clausen T, Kaiser JT, Steegborn C, Huber R, Kessler D (2000) Proc Natl Acad Sci 97:3856–3861
79. Lima CDJ (2002) Mol Biol 315:1199–1208
80. Gascon JA, Batista VS (2004) Biophys J 87: 2931–2941
81. Blomgren F, Larsson S (2005) J Phys Chem B 109:9104–9110
82. Yamada A, Ishikura K, Yamato T (2004) Proteins 55:1063–1069
83. Pelmenschikov V, Siegbahn PEM (2002) Inorg Chem 41:5659–5666
84. Cross JB, Vreven T, Meroueh SO, Mobashery S, Schlegel HB (2005) J Phys Chem B 109:4761–4769
85. Cerqueira NMFSA, Fernandes PA, Eriksson LA, Ramos MJ (2006) Biophys J 90:2109–2119
86. Sousa SF, Fernandes PA, Ramos MJ (2007) J Comput Chem 28:1160–1168
87. Yao L, Han Y, Cukier R (2006) J Phys Chem B 110:26320–26326
88. Kahn K, Bruice TC (2000) J Am Chem Soc 122:46–51
89. Kamachi T, Yoshizawa K (2005) J Am Chem Soc 127:10686–10692
90. Leopoldini M, Russo N, Toscano M, Dulak M, Wesolowski TA (2005) Chem Eur J 12:2532–2541
91. Wu X-H, Quan J-M, Wu Y-D (2007) J Phys Chem B 111:6236–6244
92. Matsubara T, Dupuis M, Aida M (2007) Chem Phys Lett 437:138–142

# CHAPTER 3

# COMPARISON OF REACTION BARRIERS IN ENERGY AND FREE ENERGY FOR ENZYME CATALYSIS

## G. ANDRÉS CISNEROS[1,2] AND WEITAO YANG[2]

[1]*Laboratory of Structural Biology, National Institute of Environmental Health Sciences, RTP, NC 27707, USA, e-mail: cisnero1@niehs.nih.gov*
[2]*Department of Chemistry, Duke University, Box 90346, Durham, NC 27708, USA, e-mail: weitao.yang@duke.edu*

**Abstract:**    Reaction paths on potential energy surfaces obtained from QM/MM calculations of enzymatic or solution reactions depend on the starting structure employed for the path calculations. The free energies associated with these paths should be more reliable for studying reaction mechanisms, because statistical averages are used. To investigate this, the role of enzyme environment fluctuations on reaction paths has been studied with an ab initio QM/MM method for the first step of the reaction catalyzed by 4-oxalocrotonate tautomerase (4OT). Four minimum energy paths (MEPs) are compared, which have been determined with two different methods. The first path (path A) has been determined with a procedure that combines the nudged elastic band (NEB) method and a second order parallel path optimizer recently developed in our group. The second path (path B) has also been determined by the combined procedure, however, the enzyme environment has been relaxed by molecular dynamics (MD) simulations. The third path (path C) has been determined with the coordinate driving (CD) method, using the enzyme environment from path B. We compare these three paths to a previously determined path (path D) determined with the CD method. In all four cases the QM/MM–FE method (Y. Zhang et al., JCP, 112, 3483) was employed to obtain the free energy barriers for all four paths. In the case of the combined procedure, the reaction path is approximated by a small number of images which are optimized to the MEP in parallel, which results in a reduced computational cost. However, this does not allow the FEP calculation on the MEP. In order to perform FEP calculations on these paths, we introduce a modification to the NEB method that enables the addition of as many extra images to the path as needed for the FEP calculations. The calculated potential energy barriers show differences in the activation barrier between the calculated paths of as much as 5.17 kcal/mol. However, the largest free energy barrier difference is 1.58 kcal/mol. These results show the importance of the inclusion of the environment fluctuation in the calculation of enzymatic activation barriers

**Keywords:**    QM/MM, 4-oxalocrotonate tautomerase, Free energy Perturbation, Enzyme catalysis

## 3.1.    INTRODUCTION

A key question in the action of enzymes is the understanding of the mechanisms by which they attain their catalytic rate enhancement relative to the uncatalyzed reactions. Some enzymes have been shown to produce rate accelerations as large as $10^{19}$ [1]. The theoretical determination of the reaction mechanisms by which enzymes carry out the chemical reactions has been an area of great interests and intense development in recent years [2–11]. A common approach for the modeling of enzyme systems is the QM/MM method proposed by Warshel and Levitt [12]. In this method the enzyme is divided into two parts. One part includes the atoms or molecules that participate in the chemical process, which are treated by quantum mechanical calculations. The other contains the rest of the enzyme and the solvent, generally thousands of atoms, which is treated by molecular mechanics methods.

In this way, the catalytic mechanism of enzymes can be modeled by calculating the reaction paths in the QM/MM potential energy surface (PES) and the free energy changes associated with the reactions. In the case of enzymes, extensive sampling of the PES is required to obtain reliable results, due to the large number of degrees of freedom in these systems. However, the sampling of multiple enzymatic reaction paths is computationally expensive because of the size of the systems. In this case, generally one, or at most a few minimum energy paths (MEPs) can be considered representative. In addition, the free energy changes associated with the reactions are not only better defined, but also characterize the reactions more accurately than potential energies when compared with experimental data.

One approach for calculating the free energies (FE) associated with reactions calculated with QM/MM methods is QM/MM-FE [13]. This method is an extension of the QM-FE method initially developed by Jorgensen [14] and by Kollman and coworkers [15]. The QM/MM-FE approach is carried out in two stages: QM/MM optimization of the reaction path, followed by the calculation of the free energies for the path using free energy perturbation (FEP) methods. The advantage of using the QM/MM calculated reaction path is that it has been determined in the enzyme environment, which includes the polarization effects of the enzyme environment on the QM part. Recently two other methods have been proposed for the determination of FE profiles in enzymes: QM/MM minimum free energy path (QM/MM-MFEP) [11, 16, 17] and QM/MM-MD [18]. In the former, the thermodynamics of a complex reaction system are described by a potential of mean force (PMF) of the QM subsystem. The latter employs on-the-fly Born–Oppenheimer MD simulations with the ab initio QM/MM approach using the umbrella sampling method [18].

As explained above, the QM/MM-FE method requires the calculation of the MEP. The MEP for a potential energy surface is the steepest descent path that connects a first order saddle point (transition state) with two minima (reactant and product). Several methods have been recently adapted by our lab to calculate MEPs in enzymes. These methods include: coordinate driving (CD) [13, 19], nudged elastic band (NEB) [20–25], a second order parallel path optimizer method [25, 26], a procedure that combines these last two methods in order to improve computational efficiency [27],

the superlinearly convergent quadratic string method (QSM) [28], and a sequential quadratic programming (SQP) method [29].

The NEB, QSM, SQP and parallel path optimizer methods are types of "chain of states" methods [20, 21, 25, 26, 30, 31]. In these methods the MEPs are represented by a discrete number of structures that form a chain connecting the reactant to the product. In general, the intermediate structures are obtained by a linear interpolation between the end points, and subsequently all points are optimized to the MEP simultaneously. Since the QM/MM calculations for enzyme systems are computationally expensive, only a small number of images are employed to represent the path. However, the use of a small number of images along the path for the calculation of the MEP prevents the determination of the free energy with the FEP method on these paths. The reason for this is that for FEP calculations, the averaging of the systems has to be performed with care. If the systems on which the sampling is being performed are too different, the convergence may be slow [14]. This is due to the fact that if the perturbation between two images is too big, the sampling will not be realistic.

Our goal in the present work is to assess the role of the MM fluctuations to the catalytic mechanism of 4OT [32–36]. We compare the potential and free energy barriers for the first step of the reaction to determine whether potential energies alone are enough to characterize the enzymatic activation barriers. Additionally, we present the development of a method that allows us to perform FEP calculations on paths with a small number of images describing the MEP. Here we propose adding as many extra images as needed between the converged points on paths obtained with the parallel iterative path method or the combined procedure, to obtain a proper sampling. After the extra points have been added, a modified NEB calculation is performed in which these images are optimized to the MEP. During this optimization, the original optimized images are held fixed. When the new images have been optimized to the MEP, free energy sampling may be carried out on these reaction paths.

The organization of the paper is as follows. In the Methods Section we present a brief review of the procedures employed to determine the MEP, namely QM/MM geometry optimizations and the CD method (Section 3.2.1). NEB and the second order parallel path optimizer methods are reviewed in Section 3.2.2. In Section 3.2.3 we present a modified NEB method which allows the addition of extra images to a converged path. Finally we describe the FEP method developed in our lab in Section 3.2.4. In the Computational Details Section we present a description of the procedures employed for the determination of all four paths. Subsequently, we analyze the results obtained from all four paths and conclude with closing remarks.

## 3.2. METHODS

### 3.2.1. QM/MM Optimization and the CD Method

In the QM/MM potential energy model, the total energy of the system is

$$E_{Total} = E_{MM} + E_{QM} + E_{QM/MM}. \tag{3-1}$$

The QM/MM interactions ($E_{QM/MM}$) are taken to include bonded and non-bonded interactions. For the non-bonded interactions, the subsystems interact with each other through Lennard–Jones and point charge interaction potentials. When the electronic structure is determined for the QM subsystem, the charges in the MM subsystem are included as a collection of fixed point charges in an effective Hamiltonian, which describes the QM subsystem. That is, in the calculation of the QM subsystem we determine the contributions from the QM subsystem ($E_{QM}$) and the electrostatic contributions from the interaction between the QM and MM subsystems as explained by Zhang et al. [13].

Geometry optimizations are carried out by an iterative minimization procedure as described by Zhang et al. [13] In this procedure one iteration consists of a complete optimization of the QM subsystem, followed by a complete optimization of the MM subsystem. At each point the subsystem not being optimized is held fixed at the geometry obtained from the previous iteration; QM/MM interactions are also included at each iteration. The iterations are continued until the geometries of both systems no longer change.

When the MM subsystem is being optimized, or a molecular dynamics simulation is being carried out on the MM subsystem, the QM/MM electrostatic interactions are approximated with fixed point charges on the QM atoms which are fitted to reproduce the electrostatic potential (ESP) of the QM subsystem [37].

One of the procedures employed for the determination of the MEP is the CD method [19]. This method introduces a harmonic restraint on the reaction coordinate, which is a linear combination of the distances between the atoms involved in the reaction to perform an optimization along a proposed reaction path. In this case the reaction coordinate is given by the expression:

$$R = \sum_{i=1}^{n} a_i r_i, \tag{3-2}$$

where $r_i$ are the distances between atoms, $a_i$ is constant 1 for the distance that increases, $-1$ for the distance that decreases. The sum over $i$ includes all the distances that change throughout the course of the reaction. $R$ is included in the following energy expression:

$$E_{Restrain} = k(R - s)^2, \tag{3-3}$$

where $R$ is given by Eq. (3-2), $s$ is an adjustable parameter corresponding to the value of the reaction coordinate which is varied in a stepwise manner at each point on the PES, and $k$ is a force constant. In this case the value of $k$ was set to 2000 kcal/mol for all points. This energy is included in the total energy expression in the process of the optimization.

All the reaction paths calculated with the CD method were determined by stepping forward (from reactant to intermediate state) and backward (from intermediate

to reactant state) along the path several times until there was no change between the forward and backward paths.

### 3.2.2. Chain of States Methods

The procedure and methods for the MEP determination by the NEB and parallel path optimizer methods have been explained in detail elsewhere [25, 27]. Briefly, these methods are types of "chain of states" methods [20, 21, 25, 26, 30, 31]. In these methods the path is represented by a discrete number of images which are optimized to the MEP simultaneously. This parallel optimization is possible since any point on the MEP is a minimum in all directions except for the reaction coordinate, and thus the energy gradient for any point is parallel to the local tangent of the reaction path.

The path optimizations are carried out by an iterative optimization procedure [25]. In the case of enzyme systems, because of the large number of degrees of freedom, we partition them into a core set and an environmental set. The core set is small and contains all the degrees of freedom that are involved with the chemical steps of the reaction, while all the remaining degrees of freedom are included in the environmental set. In all the QM/MM calculations presented below, the core set is defined by the QM subsystem and the environmental set by the MM subsystem.

Initially, the path is minimized starting with the optimization of the core set of atoms (path optimization) with either the NEB or the parallel path optimizer method, followed by the optimization of the environment set of atoms. During the optimization of the environment, the path may become discontinuous in the early cycles because the initial guess to the MEP may not be accurate. This problem is overcome by performing a restrained minimization of the environment set as proposed by Xie et al. [25]. In this method, the environmental atoms are restrained during the MM optimization. These restraints are gradually reduced after each cycle until no restraints remain, ensuring a smooth change in the MM environment and avoiding sudden fluctuations.

To ensure a smooth and continuous change of the environment set during the optimization, the initial coordinates for the environment atoms are set to be the same for all the replicas along the path. Thus, the initially guessed replicas differ only in the coordinates of the core atoms. In other words, the initial environment for all the points on the path is the same and is taken to be the environment from one of the end points.

In NEB the path is simulated by an elastic band with $N + 1$ images where the end points $\mathbf{x}_0$ and $\mathbf{x}_N$ are fixed and correspond to the reactant and product of the system under study, while the $N - 1$ remaining images are optimized in Cartesian coordinates using a projected velocity Verlet algorithm by minimizing the forces acting on the images along the path [20]. A spring interaction is included between the images to ensure equal spacing along the path. A unique feature of NEB is that the forces are projected in such a way that the spring forces do not interfere with the minimization, while also ensuring that the true forces do not interfere with the image separation.

The parallel path optimizer method is a second order method based on the approach developed by Ayala and Schlegel for small molecules [26], which combines synchronous transit and quasi–Newton methods. This method has been extended by our group to study reactions in enzymes [25]. In this case, the path is also represented by a set of discrete points connecting the reactant and product. For each point the total energies and gradients are calculated and empirical Hessians are obtained. The highest energy point is chosen to optimize to the closest transition state (TS) which divides the path into two downhill segments. If the end points are not already at a minimum, they can be optimized to the closest minimum. The remaining points on the path are optimized to lie on the steepest descent path.

We have recently developed a procedure to calculate MEPs [27] which combines NEB and the parallel path optimizer. This procedure increases the performance of MEP determinations by taking advantage of the strengths of both methods. The combined procedure consists of two parts, an initial NEB optimization, and a refinement of the path with the parallel path optimizer method. Initially, the path is iteratively optimized with the NEB method for the core set followed by an optimization of the environment set. When the core set is optimized with NEB the convergence is chosen to be relatively loose to reduce the computing time. During this step, the optimization for the environment set is carried out in stages by using a restrained minimization as mentioned above [25].

Once the path has been optimized with NEB, it is used as the initial guess for the parallel path optimizer method. In this second step the path is again iteratively optimized with the parallel path optimizer method for the core set, followed by the optimization of the environment set. In this part of the calculation no restraints are imposed on the environment set during the optimization. The iterations are continued until all the convergence criteria are met and the final optimized MEP is obtained.

### 3.2.3.    Modified NEB Method

In order to perform FEP calculations on optimized paths with a small number of images, extra images need to be added on the path between the previously optimized points. Once these extra images have been added, an optimization has to be performed to minimize them to the MEP. Here we have developed a modification to our NEB QM/MM implementation [27]. This modification allows for the optimization of only selected images on the path while maintaining the points previously optimized with the parallel iterative path method or the combined procedure fixed.

Initially, the coordinates $\mathbf{x}'_i$ for the extra images added to the path are approximated by a linear interpolation between the converged points for the core set. In the case of the environment set, the initial coordinates are approximated by the environment set of the immediate neighboring converged point. That is, if only one image is added between each pair of optimized points, the initial coordinates of the core set for the image added between the optimized points $\mathbf{x}_0$ and $\mathbf{x}_1$ are given by a linear interpolation between $\mathbf{x}_0$ and $\mathbf{x}_1$. The environment coordinates are set to correspond

to that of $x_1$, and so on for all the added images. If two images are added between $x_0$ and $x_1$, the core set coordinates are also approximated by a linear interpolation between $x_0$ and $x_1$. Here the environment coordinates are set to that of $x_1$ for both new images, for all the added images. By approximating the coordinates of the extra images in this way, we ensure that the added points are not too far from the MEP. This results in faster convergence of the added images to the MEP.

Once the images have been added, they are iteratively optimized by performing a modified NEB calculation. Here, only the added images are optimized and the previously converged points from the parallel path optimizer method or combined procedure calculations are held fixed. Note that in order to calculate the tangents needed by NEB for the images along the path, energies and gradients need to be calculated for each image. Thus, in the case of the converged points, only one energy-gradient calculation must be performed in the first path optimization cycle.

In order to ensure that the added images remain equidistant along the reaction coordinate during the course of the path optimization, a harmonic restraint is included when the core set is being optimized. The reason is that the definition of the spring forces of NEB is calculated with all the degrees of freedom of the core set. However, only a small number of those degrees of freedom are involved in the reaction coordinate. If this restraint is not included during the optimization, the extra images "slide-down" towards the minima, away from the high energy regions of the path. This is because the remaining degrees of freedom in the core set, which are not involved in the reaction coordinate, could satisfy the equal-spacing constraint for the NEB optimization. Thus, the added points could optimize to lower energies along the reaction coordinate direction.

By adding a restraint which is projected exclusively along the reaction coordinate, we can ensure that all the newly added images will be evenly spaced. This restraint is defined as a linear combination of the distances between the atoms involved in the reaction, and is given by

$$R = \sum_{i=1}^{n} a_i r_i, \tag{3-4}$$

where $r_i$ are the distances between atoms, $a_i$ is a constant, set to 1 for the distances that increase and $-1$ for the distances that decrease. The sum over $i$ includes all the distances that change throughout the course of the reaction. After $R$ has been determined, it is employed to determine the restrain energy by

$$E_{Restrain} = k(R - s)^2, \tag{3-5}$$

where $R$ is given by Eq. (3-4), $s$ is an adjustable parameter corresponding to the value of the reaction coordinate at each point on the PES, and $k$ is a force constant. This energy is included in the total energy expression in the process of the optimization. In this case, $s$ is selected such that each added image is equidistant along

the reaction coordinate from the previously optimized pair of images to which it is added. For example, if only one image is added between each pair of optimized points, the parameter $s$ for the image added between the optimized points $\mathbf{x}_0$ and $\mathbf{x}_1$ is set as $s = R_0 + (\frac{abs(R_0 - R_1)}{2})$, and so on for the rest of the points. In the test calculations presented below the value of $k$ was set to 2000 kcal/mol for all points.

### 3.2.4.     The QM/MM–FE Approach

Once the MEPs have been determined with enough images along the path, free energy perturbation (FEP) calculations can be carried out to determine the changes in free energy associated with the paths. The details of the QM/MM-FE method have been described in detail previously in Ref. [13]. In this section we briefly review the method for obtaining the free energy profile of enzyme reactions.

Relative free energy changes along the reaction coordinate may be calculated via MD simulations. Subsequently, they may be combined with the QM total energies of the reacting QM subsystem to obtain the overall free energy profiles associated with the reaction steps in the following way

$$\Delta F(R_c) \approx \Delta E_{QM}(\mathbf{x}_{QM}^{min}) + \Delta F_{QM/MM}(\mathbf{x}_{QM}^{min}) \tag{3-6}$$

where $\Delta E_{QM}(\mathbf{x}_{QM}^{min})$ is the change in QM energies in the reacting QM subsystem, $\Delta F_{QM/MM}(\mathbf{x}_{QM}^{min})$ is the QM/MM interaction free energy changes, and $\mathbf{x}_{QM}^{min}$ are the degrees of freedom of the QM subsystem. In this case it is assumed that the QM subsystem only fluctuates around that path.

It is important to point out that Eq. (3-6) has the same form as Jorgensen's approach [14, 15]. However, this approach has two major differences. First, $\mathbf{x}_{QM}^{min}$ is defined as the MEP obtained from the modified NEB procedure (Section 3.2.3) or from the CD method (Section 3.2.1), where the QM subsystem interacts with the enzyme environment. In the case of the QM-FE approach, the QM subsystem is isolated.

Second, $\Delta E_{QM}(\mathbf{x}_{QM}^{min})$ is calculated as the difference between the QM subsystem energy computed form the ab initio calculation ($E_{QM}(QM)$), and the QM/MM electrostatic energy computed classically ($E_{electrostatics}(QM/MM)$) [13].

The free energy differences $\Delta F_{QM/MM}(r_{QM}^{min})$ between different images along the path are calculated using MD simulations and the FEP theory [38],

$$\Delta F_{QM/MM}^{A \to B} = F_{QM/MM}(R_C^B) - F_{QM/MM}(R_C^A)$$

$$= -\frac{1}{\beta} ln \left\langle exp \left\{ -\beta [E_{QM/MM}(\mathbf{x}_{QM}^{min}(R_C^B)) - E_{QM/MM}(\mathbf{x}_{QM}^{min}(R_C^A))] \right\} \right\rangle_{MM,A} \tag{3-7}$$

where $R_C^A$ and $R_C^B$ represent images along the MEP, and $\langle\cdots\rangle_{MM,A}$ represents an ensemble average over the MM sub-system, with the QM subsystem frozen to the $\mathbf{x}_{QM}^{min}(R_C^A)$ configuration, as defined in Ref. [13].

The calculation of the free energies of enzymes by this procedure provides several advantages. In this approach, the MEP is determined in the enzyme environment with a smooth connection between the QM and MM subsystems by means of the pseudobond QM/MM method [39]. Also, the polarization effects of the enzyme environment on the QM subsystem have been included [13].

It is important to note that in this method, the dynamic fluctuations associated with the QM subsystem are assumed to be independent of the fluctuations from the MM subsystem. Also, in this method we assume that the contributions of the fluctuations of the QM subsystem to the total free energy are the same along the reaction coordinate. We have recently addressed these approximations by developing a novel reaction path potential method where the dynamics of the system are sampled by employing an analytical expression of the combined QM/MM PES along the MEP [40].

### 3.2.5. Computational Details

In all cases the calculations were performed using QM/MM methodology that includes the pseudobond model for the QM/MM boundary [13, 39, 41]. This methodology has been implemented in a modified version of Gaussian 98 [42], which interfaces to a modified version of TINKER [43]. The AMBER94 all-atom force field parameter set [44] and the TIP3P [45] model for water were used.

In all cases the core (QM) subsystem consists of 36 atoms which include the substrate molecule (2o4hex), a water molecule and Pro-1. The remainder of the enzyme and solvation sphere is located in the environment (MM) subsystem, to give a total of 13197 atoms for the overall system. For all calculations involving the environment set (MM optimizations and MD simulations), since we do not simulate an infinite system with periodic boundary conditions, the conformational fluctuations near the boundary will not be realistic and thus were ignored. A 20 Å sphere centered around C3 of the substrate was employed, in which atoms were allowed to move; all atoms outside this sphere were held fixed. The twin range cutoff method [46] for non bonded interactions was employed, with a long-range cutoff distance of 15 Å and a short-range cutoff of 8 Å, as in our previous calculations [33].

As explained above, we compare four different reaction paths for the first step of the reaction catalyzed by 4OT. The first difference between the paths was the method employed to determine the MEP for the paths. The first two paths (A and B), were determined by the combined procedure [27, 35]. In both cases seven images were employed for the optimization of the path. Path C was calculated with the CD method. Finally, path D corresponds to the MEP obtained from our previous calculations [33].

The second major difference between the paths was the initial enzyme environment employed for the path determinations. The environment for path A corresponds to a previously optimized MEP [25], where the initial enzyme environment was

chosen to be that of the intermediate point obtained from our previous path determination [33]. In this case the environment was relaxed by running 200 ps of MD. The environment for path B was obtained by performing a 200 ps MD simulation on the intermediate and taking snapshots at 10 ps intervals as described in Ref. [35]. Each snapshot was optimized without constraints and the lowest energy structure was chosen to be the intermediate point, as well as the initial enzyme environment. In the case of path C, the optimized structure for the reactant from path B was employed to start the CD calculation.

All energy-gradient calculations required for the NEB calculations to optimize the core set of degrees of freedom of the added points on the paths were performed at the HF/3–21G level, which is the same as the optimization of the paths with the combined procedure. The core sets were considered converged when the gradient was below of 0.01 au for all images. A maximum of 10 path optimization cycles was used for each NEB iteration. In the case of the optimization for the environment set performed after each NEB calculation, all points being optimized were required to converge to 0.1 kcal/mol $Å^2$. This procedure was repeated until no change was observed in the core and environment sets. In the case of the CD method, all points were also optimized at the HF/3–21G level. The convergence criteria was the default of Gaussian 98 [42].

In order to perform FEP calculations on paths A and B, the modified NEB method presented in Section 3.2.3 was employed. In both cases two extra images were added between the optimized images. This resulted in a total of 19 images for both paths.

FEP calculations for paths A, B and C were performed with a 40 ps equilibration run prior to the sampling for all points along the path. The free energy contributions were sampled for 20 ps for each point on the MEP. In all cases a time step of 2.0 fs was employed, maintaining a constant temperature of 300 K. The SHAKE [47] algorithm was used to constrain all bonds involving hydrogen atoms.

## 3.3.    RESULTS AND DISCUSSION

In this section we present the results from our QM/MM and FEP calculations. First, we present the path obtained by using the modified NEB method from Section 3.2.3 (path A) to test it's applicability. Subsequently, we proceed to the description and analysis of the energetics and structures obtained from the comparison of the four different paths.

In order to compare the calculated potential and free energy differences on enzymatic reaction mechanisms we have chosen to study a system that has been previously calculated by our group, 4-oxalocrotonate-tautomerase (4OT). 4OT is a hexameric bacterial enzyme that catalyzes the isomerization of unconjugated $\alpha$-keto acids such as 2-oxo-4-hexenedioate (2o4hex) to its conjugated isomer 2-oxo-3-hexenedioate (2o3hex).

In a previous study [33] we confirmed the reaction mechanism proposed by Harris et al. [32] for 4OT, which is known to be a general acid–base reaction that takes place in two steps (see Figure 3-1). In the first step, one of the hydrogens from C3

*Figure 3-1.* Proposed 4OT mechanism [32, 33]

of the substrate (2o4hex) is abstracted by the nitrogen atom of Pro-1. This produces a negative charge on the carbonyl oxygen of the substrate in the transition state and intermediate structures, which is stabilized by Arg-39" and an ordered water in the active site. For the second step the proton is returned from Pro-1 to C5 of the substrate to form the product (2o3hex).

Initially, we have applied the modified NEB method to the calculation of both steps of the 4OT catalyzed reaction. The free energy profiles and relative free energies obtained with this method were compared to our previously determined profiles [33]. As we had previously shown, the calculated MEPs for Ref. [33] determined with the reaction coordinate driving method, and the MEPs for Ref. [25] calculated with the parallel path optimizer method, agree in the overall shape and relative potential energies. This provides a good starting point for our comparison.

Note that, for the comparisons of all four paths discussed below, the first step of the reaction obtained from the modified NEB calculation corresponds to path A, and the first step of our previously determined paths [33] corresponds to path D (see Figures 3-2 and 3-3).

The converged MEPs from the combined procedure were employed as a template for our modified NEB calculations [27]. The calculated MEPs for both steps of the reaction were determined with seven images for each path. In order to carry out the FEP calculation, following the procedure of Section 3.2.3, two images were added between each pair of converged images on the path. This produced overall paths with



*Figure 3-2.* Calculated forward and backward paths for path D

*Figure 3-3.* Free energy profiles for path D

19 images for each step of the reaction. The new images were optimized to the MEP with the iterative modified NEB method maintaining the original optimized images fixed.

For both reaction path calculations, it only took two modified NEB cycles to optimize the new images to the MEP. The reason for this very fast convergence is that the initial approximation for the extra images on the path is not too far from the MEP.

Figures 3-4 and 3-5 show the optimized paths with the added images and the original combined method [27] and parallel path optimizer method [25] calculated paths for the first and second steps of the reaction respectively. In both cases, the addition of extra images on the converged path, and subsequent optimization of these extra images produces a smoother path since the additional images allows for a better mapping of the potential energy surfaces (PESs).

After obtaining the converged paths for both steps of the reaction, we proceeded to perform the FEP calculations. Figure 3-6 shows the calculated potential and free energy surfaces for the first and second steps of the reaction. In both cases, the forward and backward calculated free energies ($\Delta F_{QM/MM}$) differ by less than 0.5 kcal/mol at each image.

Table 3-1 shows the comparison between our previous calculations and the present results for the potential and free energy differences for selected structures. Overall the results are in agreement. However, it is important to note that the present results are lower than the previous ones in all cases.

*Figure 3-4.* Modified NEB path (path B) for the first step of the reaction catalyzed by 4OT. (1) from Ref. [25]. (2) from Ref. [27]



*Figure 3-5.* Modified NEB path (path B) for the second step of the reaction catalyzed by 4OT. (1) from Ref. [25]. (2) from Ref. [27]

*Figure 3-6.* Free energy profiles for both steps of the reaction catalyzed by 4OT (path B)

*Table 3-1.* Calculated potential and free energy differences for path B (in kcal/mol) between the determined structure and the reactant (ES complex), where $\Delta E$ is the total HF potential energy difference, $\Delta E_{QM}$ refers to the QM energy difference between two QM subsystems. $\Delta F_{QM/MM}$ is the free energy change in the QM/MM interaction, and $\Delta F = \Delta E_{QM} + \Delta F_{QM/MM}$. Numbers without parentheses correspond to the present work and numbers in parentheses correspond to our previous determinations (path D) [33]

| Structure | $\Delta E$ | $\Delta E_{QM}$ | $\Delta F_{QM/MM}$ | $\Delta F$ |
|---|---|---|---|---|
| TS1 | 17.77 (22.02) | 35.25 (35.47) | −19.96 (−20.29) | 15.48 (15.18) |
| I | 10.11 (17.27) | 38.62 (44.47) | −34.90 (−32.87) | 3.72 (11.60) |
| TS2 | 20.98 (22.50) | 34.97 (36.08) | −19.01 (−18.75) | 15.97 (17.33) |
| P | 1.45 (0.63) | 0.35 (0.66) | −0.44 (−0.04) | −0.09 (0.62) |

Moreover, these energy barriers for the TS show a surprising result. If the potential energies are compared, a reduction of around 4 kcal/mol is observed. However, the calculated free energy differences show a maximum difference of only around 1.5 kcal/mol.

In order to better understand these results we have compared the results from the first step of these paths (paths A and D respectively) to a previously obtained path (path B) which was calculated with the combined procedure [35]. In the case of path B, the enzyme environment has been relaxed by MD simulations and QM/MM optimizations as explained in Section 3.2.5. Additionally, we have also calculated a

*Table 3-2.* Potential and free energy differences (in kcal/mol) for the TS and intermediate points of all four paths. All energies are relative to the respective reactant structure

|            | TS1          |            | Intermediate |            |
|------------|--------------|------------|--------------|------------|
| Structure  | $\Delta$ E   | $\Delta$ F | $\Delta$ E   | $\Delta$ F |
| path A     | 17.77        | 15.48      | 10.11        | 3.72       |
| path B     | 16.85        | 14.45      | 9.34         | 3.71       |
| path C     | 20.37        | 13.90      | 15.32        | 4.62       |
| path D     | 22.02        | 15.18      | 17.27        | 13.52      |

fourth path (path C) with the CD method. In this case the starting structure was the optimized reactant from path B.

The calculated potential and free energies for the first step of the four paths is presented in Table 3-2. As can be seen, the paths calculated with the combined procedure present activation energies of 17.77 and 16.85 kcal/mol for paths A and B respectively. On the other hand, the calculated potential activation energies are 20.37 kcal/mol for path C and 22.02 for path D. This is in contrast to the calculated

*Table 3-3.* Differences in individual residue contribution to TS stabilization. All values are relative to the TS from path B (lowest energy structure). Negative values mean that that particular residue helps stabilize the TS while a positive value means that the residue destabilizes the TS. Dashes mean that the absolute difference in contribution of that residue was less than 0.5 kcal/mol. Total value includes *all* residues in the enzyme

|                   | $\Delta$ E (kcal/mol) |         |         |
|-------------------|--------|--------|--------|
| Residues          | Path A | Path C | Path D |
| Arg–39″           | –      | 0.97   | 2.00   |
| Arg–61′           | 1.80   | 1.36   | 2.63   |
| Arg–62′           | –      | –      | 1.17   |
| Ser–24            | –      | –      | 1.02   |
| Ile–27            | –      | –      | −1.29  |
| Pro–34            | –      | –      | 1.08   |
| Leu–35            | –      | –      | 1.28   |
| Thr–36            | –      | –      | 1.52   |
| Ser–37            | 0.60   | –      | −0.76  |
| Val–38            | –      | –      | 1.31   |
| Arg–39            | –      | –      | –      |
| Ile–7′            | –      | –      | 1.10   |
| Leu–8′            | –      | –      | −0.94  |
| Arg–11′           | –      | −1.02  | –      |
| Phe–50′           | 0.60   | –      | –      |
| $H_2O$–5          | –      | –      | –      |
| Total             | 2.6    | 3.16   | 8.37   |

activation free energies which differ by no more than 1.6 kcal/mol in all cases. The difference in the potential energies for all four paths is explained by the fact that the MEPs for paths A and B were calculated with a method (combined procedure) that provides a much more accurate representation of the paths. Note that the paths calculated with the CD method present a higher potential energy barrier of activation. However, the inclusion of the fluctuation of the MM environment corrected the shortcomings of the paths determined with this method.

In the case of the intermediates, the calculated potential and free energies for the first three paths follow the same trend as for the TS. However, a big difference is observed with respect to path D. This energy difference is due to the fact that in the case of the first three paths, the MM environment of the intermediate structure was more relaxed for the calculations than that of path D. This relaxation in the environment resulted in the big energy differences for both the potential and free energy barriers of path D with respect to the other paths.

In order to determine where the energy differences stem from, the effects of the individual residues on the stabilization(destabilization) of the QM subsystem were analyzed as explained in Ref. [33]. These results are shown in Tables 3-3 and 3-4,

*Table 3-4.* Differences in individual residue contribution to intermediate stabilization. All values are relative to the intermediate from path B (lowest energy structure). Negative values mean that that particular residue helps stabilize the intermediate while a positive value means that the residue destabilizes the intermediate. Dashes mean that the absolute difference in contribution of that residue was less than 0.5 kcal/mol. Total value includes *all* residues in the enzyme

| Residues | $\Delta E$ (kcal/mol) | | |
|---|---|---|---|
| | Path A | Path C | Path D |
| Arg–39'' | −0.84 | – | 2.45 |
| Arg–61' | 0.87 | 2.29 | 5.41 |
| Arg–62' | – | – | 2.47 |
| Ser–24 | – | – | 1.37 |
| Ile–27 | – | – | −1.57 |
| Pro–34 | – | – | 1.66 |
| Leu–35 | – | – | 3.38 |
| Thr–36 | – | – | 2.09 |
| Ser–37 | – | 1.05 | −3.72 |
| Val–38 | −0.85 | −0.73 | 2.28 |
| Arg–39 | – | – | −0.81 |
| Ile–7' | – | – | −1.66 |
| Leu–8' | – | – | 2.09 |
| Arg–11' | – | 1.45 | – |
| Phe–50' | – | – | 1.01 |
| $H_2O$–5 | −0.99 | – | – |
| Total | −2.4 | 0.56 | 14.87 |

in both cases the residue contribution differences are reported relative to path B. The residue analysis shows that in the case of paths A and C only a small number of residues show a significant change in contribution with respect to path B. In the case of path D however, a big difference is observed in a large number of residues, especially on part of the chain with the catalytic Pro residue.

As can be seen from Figures 3-7 and 3-8, the structural changes in a sphere of 12 Å around the active site remains relatively similar for all four calculated structures for the TSs as well as the intermediates. The segment of residues (24–39) that show a big deviation between the intermediates of path B and path D are shown in Figure 3-9. As can be seen, the structural differences are more notable among the structures which also helps explain the big energy difference between path D and the rest of the calculated paths.



*Figure 3-7.* Superposition of TS structures (active site) for all 4 paths. Path A: *red*, path B: *blue*, path C: *grey*, path D: *yellow*

*Figure 3-8.* Superposition of intermediate structures (active site) for all 4 paths. Path A: *red*, path B: *blue*, path C: *grey*, path D: *yellow*

## 3.4.     CONCLUSIONS

We have compared the potential and free energies from four different calculated paths for the first step of the isomerization of 2o4hex catalyzed by 4OT. Two of these paths were determined with a combined "chain-of-states" method. The remaining two paths were obtained with the coordinate driving method.

Our results show that the calculated potential energies for the TS obtained from the combined procedure are around 4 kcal/mol lower than the corresponding ones calculated with the CD method. In contrast, the calculated free energies of activation for all four paths differ by not more than 1.6 kcal/mol. These results show that the inclusion of the fluctuation of the MM environment dramatically improves the results of the calculations of enzymatic catalysis, even if the calculated PES is not highly accurate. In addition, the calculation of free energies for multiple paths using the QM/MM-FE method can serve as an alternative to more expensive sampling methods such as QM/MM-MFEP and QM/MM-MD.

We have also presented the development of a method that enables the determination of free energy profiles associated with calculated MEPs for enzymatic reaction

*Figure 3-9.* Superposition of residues 24–39 for the intermediate structures of all 4 paths. Path A: *red*, path B: *blue*, path C: *grey*, path D: *yellow*

mechanisms with a small number of points along the path. In this method extra images are added to the converged MEP between the minimized images and optimized with a modified NEB implementation that maintains the original images fixed. The addition of enough images on the path allows for the proper sampling of the associated free energy profile.

As explained above, two other methods have been recently introduced for the determination of free energy barriers in QM/MM: QM/MM-MFEP [11, 16, 17] and QM/MM-MD [18]. While both are more robust and advanced approaches to determining free energy barriers, they are more expensive computationally. In QM/MM-FEP this is due to the associated sampling and optimization required to obtain the PMF for the QM subsystem. In the case of the QM/MM-MD the expense results in the evaluation of the energy and gradient of the entire QM/MM system for each MD step. Therefore, the calculations of free energies on multiple paths with the QM/MM-FE method can provide meaningful insights into the mechanisms with reduced computational expense.

## ACKNOWLEDGMENTS

## REFERENCES

1. Wolfenden R, Snyder M (2001) Acc Chem Res 34: 938
2. Villa J, Warshel A (2001) J Phys Chem B 105: 7887
3. Kollman P et al. (2001) Acc Chem Res 34: 72
4. Warshel A (2002) Acc Chem Res 35: 385
5. Gao J, Truhlar D (2002) Ann Rev Phys Chem 53: 467
6. Monard G, Prat-Resina X, González-Lafont A, Lluch J (2003) Int J Quantum Chem 93: 229
7. Warshel A (2003) Annu Rev Biophys Biomol Struct 32: 425
8. Hermann JC, Ridder L, Höltje H-D, Mulholland AJ (2006) Org Biomol Chem 4: 206
9. Zhang Y (2006) Theo Chem Acc 116: 43
10. Senn H, Thiel W (2007) "QM/MM methods for biological systems" in Atomistic approachs in modern biology, Vol. 268 of Topics in Current Chemistry, Springer Berlin/Heidelberg, Berlin, Germany, pp 173–290
11. Hu H, Yang W (2008) Ann Rev Phys Chem 59: 573
12. Warshel A, Levitt M (1977) J Mol Biol 103: 227
13. Zhang Y, Liu H, Yang W (2000) J Chem Phys 112: 3483
14. Jorgensen W (1989) Acc Chem Res 22: 184
15. Stanton R, Perakyla M, Bakowies D, Kollman P (1998) J Am Chem Soc **120**: 3448
16. Hu H, Lu Z, Yang W (2007) J Chem Theo Comp 3: 309
17. Hu H et al. (2008) J Chem Phys 128: 034105
18. Wang S, Hu P, Zhang Y (2007) J Phys Chem B 111: 3758
19. Williams I, Maggiora G (1982) J Mol Struct 89: 365
20. Jónsson H, Mills G, Jacobsen K (1998) In: Berne BJ, Ciccotti G, Coker DF (eds) "Nudged Elastic Band Method", in Classical and quantum dynamics in condensed phase simulations, World Scientific, Singapore, pp 387–404
21. Henkelman G, Jónsson H (1999) J Chem Phys 111: 7010
22. Henkelman G, Jónsson H (2000) J Chem Phys 113: 9978
23. Maragakis P et al. (2002) J Chem Phys 117: 4651
24. Jhih-Wei C, Trout B, Brooks B (2003) J Chem Phys 119: 12708
25. Liu H, Lu Z, Cisneros GA, Yang W (2004) J Chem Phys 121: 697
26. Ayala P, Schlegel H (1997) J Chem Phys 107: 375
27. Cisneros GA, Liu H, Lu Z, Yang W (2005) J Chem Phys 122: 114502
28. Burger SK, Yang W (2006) J Chem Phys 124: 054109
29. Burger SK, Yang W (2006) J Chem Phys 124: 224108
30. Woodcock HL et al. (2003) Theo Chem Acc 109: 140
31. Elber R, Karplus M (1987) Chem Phys Lett 139: 375
32. Harris T et al. (1999) Biochemistry 38: 12343
33. Cisneros GA, Liu H, Zhang Y, Yang W (2003) J Am Chem Soc 134: 10348
34. Cisneros GA et al. (2004) Biochemistry 43: 6885

35. Cisneros GA et al. (2006) J Phys Chem A 110: 700
36. Tuttle T, Keinan E, Thiel W (2006) J Phys Chem B 110: 19685
37. Besler BH, Merz KM Jr, Kollman P (1990) J Comp Chem 11: 431
38. Zwanzig RW (1954) J Chem Phys 22: 1420
39. Zhang Y, Lee T, Yang W (1999) J Chem Phys 110: 46
40. Lu Z, Yang W (2004) J Chem Phys 121: 89
41. Zhang Y, Liu H, Yang W (2002) In: Schlick T, Gan HH (eds) "Ab Initio QM/MM and Free Energy Calculations of Enzyme Reactions", in Computational Methods for Macromolecules–Challenges and Applications, Springer Verlag, Heidelberg, Germany, pp 332–354
42. Frisch MJ et al. (1998) Gaussian 98, Revision A.8, Gaussian, Inc., Pittsburgh PA
43. Ponder J (1998) TINKER, Software Tools for Molecular Design, Version 3.6: the most updated version for the TINKER program can be obtained from J.W. Ponder's WWW site at http://dasher.wustl.edu/tinker, Washington University, St. Louis
44. Cornell WD et al. (1995) J Am Chem Soc 117: 5179
45. Jorgensen W et al. (1983) J Chem Phys 79: 926
46. van Gunsteren W et al. (1984) J Comp Chem 5: 272
47. Ryckaert J, Ciccotti G, Berendsen H (1977) J Comp Phys 23: 327

CHAPTER 4

# QUANTUM MECHANICAL METHODS FOR BIOMOLECULAR SIMULATIONS

KIN-YIU WONG, LINGCHUN SONG, WANGSHEN XIE, DAN T. MAJOR, YEN-LIN LIN, ALESSANDRO CEMBRAN, AND JIALI GAO

*Department of Chemistry, Digital Technology Center, and Minnesota Supercomputing Institute, University of Minnesota, MN 55455, USA, e-mail: gao@jialigao.org*

**Abstract:**     We discuss quantum mechanical methods for the description of the potential energy surface and for the treatment of nuclear quantum effects in chemical and biological applications. Two novel electronic structure methods are described, including an electronic structure-based explicit polarization (X-Pol) force field and an effective Hamiltonian molecular orbital and valence bond (EH-MOVB) theory. In addition, we present two path integral techniques to treat nuclear quantum effects, which include an analytical path-integral method based on Kleinert's variational perturbation theory, and integrated path-integral free-energy perturbation and umbrella sampling (PI-FEP/UM) simulation. Studies have shown that quantum mechanics can be applied to biocatalytic systems in a variety of ways and scales. We hope that the methods presented in this article can further expand the scope of quantum mechanical applications to biomolecular systems

## 4.1.     INTRODUCTION

Quantum mechanics is essential for studying enzymatic processes [1–3]. Depending on the specific problem of interest, there are different requirements on the level of theory used and the scale of treatment involved. This ranges from the simplest cluster representation of the active site, modeled by the most accurate quantum chemical methods, to a hybrid description of the biomacromolecular catalyst by quantum mechanics and molecular mechanics (QM/MM) [1], to the full treatment of the entire enzyme-solvent system by a fully quantum-mechanical force field [4–8]. In addition, the time-evolution of the macromolecular system can be modeled purely by classical mechanics in molecular dynamicssimulations, whereas the explicit incorporation

of nuclear quantum effects, including zero-point energy, non-diabatic coupling and quantum mechanical tunneling, is sometimes needed to adequately understand the dynamics and function of an enzyme, particularly when reactions involving electron or light atom transfer are considered [2]. In this article, we highlight some of the approaches that have been developed in our laboratory for the study of enzymatic reactions.

Because of the complexity of biomolecular systems, it is important to analyze the individual situation to achieve the best possible result for a specific question. More often than not, one needs to balance the level of theory to obtain an accurate representation of the potential energy surface and the time scale in dynamics simulation to gain adequate sampling of the conformational space [1]. In principle, a higher level of quantum mechanical theory can yield a more reliable description of the potential energy surface, but its accuracy is limited by the scale of dynamics sampling of the enzyme conformational space. On the other hand, a well-calibrated lower-level model, specifically parameterized to describe a given chemical reaction, although not systematic but accurate, can be used to extensively sample the dynamic fluctuations of the enzyme environment to yield the potential of mean force, or the free energy reaction profile, along a reaction pathway. In some cases in which the change in protein conformation is relatively small and the dominant factor that determines the chemical reactivity is the local environment of the active site, a simple cluster model may be sufficient to probe the reaction mechanism. However, in most situations where the involvement of the protein dynamics is an integral part of the reaction pathway, the determination of the potential of mean force becomes essential [1–3]. We present a series of hierarchical quantum mechanical methods that can be used to study the dynamics and mechanism of enzymatic reactions. Our treatment features the explicit representation of part or all of the enzyme-solvent system by quantum mechanical method in which the potential energy and forces are determined on the fly during molecular dynamics simulations. This distinguishes from other models which employ a molecular mechanics force field, parameterized from QM/MM calculations of part or the entire system, and follow a predefined reaction path that do not vary during free energy simulations.

In this chapter, we first introduce two new electronic structure methods which can increase the accuracy in the description of the potential energy surface for chemical reactions in biocatalysis. The first is a novel mixed molecular orbital-valence bond (MOVB) theory [9–11] and the other features a next-generation force field based on electronic structure theory [4–8]. Next, we present two methodologies for the treatment of nuclear quantum effects, including an analytical and automated integration-free path-integral (AIF-PI) method based on Kleinert's variational perturbation theory [12, 13]. This path-integral methods provide a systematic way of computing zero-point energies and tunneling effects. The second path-integral method is an integrated path-integral free-energy perturbation and umbrella sampling (PI-FEP/UM) approach for computing kinetic isotope effects. Finally this paper is concluded with highlights of major findings.

## 4.2. THEORETICAL BACKGROUND – PATH INTEGRAL QUANTUM TRANSITION STATE THEORY

We begin our discussion with path integral quantum transition state theory (QTST) [14], which is the theoretical model that we use to model enzymatic reactions. In QTST, the exact rate constant is expressed by the QTST rate constant, $k_{QTST}$, multiplied by a transmission coefficient $\gamma_q$:

$$k = \gamma_q \cdot k_{QTST} \tag{4-1}$$

where the QTST rate constant is given by

$$k_{QTST} = \frac{1}{2} <|\dot{z}|>_{z^{\neq}} e^{-\beta w(z^{\neq})} / \int_{-\infty}^{z^{\neq}} dz \, e^{-\beta w(z)} \tag{4-2}$$

where $\beta = 1/k_B T$ with $k_B$ being Boltzmann's constant and $T$ the temperature. In Eq. (4-2), $w(z)$ is the potential of mean force (PMF) as a function of the centroid reaction coordinate $z[\bar{\mathbf{r}}]$, $z^{\neq}$ is the value of $z[\bar{\mathbf{r}}]$ at the maximum of the PMF, and $<|\dot{z}|>_{z^{\neq}} = (2/\pi \beta M_{eff})^{1/2}$ is a dynamical frequency factor approximated by the velocity for a free particle of effective mass $M_{eff}$ along the reaction coordinate $z[\bar{\mathbf{r}}]$ direction.

In Eq. (4-2), we have introduced a key quantity, the centroid reaction coordinate, $z[\bar{\mathbf{r}}]$, which is a function of the centroid positions of quantized particles in the discrete Feynman path integral representation. In this approach, a quantized nucleus, $(n)$, is represented by a ring of $P$ quasiparticles, whose coordinates are denoted as $\mathbf{r}^{(n)} = \{\mathbf{r}_i^{(n)}; i = 1, \cdots, P\}$ and they are connected harmonically with their neighbors by a force constant of $k^{(n)} = 2\pi P/\beta \lambda_n^2$ where $\lambda_n = (2\pi \beta \hbar^2/M_n)^{1/2}$ is the de Broglie thermal wavelength of atom $n$ with a mass of $M_n$. The classical correspondence of atomic coordinates is the geometrical center of the quasiparticles, i.e., the centroid position:

$$\bar{\mathbf{r}}^{(n)} = \frac{1}{P} \sum_{i=1}^{P} \mathbf{r}_i^{(n)} \tag{4-3}$$

where the superscript $(n)$ specifies the $n$th quantized atom. In practice, the discretization parameter $P$ is chosen to be sufficiently large such that the numerical results converge to the quantum limit.

Although Eqs. (4-1) and (4-2) have identical expressions as that of the classical rate constant, there is no variational upper bound in the QTST rate constant because the quantum transmission coefficient $\gamma_q$ may be either greater than or less than one. There is no practical procedure to compute the quantum transmission coefficient $\gamma_q$. For a model reaction with a parabolic barrier along the reaction coordinate coupled to a bath of harmonic oscillators, the quantum transmission

coefficient is the Grote-Hynes (GH) classical transmission coefficient $\kappa_{GH}$. Often, the classical $\gamma_q$ is used to approximate the quantum transmission coefficient; however, there is no correspondence between classical and quantum dynamic trajectories and the effects of tunneling may greatly affect reaction dynamics near the barrier top.

As in classical transition state theory, the PMF, $w(z)$, can be computed from the equilibrium average:

$$\Delta F(z) = w(z) - w(z_R) = -\frac{1}{\beta} \ln \frac{<\delta(z[\bar{\mathbf{r}}] - z)>}{<\delta(z[\bar{\mathbf{r}}] - z_R)>} \tag{4-4}$$

Where $\Delta F(z)$ is the free energy at $z$ relative to that at the reactant state minimum $z_R$, and the ensemble average $< \cdots >$ is obtained by a quantum mechanical effective potential [15]. Note that the inherent nature of quantum mechanics is at odds with a potential of mean force as a function of a finite reaction coordinate. Nevertheless, the reaction coordinate function $z[\bar{\mathbf{r}}]$ can be evaluated from the path centroids $\bar{\mathbf{r}}$, first recognized by Feynman and Hibbs as the most classical-like variable in quantum statistical mechanics and later explored by many researchers [14, 15].

The calculation of the potential of mean force, $\Delta F(z)$, along the reaction coordinate $z$, requires statistical sampling by Monte Carlo or molecular dynamics simulations that incorporate nuclear quantum effects employing an adequate potential energy function. In our approach, we use combined QM/MM methods to describe the potential energy function and Feynman path integral approaches to model nuclear quantum effects.

## 4.3.    POTENTIAL ENERGY SURFACE BASED ON QM METHODS

The potential energy function describes the energetic change as a function of the variation in atomic coordinates as a result of the Born-Oppenheimer approximation, which is adequate for treating almost all enzymatic reactions [3]. The potential energy function for a given chemical reaction can be modeled by analytical functions, as in molecular mechanics, fitted to reproduce key energetic, structural, and force constant data, from either experiments or high-level ab initio calculations. Alternatively, quantum mechanical methods can be used to model the chemical transformation directly, and depending on the problem of interest, combined QM/MM or fully quantum mechanical methods can be used. We first, discuss the use of a combined QM/MM method employing a mixed molecular orbital and valence bond (MOVB) theory to represent the active site and a MM force field to approximate the rest of the system [9–11]. Then, we consider a quantum mechanics-based explicit polarization (X-Pol) potential in which an electronic structure method is used to treat the entire protein solvent system [4–8].

### 4.3.1. Combined QM/MM Potentials

In combined QM/MM potentials, the system is divided into a QM region and an MM region. The QM region typically includes atoms that are directly involved in the chemical step and they are treated explicitly by a quantum mechanical electronic structure method. The MM region consists of the rest of the system and is approximated by an MM force field. The QM/MM potential is given by:

$$U_{tot} = \langle \Psi(S) | H_{qm}^o(S) + H_{qm/mm}(S) | \Psi(S) \rangle + U_{mm} \tag{4-5}$$

where $H_{qm}^o(S)$ is the Hamiltonian of the QM-subsystem (typically consisting of the substrate and key amino acid residues), $U_{mm}$ is the classical (MM) potential energy of the remainder of the system, and $H_{qm/mm}(S)$ is the QM/MM interaction Hamiltonian between the two regions. In Eq. (4-5), $\Psi(S)$ is the molecular wavefunction of the QM-subsystem that minimizes the energy of the effective Hamiltonian $H_{qm}^o(S) + H_{qm/mm}(S)$ involving electronic degrees of freedom.

Equation (4-5) can be directly utilized in statistical mechanical Monte Carlo and molecular dynamics simulations by choosing an appropriate QM model, balancing computational efficiency and accuracy, and MM force fields for biomacromolecules and the solvent water. Our group has extensively explored various QM/MM methods using different quantum models, ranging from semiempirical methods to ab initio molecular orbital and valence bond theories to density functional theory, applied to a wide range of applications in chemistry and biology. Some of these studies have been discussed before and they are not emphasized in this article. We focus on developments that have not been often discussed.

#### 4.3.1.1. Dual-Level Method

It is useful to rewrite Eq. (4-5) as follows:

$$U_{tot} = E_{qm}^o(S) + \Delta E_{qm/mm}(S) + U_{mm} \tag{4-6}$$

where $E_{qm}^o(S)$ is the energy of an isolated QM-subsystem,

$$E_{qm}^o(S) = \langle \Psi^o(S) | H_{qm}^o(S) | \Psi^o(S) \rangle \tag{4-7}$$

and $\Delta E_{qm/mm}(S)$ is the *interaction energy* between the QM and MM regions, corresponding to the energy change of transferring the QM subsystem from the gas phase into the enzyme environment:

$$\Delta E_{qm/mm}(S) = \langle \Psi(S) | H_{qm}^o(S) + H_{qm/mm}(S) | \Psi(S) \rangle - E_g^o(S) \tag{4-8}$$

Equation (4-6) is especially useful in that the total energy of a hybrid QM and MM system is separated into two "independent" terms – the gas-phase energy and

the interaction energy – which can now be evaluated using different methods or at different levels of QM theory. If we use a high-level (HL) theory to describe the intrinsic energies of the QM system ($E_{qm}^{HL}(S) = E_{qm}^o(S)$), and a low-level (LL) method to carry out molecular dynamics simulations to determine the QM-MM interaction energies ($\Delta E_{qm/mm}^{LL}(S) = \Delta E_{qm/mm}(S)$), we can obtain a highly accurate potential surface for the enzyme reaction to be investigated:

$$U_{tot}^{DL} = E_{qm}^{HL}(S) + \Delta E_{qm/mm}^{LL}(S) + U_{mm} \qquad (4\text{-}9)$$

where $U_{tot}^{DL}$ is the dual-level (DL) total potential energy. In practice, it is more convenient to obtain the total QM/MM energy at a given level, $E_g^o(S) + \Delta E_{qm/mm}^{LL}(S)$. Thus, the following expression is used, which effectively makes an HL correction to the LL quantum mechanical energy.

$$U_{tot}^{DL} = \Delta E_{qm}^{HL}(S) + U_{tot}^{LL} \qquad (4\text{-}10)$$

where $\Delta E_{qm}^{HL}(S) = E_{qm}^{HL}(S) - E_g^o(S)$, which may be used as a post-priori correction to an LL combined QM/MM simulation.

Usually, the semiempirical Austin Model 1 (AM1), parameterization model 3 (PM3), or other semiempirical approaches such as the tight-binding density functional theory (DFTB) model are used in connection with a force field. We have used the standard CHARMM22 force field in our applications. The higher-level method is chosen on the basis of the size of the QM subsystem and the accuracy of the QM model for a particular system. Typically, at least MP2 or G2 quality of method or a tested DFT model is used with a large basis set. The good news is that only a small number of QM calculations are needed to obtain $\Delta E_{qm}^{HL}(S)$, whereas millions of energy and gradients evaluations are need to carry out molecular dynamics simulations using $U_{tot}$ at tractable computational costs.

It should be emphasized that the term $\Delta E_{qm/mm}(S)$ in a combined QM/MM potential contains empirical parameters, and should be optimized to describe accurately QM/MM interactions. By systematically optimizing the associated van der Waals parameters for the "QM-atoms", both semiempirical and ab initio (Hartree-Fock) QM/MM potentials can yield excellent results for hydrogen-bonding and dispersion interactions in comparison with experimental data.[1] The use of semiempirical methods such as the Austin Model 1 (AM1) or Parameterized Model 3 (PM3) in QM/MM simulations has been validated through extensive studies of a variety of properties and molecular systems, including computations of free energies of solvation and polarization energies of organic compounds, the free energy profiles for organic reactions, and the effects of solvation on molecular structures and on electronic transitions.

### 4.3.2. The MOVB Method

An alternative approach that was developed in our laboratory to model the potential energy surface for chemical reactions is the mixed molecular orbital and valence bond (MOVB) theory, in which effective diabatic states are constructed by a block-localized wave function (BLW) method [9–11]. In this approach, molecular orbitals (MOs) are strictly localized within individual fragments of a molecular system based on the reactant or product state configuration. The adiabatic ground state, and excited states if desired, potential energy surface is obtained by avoided coupling of the two or more diabatic states. Key features of the MOVB theory include (1) that the MOs within each fragment are orthogonal, which makes computation efficient, and (2) that the MOs between different fragments are nonorthogonal, which retains important characteristics of valence bond (VB) theory.

In the MOVB method, we use one Slater determinant with block-localized molecular orbitals to define individual VB configuration, called diabatic state. For example, the reactant state of the $S_N2$ reaction between $HS^-$ and $CH_3Cl$ is defined as the Lewis bond structure of the substrate {$CH_3Cl$}:

$$[HS^-]\ [CH_3Cl]; \qquad \Psi^R_{MOVB} = \hat{A}\{\chi^R_{HS}\ \chi^R_{CH_3Cl}\} \qquad (4\text{-}11)$$

where $\Psi^R_{MOVB}$ is the wave function of the reactant diabatic state, $\hat{A}$ is an antisymmetrizing operator, $\chi^R_{HS}$ and $\chi^R_{CH_3Cl}$ are products of molecular orbitals for fragments [$HS^-$] and [$CH_3Cl$], respectively. The product state is similarly defined by

$$[HSCH_3]\ [Cl^-]; \qquad \Psi^P_{MOVB} = \hat{A}\{\chi^P_{HSCH_3}\chi^P_{Cl}\} \qquad (4\text{-}12)$$

where $\chi^P_{HSCH_3}$ and $\chi^P_{Cl}$ are products of molecular orbitals for fragments [$HSCH_3$] and [$Cl^-$].

The MOVB wave function for the adiabatic ground state is written as a linear combination of the diabatic states in Eqs. (4-11) and (4-12):

$$\Phi_{MOVB} = a^R \Psi^R_{MOVB} + a^P \Psi^P_{MOVB} \qquad (4\text{-}13)$$

where $a^R$ and $a^P$ are the configurational coefficients for the reactant and product diabatic state, respectively. The reactant and product diabatic states $\Psi^R_{MOVB}$ and $\Psi^P_{MOVB}$ can be individually optimized, giving rise to the variational diabatic configurations (VDC) [9, 10]. When the wave function of Eq. (4-13) is variationally optimized to yield the minimum energy of the adiabatic ground state, the resulting MOVB diabatic states are called the consistent diabatic configurations (CDC) [11].

The CDC-MOVB method is the appropriate computational approach for studying properties associated with the adiabatic ground state such as the reaction barrier for a chemical reaction and the solvent reorganization energy.

While MOVB can yield reasonable energetic results and an excellent description of the overall potential energy surface on diabatic states and the adiabatic ground

state, the computed barrier is typically a few kilocaleries per mole higher than high-level ab initio results. To overcome this problem in MOVB, we introduce a scaling parameter to adjust the off-diagonal Hamiltonian matrix element $H_{12}$.

$$H_{12}^{\text{EH}} = \beta H_{12} \qquad (4\text{-}14)$$

where $H_{12}^{\text{EH}}$ the effective Hamiltonian (EH) off-diagonal matrix element, and $\beta$ is the diabatic coupling scaling constant to account for static correlations that are not fully accounted for in MOVB states and dynamic correlations that are not included. It is also possible to match the experimental value in this EH-MOVB approach by introducing a shift-parameter $\Delta\varepsilon$ in the diagonal matrix element $H_{22}^{P}$.

In the EH-MOVB model, the energy of the diabatic ground state is determined by using the modified secular equation:

$$\begin{vmatrix} H_{11}^{\text{R}} - V & \beta H_{12} - S_{12} V \\ \beta H_{12} - S_{12} V & H_{22}^{\text{P}} + \Delta\varepsilon - V \end{vmatrix} = 0 \qquad (4\text{-}15)$$

With the introduction of two parameters in Eq. (4-15), the EH-MOVB method can be calibrated to reproduce exactly the experimental barrier height and the desired reaction energy.

### 4.3.3.  The Electronic Structure-Based Explicit Polarization (X-Pol) Potential

Essentially all molecular dynamics simulations of biomacromolecules utilize molecular mechanics force fields to approximate the potential energy surface of the system. In fact, the formalisms of the current force fields for biomolecular systems were already established by Lifson's group in the 1960s [16, 17], which have hardly changed in the past 40 years. Despite the success of molecular mechanics in biomacromolecular modeling, there are also many shortcomings, including redundancy of empirical parameters and a lack of unified treatment of electronic polarization. Combined QM/MM methods allow chemical processes to be readily studied, but its accuracy is ultimately restricted by the use of molecular mechanics force field. Going beyond the molecular mechanical representation of biomolecular systems, we have introduced an explicit polarization (X-Pol) model based on quantum mechanics as a framework for development of next-generation force fields [4–8].

In the X-Pol potential, a molecular system is partitioned into fragments, such as an individual solvent molecule or a peptide unit or a group of such entities. The electronic interactions within each fragment are treated using electronic structure theory described by a Slater determinant wave function, while the inter-fragment electrostatic interactions are treated using a quantal analog of the combined QM/MM approach. The X-Pol potential is designed with a series of hierarchical approximations. The first approximation is that the wavefunction of the full system ($\Phi$) is represented by a Hartree product of the antisymmetric wavefunctions of the individual fragment:

$$\Phi = \prod_{I=1}^{N} \Psi_I \tag{4-16}$$

where the individual molecular wave function $\{\Psi_I; I = 1, \cdots, N\}$ is a Slater determinant. This is equivalent to a subsystem division in density functional theory. This assumption made in Eq. (4-16) neglects the exchange correlation interactions between fragments, and the entire system does not satisfy the Pauli exclusion principle. To account for the short-range exchange repulsion and the long-range dispersion interactions, the popular Lennard-Jones potential is used

$$E_{IJ}^{vdW} = \sum_{\alpha=1}^{A} \sum_{\beta=1}^{B} 4\varepsilon_{\alpha\beta} \left[ \left( \frac{\sigma_{\alpha\beta}}{R_{\alpha\beta}} \right)^{12} - \left( \frac{\sigma_{\alpha\beta}}{R_{\alpha\beta}} \right)^{6} \right] \tag{4-17}$$

where $A$ and $B$ are the number of atoms in fragments $I$ and $J$, and the parameters $\varepsilon_{\alpha\beta}$ and $\sigma_{\alpha\beta}$ are obtained using standard combining rules such that $\varepsilon_{\alpha\beta} = (\varepsilon_\alpha \varepsilon_\beta)^{1/2}$ and $\sigma_{\alpha\beta} = (\sigma_\alpha + \sigma_\beta)/2$, in which $\varepsilon$ and $\sigma$ are atomic empirical parameters that are treatment as in a typical MM force field. The total potential energy of the system is

$$E_{tot} = <\Phi|\hat{H}|\Phi> - \sum_{I=1}^{N} E_I^o \tag{4-18}$$

where $\hat{H}$ is the Hamiltonian of the system defined below, $E_I^o$ is the energy of an isolated fragment $I$ whose wavefunction is $\Psi_I^o$. $E_I^o$ has a constant value and is used for setting the zero of energy of the system corresponding to that of non-interacting fragments.

The Hamiltonian of the system is given as follows:

$$\hat{H} = \sum_{I=1}^{N} \hat{H}_I^o + \frac{1}{2} \sum_{I=1}^{N} \sum_{J \neq I}^{N} \hat{H}_{IJ} \tag{4-19}$$

where $\hat{H}_I^o$ is the electronic Hamiltonian of fragment $I$, and $\hat{H}_{IJ}$ describe the interactions between fragments $I$ and $J$, which is given below:

$$\hat{H}_{IJ}(\Psi_J) = - \sum_{i=1}^{2M} V_i(\Psi_J) + \sum_{\alpha=1}^{A} [Z_\alpha(I) V_\alpha(\Psi_J) + E_{IJ}^{vdW}] \tag{4-20}$$

Here, $Z_\alpha(I)$ is the nucleus charge of atom $\alpha$ and $2M$ is the number of electrons in molecule $I$, $V_t(\Psi_J)$ is the electrostatic potential of molecular $J$ at either the electronic ($t = i$) or nuclear ($t = \alpha$) positions of molecule $I$. The electrostatic potential due to molecule $J$ is defined as follows:

$$V_t(\Psi_J) = -\int \frac{d\mathbf{r}\rho_J(\mathbf{r})}{|\mathbf{r}_t - \mathbf{r}|} + \sum_{\beta=1}^{B} \frac{Z_\beta(J)}{|\mathbf{r}_t - \mathbf{R}_\beta(J)|} \tag{4-21}$$

where $\rho_J(\mathbf{r})$ is the electron density of molecule $J$, derived from the molecular wavefunction, $\rho_J(\mathbf{r}) = |\Psi_J(\mathbf{r})|^2$.

At this point, the theory is completely general in that both molecular orbital theory (MOT) and density functional theory (DFT), ab initio or semiempirical, can be used to obtain the energy in Eq. (4-18). In the latter case, Eq. (4-18) ought be replaced by a DFT energy expression. It is clear that it is particularly straightforward to use DFT theory to define the Hamiltonian expressed in Eqs. (4-19–4-21). Without further approximation, it is necessary to compute the two-electron integrals arising from different molecules, which would be too expensive for a force field designed for condensed phase simulations.

Consequently, we introduce the second approximation which is to use an approximate electrostatic potential in Eq.(4-21) to determine inter-fragment electronic interaction energies. Thus, the electronic integrals in Eq. (4-21) are expressed as a multipole expansion on molecule $J$, whose formalisms are not detailed here. If we only use the monopole term, i.e., partial atomic charges, the interaction Hamiltonian is simply given as follows:

$$\hat{H}_{IJ}(\Psi_J) = -\sum_{i=1}^{2M} \sum_{\beta=1}^{B} \frac{e \cdot q_\beta(\Psi_J)}{|\mathbf{r}_i - \mathbf{R}_\beta(J)|} + \sum_{\alpha=1}^{A} \sum_{\beta=1}^{B} \frac{Z_\alpha(I)q_\beta(\Psi_J)}{R_{\alpha\beta}} + E_{IJ}^{vdW} \tag{4-22}$$

where $q_\beta(\Psi_J)$ is the partial atomic charge on atom $\beta$ of molecule $J$, fitted to reproduce the electrostatic potential of Eq. (4-21) from the wavefunction $\Psi_J$, and $R_{\alpha\beta}$ is the distance between two atoms. Alternatively, we have shown that intermolecular interactions can be adequately described simply by scaling the Mulliken population charges in the simulation of liquid water. Obviously, Eq. (4-22) is the familiar QM/MM interaction Hamiltonian employing a fixed-charge MM force field.

The third approximation is the specific quantum mechanical model to be used for a given problem and a specific purpose. If one has a large amount of dispensable computer time, one would chose to use ab initio MOT or a DFT model to carry out tens to hundreds of millions energy and gradient evaluations for systems of 10,000–100,000 atoms. Unfortunately, it is unlikely to be feasible in the near future. Alternatively, one can adopt a semiempirical approach such as models based on the neglect diatomic differential overlap (NDDO) approximation in MOT or the extended Huckel theory/tight-binding formulation of DFT (DFTB). We choose the NDDO approach in this discussion because there is a well-defined and systematic procedure to evaluate and parameterize the electronic integrals through the Dewar-Thiel type of multipole expansion.

For proteins, the fragment is defined by the peptide unit convention as recommended by the International Union of Pure and Applied Chemistry, although the

*Figure 4-1.* Schematic illustration of the partition between two residues along a polypeptide chain in the X-pol force field

residue convention is typically used in MM force fields. Figure 4-1 shows the division of a peptide chain into peptide units at a $C_\alpha$ carbon; each peptide unit is defined as a quantum mechanical (QM) fragment in the present calculation, and the $C_\alpha$ atom is called a boundary atom. Only the valence electrons of the boundary atom are treated explicitly, so the effective nuclear charge is four and the associated number of electrons is also four. Both the nuclear charge and electrons are divided equally into the two neighboring fragments. The boundary atom has four hybrid bonding orbitals, such that each fragment has two of them as active orbitals and the other two as auxiliary orbitals. This partition of a polypeptide results in two "pseudo atoms", which have identical coordinates, and each of which is half of a boundary $C_\alpha$ atom.

The X-Pol potential has been tested and applied to the simulation of liquid water[5] and liquid hydrogen fluoride, and it has been recently extended to treat fragments that are covalently bonded to one another [6]. An analytic gradient technique has been developed and implemented into the program CHARMM within the NDDO-MOT framework [7], which was used to carry out a molecular dynamical simulation of a small protein, the bovine pancreatic trypsin inhibitor (BPTI) in aqueous solution with periodic boundary conditions. We note that the total computational time is linear scaling by $S \times N \times O(N_{max}^3)$, where $S$ is the number of iterations in system SCF, and $O(N_{max}^3)$ is the computing efforts for the largest residue. The difference between electronic structure calculations for a molecule of the size of $\sum_I^N N_I$ orbitals and that of $N$ separate calculations of the size of $N_{max}$ is obvious because the former would scale as $O([\sum_I^N N_I]^3)$ due at least to diagonalization. These studies have demonstrated the feasibility that a fully QM-based explicit polarization (X-Pol) force field can be used to model biomolecular systems.

## 4.4. PATH INTEGRAL METHODS FOR THE TREATMENT OF NUCLEAR QUANTUM EFFECTS

Proton, hydride and hydrogen atom transfer reactions are ubiquitous in biological processes, and because of their relatively small mass, zero-point energy and quantum tunneling are significant in determining free energy reaction barriers. The incorporation of nuclear quantum effects (NQE) is also important for reactions involving

heavy atoms since one of the most direct experimental assessment of the transition state and the mechanism of a chemical reaction is by measurements of kinetic isotope effects (KIE) [18], which is of quantum mechanical origin. The challenge to theory is the difficulty to accurately determine the small difference in free energy of activation due to isotope replacements. This is further exacerbated by the complexity and size of an enzyme system that requires statistical averaging. In this section, we present two path-integral methods for the treatment of nuclear quantum effects. The first algorithm employs Kleinert's variational perturbation theory at the second order in which the path integrals have been integrated analytically [12, 13]. In the second approach, an integrated path-integral free-energy perturbation and umbrella-sampling (PI-FEP/UM) approach is described to compute KIEs for chemical reactions [19, 20].

### 4.4.1.    The Kleinert Variational Perturbation Theory

The path-integral (PI) representation of the quantum canonical partition function $Q_{QM}$ for a quantized particle can be written in terms of the effective centroid potential $W$ as a classical configuration integral:

$$Q_{QM} = \left( \sqrt{\frac{Mk_BT}{2\pi\hbar^2}} \right)^3 \int_{-\infty}^{\infty} e^{-\beta W(\bar{\mathbf{r}})} d\bar{\mathbf{r}}, \tag{4-23}$$

where $M$ is mass, $\hbar$ is Planck's constant divided by $2\pi$, and $\bar{\mathbf{r}}$ is a point in space. Given $W$, thermodynamic and quantum dynamic quantities can be accurately estimated. Feynman and Kleinert, and Giachetti and Tognetti (FKGT) independently described a variational approach, which yields an upper bound of $W$. The method is highly accurate and has been applied to a variety of systems. It was recently used to model quantum dynamic processes in condensed phases such as water and helium. Kleinert's variational perturbation theory (KP theory) is a more complete and systematic approach, in which the FKGT method is just the first order term of KP theory. Kleinert showed in several model systems that KP theory is exponentially and uniformly convergent. For example, the electronic ground state energy of a hydrogen atom can be determined by computing the $W$ at the zero-temperature limit, and the accuracy of the first three orders of the KP expansion is 85, 95, and 98%, respectively.

Details of the derivation for Kleinert's variational perturbation (KP) theory can be found elsewhere [21]. The $n$th-order KP approximation $W_\Omega^n(\bar{\mathbf{r}})$ to the centroid potential $W(\bar{\mathbf{r}})$ is given by

$$-\beta W_\Omega^n(\bar{\mathbf{r}}) = \ln Q_\Omega(\bar{\mathbf{r}}) - \frac{1}{\hbar} \left\langle V_p^{\bar{\mathbf{r}}}[\mathbf{r}(\tau_1)] \right\rangle_{\Omega,c}^{\bar{\mathbf{r}}} + \frac{1}{2!\hbar^2} \left\langle V_p^{\bar{\mathbf{r}}}[\mathbf{r}(\tau_1)] V_p^{\bar{\mathbf{r}}}[\mathbf{r}(\tau_2)] \right\rangle_{\Omega,c}^{\bar{\mathbf{r}}}$$
$$+ \cdots + \frac{(-1)^n}{n!\hbar^n} \left\langle \prod_{k=1}^{n} V_p^{\bar{\mathbf{r}}}[\mathbf{r}(\tau_k)] \right\rangle_{\Omega,c}^{\bar{\mathbf{r}}} \tag{4-24}$$

where $\tau$ is an imaginary time, the angular frequency $\Omega$ is a variational parameter, which is introduced to define the perturbation potential $V_p^{\bar{r}}[\mathbf{r}] = U(\mathbf{r}, S) - U_\Omega^{\bar{r}}(\mathbf{r})$ about the reference state at $\bar{\mathbf{r}}$ with the harmonic potential $U_\Omega^{\bar{r}}(\mathbf{r}) = (3M/2)\Omega^2 (\mathbf{r} - \bar{\mathbf{r}})^2$. The quantum partition function of the harmonic reference state $Q_\Omega(\bar{\mathbf{r}})$ is given as follows:

$$Q_\Omega(\bar{\mathbf{r}}) = \left[ \frac{\beta\hbar\Omega(\bar{\mathbf{r}})/2}{\sinh(\beta\hbar\Omega(\bar{\mathbf{r}})/2)} \right]^3 \tag{4-25}$$

The remaining terms in Eq. (4-24) are the $n$th-order corrections to approximate the real system, in which the expectation value $\langle \cdots \rangle_{\Omega,c}^{\bar{r}}$ is called cumulant, which can be written in terms of the standard expectation value $\langle \cdots \rangle_\Omega^{\bar{r}}$ by cumulant expansion in terms of Gaussian smearing convolution integrals:

$$\left\langle \prod_{k=1}^n F[\mathbf{r}(\tau_k)] \right\rangle_\Omega^{\bar{r}} = \frac{\prod_{j=1}^n \int_0^{\beta\hbar} d\tau_j \prod_{k=1}^n \int_{-\infty}^\infty d\mathbf{r}_k F(\mathbf{r}_k(\tau_k))}{\{(2\pi)^n Det[a_{\tau_k\tau_{k'}}^2(\Omega)]\}^{3/2}}$$
$$\times \exp\left\{ -\frac{1}{2} \sum_{\substack{k=1 \\ k'=1}}^n (\mathbf{r}_k - \bar{\mathbf{r}})a_{\tau_k\tau_{k'}}^{-2}(\Omega)(\mathbf{r}_{k'} - \bar{\mathbf{r}}) \right\} \tag{4-26}$$

where $Det[a_{\tau_k\tau_{k'}}^2(\Omega)]$ is the determinant of the n×n matrix consisting of the Gaussian width $a_{\tau_k\tau_{k'}}^2(\Omega)$, $a_{\tau_k\tau_{k'}}^{-2}(\Omega)$ is an element of the inverse matrix of the Gaussian width:

$$a_{\tau_k\tau_{k'}}^2(\Omega) = \frac{1}{\beta M\Omega^2} \left\{ \frac{\beta\hbar\Omega \cosh[(|\tau_k - \tau_{k'}| - \beta\hbar/2)\Omega]}{2\sinh[\beta\hbar\Omega/2]} - 1 \right\} \tag{4-27}$$

As $n$ tends to infinity, $W_n^\Omega(x_0)$ is independent of the $\Omega$. At a given order $n$, the optimal frequency for $\Omega$ is given by the least dependence of $W_n^\Omega(x_0)$ on $\Omega$, i.e., zeroing the lowest order derivative of $W_n^\Omega(x_0)$ w.r.t. $\Omega$. Hence $\Omega$ is a variational parameter.

An especially attractive feature of Eq. (4-24) is that if the real system potential is expressed as a series of polynomials or Gaussians, analytic expressions can be obtained, making the computation extremely efficient because the time-demanding Monte Carlo samplings could be avoided (hereafter, the level of calculations up to $n$th order KP expansion for an $m$th-order-polynomial potential is denoted as KP$n$/P$m$) [12]. For other potentials, KP$n$ theory still involves elaborate $n$-dimensional space-time ($2n$ degrees of freedom) smearing integrals in Eq. (4-27). The intricacy of the smearing integrals increases tremendously for multidimensional potentials, where $\Omega$ becomes a $3N \times 3N$ matrix $\Omega_{ij}$ for $N$ nuclei. This complexity is a major factor limiting applications of the KP theory beyond KP1, the original FK approach.

### 4.4.1.1.    *An Automated Computational Procedure*

To render the KP theory feasible for many-body systems with $N$ particles, we make the approximation of independent instantaneous normal mode (INM) coordinates $\{q^{x_0}\}^{3N}$ for a given configuration $\{x_0\}^{3N}$ [12, 13]. Hence the multidimensional $V$ effectively reduces to $3N$ one-dimensional potentials along each normal mode coordinate. Note that INM are naturally decoupled through the 2nd order Taylor expansion. The INM approximation has also been used elsewhere. This approximation is particularly suited for the KP theory because of the exponential decaying property of the Gaussian convolution integrals in Eq. (4-26). The total effective centroid potential for $N$ nuclei can be simplified as:

$$W^n_{\Omega_i}\left(\{\bar{\mathbf{r}}_i\}^{3N}\right) \approx U\left(\{\bar{\mathbf{r}}_i\}^{3N}, S\right) + \sum_{i=1}^{3N} W^n_{\Omega_i}\left(\{q_i\}^{3N}, S\right) \qquad (4\text{-}28)$$

where $W^n_{\Omega_i}(\{q_i\}^{3N}, S)$ is the centroid potential for normal mode $i$. Although the INM approximation sacrifices some accuracy, in return, it allows analyses of quantum mechanical vibration and tunneling, and their separate contributions to the $W$. Positive and negative values of $w_i$ raise (vibration) and lower (tunneling) the original potential, respectively.

   To obtain analytical expressions for the expectation values in Eq. (4-26), we use an $m$th order polynomial (P$m$) to approximate the potential along $q_i$. Note that analytical results for P4 have been used by Kleinert for a quadratic-quartic anharmonic potential and a double-well potential [21]; however, higher order polynomials are needed to achieve the desired accuracy in real systems. We have thus derived the analytical closed forms of Eq. (4-26) up to P20 [12, 13]. Consequently, the $W$ as a function of an arbitrary $\Omega$ can be promptly obtained. This provides a convenient way to determine the least dependent $\Omega$ value without computing the complicated smearing integrals Eq. (4-26) iteratively for different trial values of $\Omega$ by Monte Carlo simulations. In fact, after the interpolating potential along each instantaneous normal mode is determined, there is little computational cost for obtaining the $W$. Thereby, high level ab initio or density-functional (DFT) methods can be used to evaluate the potential energy function for path-integral calculations [12].

   The computational procedure for obtaining the first and second order KP approximations to the centroid potential is summarized below:

1. For each classical configuration $\{x_0\}^{3N}$, the mass-scaled Hessian matrix is diagonalized to obtain a set of normal mode coordinates $\{q^{x_0}\}^{3N}$.
2. The real potential $V$ is scanned from the configuration $\{x_0\}^{3N}$ along each $q_i^{x_0}$ for 10 points both in the forward and backward directions. We found that a step size of 0.1 Å is usually a reasonable choice to yield $W$ within a few percent of the exact.

3. Each centroid potential $w_i^{\Omega}(q_i^{x_0})$ as a function of $\Omega_i$ is readily obtained using the analytical expressions of KP1/P20 or KP2/P20. Note that the path integrals for these polynomials have been analytically integrated.

4. The values of $w_{i,n}^{\Omega}(q_i^{x_0})$ are determined by numerically locating the least dependence of $w_{i,n}^{\Omega}(q_i^{x_0})$ on $\Omega_i$, i.e., zeroing the lowest order derivative of $w_{i,n}^{\Omega}(q_i^{x_0})$ w.r.t. $\Omega_i$ (1st derivative for KP1 and usually 2nd derivative for KP2).

The procedure presented above is integration-free and essentially automated. We hope it could be used by non-path-integral experts or experimentalists as a "black-box" for any given system [12, 13].

### 4.4.2. The Integrated Path-Integral Free-Energy Perturbation and Umbrella-Sampling (PI-FEP/UM) Method

The discrete Feynman path integral method has been used in a variety of applications since it offers an efficient and general approach for treating nuclear quantum effects in condensed phase simulations. In principle, centroid path integral simulations can be directly used to determine KIEs by carrying out two separate calculations for the heavy and light isotope, respectively; however, the convergence of the computed free energy barrier from dynamics simulations is typically not sufficient to ensure the desired accuracy for KIE, especially when heavy isotopes and secondary effects are involved. The integrated path-integral free-energy perturbation and umbrella-sampling (PI-FEP/UM) method [19, 20] involves two computational steps, which has been explored previously in the work of Sprik et al. [22] and in the quantized classical path (QCP) method by Warshel and coworkers [23]. First, classical molecular dynamics simulation is carried out to obtain the potential of mean force along the reaction coordinate for a given reaction. Then, centroid path integral simulations are performed to determine the nuclear quantum effects. The most significant feature of these studies is that classical and quantum simulations are fully separated, making it particularly attractive and efficient for enzymatic reactions. The special feature in the PI-FEP/UM method is to use a free energy perturbation scheme to obtain accurate KIEs for chemical reactions, by changing the atomic mass from one isotope into another in path integral sampling [19].

In centroid path integral, the canonical QM partition function of a hybrid quantum and classical system, consisting of one quantized atom for convenience, can be written as follows:

$$Q_P^{qm} = \frac{1}{\Omega} \int d\mathbf{S} \int d\mathbf{s} \left( \frac{P}{\lambda_M^2} \right)^{3P/2} \int d\mathbf{R}\, e^{-\beta V_{eff}(\{\mathbf{r}\},\mathbf{S})} \qquad (4\text{-}29)$$

where $\Omega$ is the volume element of classical particles, $P$ is the number of quasiparticles, $V_{eff}(\{\mathbf{r}\}, \mathbf{S})$ is the effective potential [15], $\int d\mathbf{R} = \int d\mathbf{r}_1 \cdots \int d\mathbf{r}_P \delta(\bar{\mathbf{r}} = \mathbf{s})$ in which the delta function $\delta(\bar{\mathbf{r}} = \mathbf{s})$ is introduced to make correspondence of the centroid variable $\bar{\mathbf{r}}$ in Feynman path integral to the classical position $\mathbf{s}$. Equation (4-29) can be rewritten exactly as a double average

$$Q_{qm} = Q_{cm} <<e^{-\beta\Delta\bar{U}(\bar{\mathbf{r}},\mathbf{S})}>_{FP,\bar{\mathbf{r}}}>_U \tag{4-30}$$

where the average $<\cdots>_U$ is a purely classical ensemble average obtained according the potential $U_{tot}(\bar{\mathbf{r}}, \mathbf{S})$, and the average differential potential is given by

$$\Delta\bar{U}(\bar{\mathbf{r}}, \mathbf{S}) = \frac{1}{P}\sum_i^P \{U(\mathbf{r}_i, \mathbf{S}) - U(\bar{\mathbf{r}}, \mathbf{S})\} \tag{4-31}$$

The inner average $<\cdots>_{FP,\bar{\mathbf{r}}}$ in Eq. (4-24) represents a free-particle sampling:

$$<\cdots>_{FP,\bar{\mathbf{r}}} = \frac{\int d\mathbf{R}\{\cdots\}e^{-(\pi P/\lambda_M^2)\sum_i^P(\Delta\mathbf{r}_i)^2}}{\int d\mathbf{R}e^{-(\pi P/\lambda_M^2)\sum_i^P(\Delta\mathbf{r}_i)^2}} \tag{4-32}$$

where $\Delta\mathbf{r}_i = \mathbf{r}_i - \mathbf{r}_{i+1}$. In Eq. (4-30), the factor $Q_{cm}$ is the classical partition function given as follows [15]:

$$Q_P^{cm} = \frac{1}{\Omega}\int d\mathbf{S}\int d\mathbf{s}\, e^{-\beta U(\mathbf{s},\mathbf{S})}\left(\frac{P}{\lambda^2}\right)^{3P/2}\int d\mathbf{R}\, e^{-(\pi P/\lambda_M^2)\sum_i^P(\Delta\mathbf{r}_i)^2} \tag{4-33}$$

where we have defined the position of the quantized particle centroid to coincide with the coordinates of the corresponding classical particle $\mathbf{s} = \bar{r}$.

The double average of Eq. (4-30), which is equivalent to Eq. (4-23), is the theoretical basis in the simulation approaches of Sprik et al. [22] called the hybrid classical and path integral, of Hwang and Warshel, called QCP [23], and, of Major and Gao, called PI-FEP/UM [19, 20]. The expression of Eq. (4-30) is particularly useful because the quantum free energy of the system can be obtained first by carrying out classical trajectories according to the classical distribution, $\exp[-\beta U(\bar{\mathbf{r}}, \mathbf{S})]$, and then, by determining the quantum contributions through free particle sampling based on the distribution $\exp[-\beta(\pi P/\beta\lambda_M^2)\sum_i^P(\Delta\mathbf{r}_i)^2]$.

### 4.4.2.1.    Kinetic Isotope Effects

The ratio of the quantum partition functions (Eq. (4-29)) for two different isotopes can be obtained directly through free energy perturbation (FEP) theory by perturbing the mass from the light isotope to the heavy isotope. Consequently, only one simulation of a given isotopic reaction is performed, while the ratio of the partition function, i.e., the KIE, to a different isotopic reaction, is obtained by FEP. This is conceptually and practically an entirely different approach than that used previously [23].

Specifically, following the rate expression of QTST in Eq. (4-1) and assuming the quantum transmission coefficients the dynamic frequency factors are the same, the kinetic isotope effect between two isopotic reactions $L$ and $H$ is rewritten in terms of the ratio of the partial partition functions at the centroid reactant and transition state

and is given by [19, 20]:

$$KIE = \frac{k^L}{k^H} = \left[ \frac{Q_{qm}^L(\bar{z}_L^{\neq})}{Q_{qm}^H(\bar{z}_H^{\neq})} \right] \left[ \frac{Q_{qm}^H(\bar{z}_H^R)}{Q_{qm}^L(\bar{z}_L^R)} \right] e^{-\beta\{F_{CPI,L}^R(\bar{z}_L^R) - F_{CPI,H}^R(\bar{z}_H^R)\}} \tag{4-34}$$

where the ratio of the partition function can be written as follows:

$$\frac{Q_{qm}^H(\bar{z})}{Q_{qm}^L(\bar{z})} = \frac{<\delta(z - \bar{z}) < e^{-\frac{\beta}{P}\sum_i \Delta U_i^{L \to H}} e^{-\beta\Delta\bar{U}_L} >_{FP,L} >_U}{<\delta(z - \bar{z}) e^{-\beta[F_L(\bar{z},S) - F_{FP}^o]} >_U} \tag{4-35}$$

where the subscripts $L$ specifies that the ensemble averages are done using the light isotope. $F_{FP}^o$ is the free energy of the free particle reference state for the quantized particles, and $\Delta U_i^{L \to H} = U(\mathbf{r}_{i,H}) - U(\mathbf{r}_{i,L})$ represents the difference in "classical" potential energy at the heavy and light bead positions $\mathbf{r}_{i,H}$ and $\mathbf{r}_{i,L}$. In the bisection sampling scheme, the perturbed heavy isotope positions are related to the lighter ones by [19]

$$\frac{\mathbf{r}_{i,L}}{\mathbf{r}_{i,H}} = \frac{\lambda_{M_L}\theta_i}{\lambda_{M_H}\theta_i} = \sqrt{\frac{M_H}{M_L}}; i = 1, 2, \cdots, P \tag{4-36}$$

where $\mathbf{r}_{i,L}$ and $\mathbf{r}_{i,H}$ are the coordinates for bead $i$ of the corresponding light and heavy isotopes, $\lambda_{M_L}$ and $\lambda_{M_H}$ are isotopic masses for the light and heavy nuclei, and $\theta_i$ is the position vector in the bisection sampling scheme which depends on the previous sequence of directions and has been fully described in reference. Equation (4-36) indicates that the position vectors for the corresponding heavy and light isotope beads in the path integral simulation are identical, thereby, resulting in the relationship that beads positions are solely determined by the ratio of the square roots of masses.

In Eq. (4-35), we obtain the free energy (inner average) difference between the heavy and light isotopes by carrying out the bisection path integral sampling with the light atom and then perturbing the heavy isotope positions according to Eq. (4-36). Then, the free energy difference between the light and heavy isotope ensembles is weighted by a Boltzmann factor for each quantized configuration (outer average).

## 4.5.    ILLUSTRATIVE APPLICATIONS

### 4.5.1.    Nucleophilic Substitution Reaction of Hydrosulfide Ion and Chloromethane

The VBSCF and EH-MOVB potential energy surfaces for the nucleophilic substitution reaction of HS$^-$ and CH$_3$Cl are depicted in Figure 4-2. The energy contours determined using the EH-MOVB method (Figure 4-2A) are found to be in good accord

(A)



(B)



*Figure 4-2.* Computed potential energy surface from (**A**) ab initio valence-bond self-consistent field (VB-SCF) and (**B**) the effective Hamiltonian molecular-orbital and valence-bond (EH-MOVB) methods for the $S_N2$ reaction between $HS^-$ and $CH_3Cl$

with the VBSCF(6) results (Figure 4-2B). Clearly, the shapes of these potential energy surfaces are in good agreement between the EH-MOVB and VBSCF models, although the energy contours for product diabatic state from VBSCF(6) calculations appears to be somewhat more tightly grouped.

### 4.5.2.    Bovine Pancreatic Trypsin Inhibitor (BPTI) Simulations

The representation of the X-Pol potential is illustrated in Figure 4-3 by the electron density isosurface of BPTI in water. For clarity, water molecules have been deleted. In X-Pol, the potential energy is determined by quantum mechanics. The main result of the BPTI simulations is the extent of electronic polarization and intramolecular charge transfer in the solvated protein. The net charge from Mulliken population analysis of the wave function for each carbonyl group (C=O) in the protein backbone is calculated and averaged over MD simulations. The average net charge on the backbone carbonyl (C=O) group of each residue along the peptide chain is shown in Figure 4-4. We found that all carbonyl groups bear a negative net charge, which is reasonable since C=O is a strong electron withdrawing group. The average net charges on the carbonyl groups range from $-0.04$ to $-0.15$ (all partial charges are in units of a proton charge), with 20 of them more negative than $-0.10$. In comparison, the CHARMM22 force field employs fixed partial atomic charges with the



*Figure 4-3.* Illustration of the X-pol potential represented by the electron density isosurface of BTPI. The color scheme is used purely for distinction of different residues

*Figure 4-4.* Average partial charges (e) on the carbonyl group of along the peptide chain of BPTI

convention that the net group charge for each carbonyl unit (C=O) is zero in the protein backbone. Since it is computationally efficient for each group charge to be zero, this can only be remedied in conventional molecular mechanics calculations by using larger units as groups. However, even if that is done, the charge on each carbonyl group would be independent of time and environment, neither of which is found to be the case in the X-Pol calculations.

### 4.5.3.    Collinear H+H$_2$ Reaction

Both tunneling and vibrational quantum effects are important in the collinear H+H$_2$ reaction. The accuracy of the present KP2 approach is presented as a quantum correction factor $\kappa$, defined as the ratio of the quantum rate constant to that of classical transition state theory with quantum vibrational partition functions but neglecting tunneling effects. We estimated the values of $\kappa$ by applying the path-integral quantum transition state theory with no correction for the transition state re-crossing. Exact values of $\kappa_{QM}$ for a temperature range from $T = 200$ to $1000\,\text{K}$ have been reported on the Porter-Karplus potential. The exact $\kappa_{QM}$, and the present $\kappa_{KP1}$ and $\kappa_{KP2}$ values are listed in Table 4-1 [12, 13]. Both KP1 and KP2 results are in reasonable accord with the exact results at high temperatures, whereas KP2 is still good at lower temperature 200 K, at which $\kappa_{KP1}$ has much greater deviations.

*Table 4-1.* Quantum correction factor $\kappa$ for the collinear reaction between H and H$_2$

| Temperature (K) | $\kappa$ | | |
| --- | --- | --- | --- |
| | Exact | KP1 | KP2 |
| 1000 | 1.5 | 1.2 | 1.2 |
| 600 | 2.5 | 1.7 | 1.7 |
| 300 | 8.7 | 5.3 | 6.6 |
| 200 | 46 | 15 | 55 |

### 4.5.4. The Decarboxylation of N-Methyl Picolinate

Both the $^{12}C/^{13}C$ primary KIE and the $^{14}N/^{15}N$ secondary KIE have been determined (Table 4-2) [19, 20], with the immediate adjacent atoms about the isotopic substitution site quantized as well. To our knowledge, we are not aware of any such simulations prior to our work for a condensed phase reaction with converged secondary heavy isotope effects. This demonstrates the applicability and accuracy of the PI-FEP/UM method.

*Table 4-2.* Computed and experimental primary $^{12}C/^{13}C$ and secondary $^{14}N/^{15}N$ kinetic isotope effects for the decarboxylation of N-methyl picolinate at 25 °C in water

|  | $^{12}k/^{13}k$ | $^{14}k/^{15}k$ |
| --- | --- | --- |
| Exp (120°C) | 1.0212 ± 0.0002 | 1.0053 ± 0.0002 |
| Exp (25°C) | 1.0281 ± 0.0003 | 1.0070 ± 0.0003 |
| PI-FEP/UM | 1.0345 ± 0.0028 | 1.0083 ± 0.0016 |

## 4.6. CONCLUDING REMARKS

In this chapter, we have discussed two new electronic structure methods, which are the X-POL [4–8] and the EH-MOVB methods [9–11], and two path integral techniques to treat nuclear quantum effects, which include an analytical path-integral method AIF-PI [12, 13], and in integrated path-integral free-energy perturbation and umbrella sampling (PI-FEP/UM) simulation [19, 20]. We have demonstrated the applicability of the X-Pol potential to a solvated protein in water. The calculation involves a molecular dynamics simulation of a 14281-atom system, consisting of about 30000 basis functions, with direct dynamics based on an explicit quantum mechanical electronic wave function for the entire system plus a van der Waals term for interfragment exchange repulsion and dispersion forces.

In the EH-MOVB method, we introduced a diabatic coupling scaling factor to uniformly scale the ab initio off-diagonal matrix element $H_{12}$ such that the computed energy of reaction from the EH-MOVB method can be adjusted in exact agreement with the target value, either directly from experiment or from high-level ab initio calculations. Furthermore, the relative energy between the reactant and product diabatic states in the EH-MOVB method can be improved by adding a constant value to the potential energy surface of the diagonal matrix element. This method was illustrated

by the nucleophilic substitution reaction between hydrosulfide and chloromethane in comparison with results from ab initio valence bond self-consistent field calculations.

In the AIF-PI method, a major achievement is to use a polynomial interpolation of the potential on each instantaneous normal mode (INM) coordinate for many-body systems to derive analytic expressions for the path integrals. This makes the KP theory to be efficient beyond the first-order approximation and applicable to realistic systems. The implementation is sufficiently general for any systems described by smooth internuclear potential energy functions. In addition, we presented a free energy perturbation approach on nuclear masses in path integral simulations to obtain accurate and converged kinetic isotope effects for condensed phase reactions. These studies show that our methods are accurate and systematic for computing zero-point energy, quantum partition function, and tunneling effect in fundamental systems. These methods have also been applied to reproduce a series of experimental KIEs. Although the INM approximation in the AIF-PI method neglects correlations between normal modes, it provides further insights into quantum contributions from vibration and tunneling. To go beyond the INM approximation, we are developing a formalism to systematically couple the INM.

Clearly, quantum mechanics can be applied to biocatalytic systems in a variety of ways and scales. We hope that the methods presented in this article can further expand the scope of applications to biomolecular systems.

## ACKNOWLEDGMENTS

## REFERENCES

1. Gao J, Ma S, Major DT, Nam K, Pu J, Truhlar DG (2006) Mechanisms and free energies of enzymatic reactions. Chem Rev 106(8):3188–3209
2. Pu J, Gao J, Truhlar DG (2006) Multidimensional tunneling, recrossing, and the transmission coefficient for enzymatic reactions. Chem Rev 106(8):3140–3169
3. Garcia-Viloca M, Gao J, Karplus M, Truhlar DG (2004) How enzymes work: Analysis by modern rate theory and computer simulations. Science (Washington, DC) 303(5655):186–195
4. Gao J (1997) Toward a molecular orbital derived empirical potential for liquid simulations. J Phys Chem B 101(4):657–663
5. Gao J (1998) A molecular-orbital derived polarization potential for liquid water. J Chem Phys 109(6):2346–2354
6. Xie W, Gao J (2007) Design of a next generation force field: the X-POL potential. J Chem Theory Comput 3(6):1890–1900
7. Xie W, Song L, Truhlar DG, Gao J (2008) The variational explicit polarization potential and analytical first derivative of energy: towards a next generation force field. J Chem Phys 128(23):234108

8. Xie W, Song L, Truhlar DG, Gao J (2008) Incorporation of QM/MM buffer zone in the variational double self-consistent field method. J Phys Chem B 112(45):14124–14131

9. Mo Y, Gao J (2000) Ab initio QM/MM simulations with a molecular orbital-valence bond (MOVB) method: application to an SN2 reaction in water. J Comput Chem 21(16):1458–1469

10. Mo Y, Gao J, (2000) An ab initio molecular orbital-valence bond (MOVB) method for simulating chemical reactions in solution. J Phys Chem A 104(13):3012–3020

11. Song L, Gao J (2008) On the construction of diabatic and adiabatic potential energy surfaces based on ab initio valence bond theory. J Phys Chem A ASAP

12. Wong K-Y, Gao J (2007) An automated integration-free path-integral method based on Kleinert's variational perturbation theory. J Chem Phys 127(21): 211103

13. Wong K-Y, Gao J (2008) Systematic approach for computing zero-point energy, quantum partition function, and tunneling effect based on Kleinert's variational perturbation theory. J Chem Theory Comput 4(9):1409–1422

14. Jang S, Voth GA (2001) A relationship between centroid dynamics and path integral quantum transition state theory. J Chem Phys 112(8747–8757): Erratum: 114, 1944

15. Feynman RP, Hibbs AR (1965) Quantum Mechanics and Path Integrals. McGraw-Hill: New York, p xiv, 365 p. For the applications in quantum statistics, see chapters 10 and 11; Corrections to the errata in the book: http://www.oberlin.edu/physics/dstyer/FeynmanHibbs/ and http://www.physik. fu-berlin.de/~kleinert/Feynman-Hibbs/

16. Bixon M, Lifson S (1967) Potential functions and conformations in cycloalkanes. Tetrahedron 23(2):769–784

17. Levitt M (2001) The birth of computational structural biology. Nat Struct Biol 8(5):392–393

18. Kohen A, Limbach H-H (2006) Isotope Effects in Chemistry and Biology. Taylor & Francis: Boca Raton, p xiv, 1074 p

19. Major DT, Gao J (2007) An integrated path integral and free-energy perturbation-umbrella sampling method for computing kinetic isotope effects of chemical reactions in solution and in enzymes. J Chem Theory Comput 3:949–960

20. Gao J, Wong K-Y, Major DT (2008) Combined QM/MM and path integral simulations of kinetic isotope effects in the proton transfer reaction between nitroethane and acetate ion in water. J Comput Chem 29:514–522

21. Kleinert H (2004) Path integrals in quantum mechanics, statistics, polymer physics, and financial markets. 3rd edition.; World Scientific: Singapore; River Edge, NJ, p xxvi, 1468 p. For the quantum mechanical integral equation, see Section 1.9; For the variational perturbation theory, see Chapters 3 and 5

22. Sprik M, Klein ML, Chandler D (1985) Phys. ReV. B: Condens. Matter Mater. Phys. 31:4234–4244

23. Hwang J-K, Warshel A (1996) J Am Chem Soc 118:11745–11751

# Part II
# Fast Quantum Models with Empirical Treatments

CHAPTER 5

# TOWARDS AN ACCURATE SEMI-EMPIRICAL MOLECULAR ORBITAL TREATMENT OF COVALENT AND NON-COVALENT BIOLOGICAL INTERACTIONS

JONATHAN P. MCNAMARA AND IAN H. HILLIER

*School of Chemistry, University of Manchester, Oxford Road, Manchester, M13 9PL, UK,*
*e-mail: Ian.Hillier@manchester.ac.uk*

**Abstract:**     Recent developments are described which have allowed computationally rapid semi-empirical molecular orbital methods to make significant advances in modelling biological interactions. Efficient strategies for obtaining the best parameters for use in models such as PM3 and AM1 are discussed. Examples of the use of such parameterised methods to understand phosphoryl transfer reactions, the conformational energetics of carbohydrates and hydrogen tunnelling in enzyme catalysed reactions are described. Recent advances in the development of suitable parameters for transition metals, particularly iron, are described, with associated applications to iron containing proteins such as rubredoxin. The recent development and use of a parameterised PM3 model which includes an empirical correction for dispersive interactions (PM3-D) which is designed to study protein structure-function relationships, is described

**Keywords:**     Semi-empirical MO, PM3, QM/MM. parameters, phosphoryl transfer, carbohydrate, enzyme catalysis, hydrogen tunnelling, iron-sulfur proteins, non-covalent interactions, biomolecules, PM3-D

## 5.1.     INTRODUCTION

Computational chemistry employing both quantum mechanics (QM) and molecular dynamics (MD) is one of the many techniques now being used to relate the structure of enzymes and proteins to their biological function. The motivation behind developing more realistic models of these complex systems is driven by the need to understand their structure and reactivity at real temperatures, often requiring the accurate prediction of complex phenomena such as non-covalent interactions [1], hydrogen tunnelling [2], electron transfer [3], as well as more traditional over-barrier crossings, all averaged over the motions of the macromolecule. In the condensed phase, ensembles of configurations and the associated potential energy surfaces (PES) and reaction pathways, must be considered. In view of this, new techniques such as transition

path sampling [4, 5] which allow the rates of such processes to be estimated (at the expense of sampling many trajectories), are being developed. Indeed, across all areas of biological modelling there is an increasing trend to attempt to calculate actual rate constants [6, 7] thereby making direct contact with experiment.

Central to understanding these processes is the use of QM models to understand the accompanying changes in electronic structure. Even the accurate QM modelling of small gas phase reactions can be far from straightforward, and is always computationally demanding. The extension of these techniques to larger and more complex molecules in the condensed phase is usually *via* a hybrid or multilevel approach (e.g. QM/MM [8, 9] or ONIOM [10]) where the reactive centre is treated at a suitable level of QM and those regions deemed to be less important are treated at a lower level of QM, or using molecular mechanics (MM). For example, within the two-level ONIOM formalism the extrapolated energy $E_{\text{ONIOM}}$, which represents the total energy of the system, is obtained from three independent calculations (Eq. 5-1):

$$E_{\text{ONIOM}} = E_{\text{high,model}} + E_{\text{low,real}} - E_{\text{low,model}} \qquad (5\text{-}1)$$

where $E_{\text{low,real}}$ denotes the energy of the entire system calculated by the low level method and $E_{\text{low,model}}$ and $E_{\text{high,model}}$ denote the energies of the model system determined at the low and high level of theory, respectively. These models of the condensed phase require a balanced approach as far as the accuracy of the various components are concerned. Thus, a high-level treatment of a small part of the system, for example the active site of an enzyme, may not be appropriate if the remainder of the system is not properly modelled. The widely used QM/MM method corresponds to the ONIOM scheme where the high level of theory is QM and the low level is MM. It is usually important to allow the point charges of the MM region to polarize the QM wavefunction, and within ONIOM this is termed electronic embedding (EE) [11]. The variant of this model in which such polarization is ignored is labelled mechanical embedding (ME). ME is appropriate if only steric, rather than electrostatic effects are important. The QM/MM approach, with a variety of different levels of QM, particularly semi-empirical and density functional theory (DFT) methods, has given important insight into an increasing number of enzymatic reactions, which have been extensively reviewed [12].

Central to the construction of these QM/MM models is deciding both the size of the high-level region, and the level of QM to be employed, bearing in mind that the calculation might be required to be performed for a very large number of configurations. For complex biological systems, the ideal strategy would involve a high-level ab initio method (with explicit inclusion of electron correlation effects) to treat the chemically active part of the macromolecule. However, even for the modelling of prototype non-covalent interactions, such as π-stacked structures in DNA base-pairs, it has been found necessary to include a high level of electron correlation [e.g. MP2 or CCSD(T)] to obtain data of chemical accuracy [1]. Although there are developments in MP2 and related methods which promise considerable reductions in computer time [13], it is still necessary to look for less computationally demanding

methods to model both covalent and non-covalent biological interactions. The popularity of DFT methods, evident across most areas of computational chemistry, is also found in biological modelling. There is an attractive reduction in computer time compared to ab initio methods, often with little loss of accuracy. However, in areas particularly important for describing biological interactions, such as describing dispersive interactions [1] and systems involving transition metal centres, DFT methods have been found to be somewhat problematic, often requiring the selection of specific functionals [14]. Furthermore, DFT methods are still too computationally demanding for many problems involving the calculation of actual reaction rates. For this reason, in biological and other areas of condensed phase modelling, there has been a renaissance in the use of semi-empirical molecular orbital (MO) methods [15]. The computational speed of semi-empirical methods (3–4 orders of magnitude faster than DFT) means these methods are amenable to molecular dynamics simulations, involving QM potentials.

Progress in developing semi-empirical methods accelerated during the 1970s when ab initio calculations were beyond the computational resources of most research groups. Nearly all semi-empirical methods developed during this time relied on variants of the zero differential overlap (ZDO) approximation, leading to a reduction in the number of two-electron integrals. The earlier semi-empirical models of Pople and co-workers, such as the CNDO [16] and INDO [17] methods, neglected a significant number of two-electron integrals, and as such, they were unable to properly describe the electronic structure of some chemical systems. These failings prompted Dewar and Thiel to develop the MNDO model (Modified Neglect of Differential Overlap) which retained more two-electron integrals [18, 19]. MNDO was parameterised against experimental data and included several new features such as a modified core-repulsion term to treat N—H and O—H hydrogen-bonding interactions.

The failure, however, of MNDO to properly describe hydrogen-bonds was realised by Dewar who subsequently modified the core-repulsion function by including a set of attractive and repulsive Gaussian functions centred at large and short range internuclear separations. Changes to the MNDO Hamiltonian thus required a different set of parameters and led to AM1 (Austin Model 1) [20]. AM1 was later refined by Stewart who developed an automated optimisation procedure for determining the necessary parameters, the new parameter set and optimisation procedure being denoted PM3 (Third Parameterisation of MNDO) [21]. Rather than to attempt to reproduce the results of Hartree-Fock calculations, AM1 and PM3, by the inclusion of experimental data in the parameter optimisation, sought to reproduce quantities such as structures, dipole moments and heats of reaction. Nowadays, these methods can be routinely used to study systems containing any of the main group elements (excluding the noble gases) whose ground-state valence electronic configuration can be described by *s* and *p* atomic orbitals alone [22].

In recent years, progress has been made in extending these methods to *d*-orbitals (e.g. MNDO/d [23, 24]) so that many of the hypervalent compounds of main group elements can now be treated using these techniques. More importantly, the addition

of *d*-orbitals now means a theoretical framework exists to apply these methods to molecules containing transition metal atoms. In spite of these advances, a complete set of parameters for the transition metals, capable of describing their diverse chemistry, is not currently available. This is probably due to the fact that the development of parameters for metals is considerably more challenging than for the main-group elements due to the complexity of their electronic states and problems of using only a minimal valence basis for their description [25, 26]. Moreover, many metals, particularly those in biological systems, are found in quite different coordination environments [3], and as such, it is unlikely that a general set of parameters (for each metal) can be obtained to be used across a range of chemical systems. To address this central problem, Rossi and Truhlar have suggested the use of specific reaction parameters (SRP) for specific chemical problems [27].

There has been considerable recent activity developing appropriate parameters to allow semi-empirical methods to describe a variety of biologically important systems, and their related properties, such as (i) enzyme reactivity, including both over- and through-barrier processes, (ii) conformations of flexible molecules such as carbohydrates, (iii) reactivity of metalloenzymes and (iv) the prediction of non-covalent interactions by addition of an empirical dispersive correction. In this review, we first outline our developing parameterisation strategy and then discuss progress that has been made in the areas outlined above.

## 5.2.    PARAMETERISATION OF SEMI-EMPIRICAL METHODS

In spite of recent advances in both computer hardware and software the development of new parameter sets for existing semi-empirical methods is a time-consuming process. Sophisticated algorithms are now available which allow large numbers of parameters to be fitted simultaneously to quite large chemical databases. Even so, the careful compilation of reference data and the evaluation of different optimised parameter sets can be time-consuming. Indeed, nearly all parameter optimisation strategies usually require the calculation of many different parameter sets, using different fitting functions and different sets of reference data, until a satisfactory parameter set is obtained. Here, the definition of a "satisfactory" parameter set can mean the one that gives the lowest value of the chosen error function or that which gives the most "balanced" chemical results. Thus the development of parameterisation strategies often requires both efficient optimisation algorithms and substantial chemical intuition. In view of this it is useful to review some of the important considerations in developing a successful semi-empirical parameterisation strategy.

### 5.2.1.    Optimisation Algorithms

Central to the development of new semi-empirical parameter sets is the construction and subsequent minimisation of an appropriate error function, *S* (Eq. 5-2). The error function contains the molecular quantities calculated at the semi-empirical level

($q_i^{semi-empirical}$) and the corresponding reference values ($q_i^{reference}$) which are usually obtained from either experiment or high-level calculation, $w_i$ being an appropriate weight.

$$S = \sum_i w_i \left( q_i^{semi-empirical} - q_i^{reference} \right) \qquad (5\text{-}2)$$

The objective of any parameterisation strategy is to minimise $S$ efficiently, and this can be achieved using a genetic algorithm (GA) or a gradient-based algorithm, and in some cases a combination of the two approaches. GAs have so far been used to obtain semi-empirical parameters for sodium [28], magnesium [29], technetium [30] and iron [31] using reference training sets containing up to a few tens of molecules. On the other hand Stewart obtained the entire PM3 parameter set using an efficient automated gradient-based algorithm, which allowed parameters for different elements to be simultaneously optimised [21, 22]. We too have developed and used a similar algorithm with appropriate modifications to obtain parameters for transition metal atoms [26, 32]. In keeping with the semi-empirical philosophy, a number of important approximations are included within the algorithm and increase the efficiency of the optimisation procedure. The procedure is based upon a modified Broyden-Fletcher-Goldfarb-Shanno (BFGS) method in which one of the central approximations is that a reference structure is taken and the parameters are adjusted until the forces on the atoms of each structure are minimised, rather than performing explicit geometry optimisation. In addition, the error function (Eq. 5-2) and its gradients can be evaluated efficiently since for small changes in $n$ parameters ($x$), any reference property $q$ (e.g. heat of formation) can be approximated using a truncated Taylor expansion centred at the initial parameter values ($x^0$).

$$q(x_1, x_2, \ldots, x_n) = q\left(x_1^0, x_2^0, \ldots, x_n^0\right) + \sum_{i=1}^{n} \left.\frac{\partial q}{\partial x_i}\right|_{x_i=x_i^0} \cdot \left(x_i - x_i^0\right) \qquad (5\text{-}3)$$

We have found that one of the major disadvantages of gradient-based algorithms, (particularly those which avoid explicit geometry optimisation) is that they require quite large databases of reference molecules in order to obtain parameter sets capable of describing different bonding situations for any given element [26, 32, 33]. Also, and perhaps more importantly, the properties of the final parameter set are quite strongly dependent on the choice of initial parameters (for the optimisation procedure) whereas when a GA is used quite different parameter sets can be obtained. For this reason, many recent semi-empirical optimisation algorithms are now making use of a combination of GAs, gradient minimisation methods and other techniques such as Monte Carlo annealing [34–36].

### 5.2.2.   Modification of the Core Repulsion Function

During the evolution of the current semi-empirical methods (AM1, PM3) a number of refinements were made to the core-repulsion function in order to improve, for example, the description of hydrogen-bonding [20, 21]. Such changes have in general led to increased flexibility within the modified semi-empirical Hamiltonians resulting in quite marked improvements in the accuracy of the parent methods.

Thus, the current semi-empirical methods (MNDO, AM1 and PM3) differ in the way in which core-repulsions are treated. Within the MNDO formalism the core-repulsion ($E_{AB}^{MNDO}$) is expressed in terms of two-centre, two-electron integrals (Eq. 5-4), where $Z_A$ and $Z_B$ correspond to the core charges, $R_{AB}$ is the internuclear separation, and $\alpha_A$ and $\alpha_B$ are adjustable parameters in the exponential term [19].

$$E_{AB}^{\mathrm{MNDO}} = Z_A Z_B \left( s^A s^A, s^B s^B \right) \left( 1 + e^{-\alpha_A R_{AB}} + e^{-\alpha_B R_{AB}} \right) \qquad (5\text{-}4)$$

When the core-repulsion function involves either hydrogen-oxygen or hydrogen-nitrogen interactions, a modified form of this function is used (A = H, B = O, N; Eq. 5-5).

$$E_{AB}^{\mathrm{MNDO}} = Z_A Z_B \left( s^A s^A, s^B s^B \right) \left( 1 + R_{AB} e^{-\alpha_B R_{AB}} + e^{-\alpha_A R_{AB}} \right) \qquad (5\text{-}5)$$

Current AM1 and PM3 methods use an alternative core repulsion function which differs from that used in MNDO in that an additional term involving one to four Gaussian functions is used (defined by parameters *a-c*, Eq. 5-6) [20, 21]. These extra terms help to reduce the excessive core-core repulsions just outside bonding distances.

$$E_{AB}^{\mathrm{AM1,PM3}} = E_{AB}^{\mathrm{MNDO}} + \frac{Z_A Z_B}{R_{AB}} \sum_{i=1}^{4} \left[ a_{iA} e^{-b_{iA}(R_{AB}-c_{iA})^2} + a_{iB} e^{-b_{iB}(R_{AB}-c_{iB})^2} \right]$$

$$(5\text{-}6)$$

Despite these modifications there remain a number of well-documented problems with the AM1/PM3 core-repulsion function [37] which has resulted in further refinements. For example, Jorgensen and co-workers have developed the PDDG (pair-wise distance directed Gaussian) PM3 and MNDO methods which display improved accuracy over standard NDDO parameterisations [38]. However, for methods which include *d*-orbitals (e.g. MNDO/d [23, 24], AM1/d [25] and AM1* [39, 40]) it has been found that to obtain the correct balance between attractive and repulsive Coulomb interactions requires an additional adjustable parameter $\rho$ (previously evaluated using the one-centre two-electron integral $G_{\mathrm{ss}}$, Eq. 5-7), which is used in the evaluation of the two-centre two-electron integrals (Eq. 5-8).

$$\rho = \frac{1}{2G_{ss}} \tag{5-7}$$

$$\left(s^A s^A, s^B s^B\right) = \frac{e^2}{\left(R_{AB}^2 + (\rho_A + \rho_B)^2\right)^{1/2}} \tag{5-8}$$

When applying these methods to the study of molybdenum complexes Voityuk and Rösch [25, 41] found that the use of the AM1 core-repulsion function (Eq. 5-6) led to some systematic deviations for some Mo—X bond lengths. To address these problems, important changes to the core-repulsion function were made by the introduction of bond-specific parameters ($\alpha_{Mo-X}$ and $\delta_{Mo-X}$, Eq. 5-9) [22, 25]. The idea of using bond-specific core repulsion parameters is not new, since the AM1 parameterisation of boron used bond-specific Gaussian functions to improve the final results [42].

$$E_{Mo-X}^{AM1/d} = Z_{Mo} Z_X \left(s^A s^A, s^B s^B\right) \left[1 + 2\delta_{Mo-X} e^{-\alpha_{Mo-X} R_{Mo-X}}\right] \tag{5-9}$$

This strategy has been found to be more efficient than using Gaussian functions and has now been used to extend the NDDO-based family of methods to the remaining main group elements [22], the AM1* parameterisation of second row elements [39] and the transition metals titanium and zirconium [40] as well as our work extending the PM3 method to iron [26, 32].

### 5.2.3. Choice of Reference Data and Construction of the Error Function

The reference data for many of the early semi-empirical parameterisations usually involved data collected from a wide range of experimental techniques [19–21]. For more recent parameterisations, particularly for SRP sets, there is now an increasing trend to replace in part, if not entirely the reference training data with information obtained from high-level ab initio [e.g. MP2 or CCSD(T)] or DFT calculations [34, 35, 43]. Such an approach is to be favoured in the absence of, or when there are concerns regarding the reliability of experimental data. The use of high-level calculations is also desirable if the semi-empirical method is required to reproduce stationary structures (intermediates, transition states) along a given reaction pathway [35, 36, 44]. Indeed, the reference data for many such parameterisations can now include actual transition structures as well as reaction energy barriers (calculated using DFT methods) for complex biological processes.

In view of the fact that recent parameterisations make use of reference data from high-level calculations, the corresponding error functions used to develop these methods can in principle involve any given property that can be calculated. Thus, in addition to structural information, the error function can involve atomic charges and spin densities, the $<S^2>$ value for the wavefunction, ionisation potentials and the relative energies of different structures within the reference database [26, 32]. Detailed information concerning the actual wavefunction can be extremely useful for

developing parameter sets for complex systems (e.g. transition metals) since atomic charges, spin densities and $<S^2>$ values can help to direct the fitting algorithm towards the desired electronic states of the reference complexes. Importantly, the performance of this approach, with respect to any chosen molecular property can be tuned by suitable adjustment of the weighting factors ($w_i$, Eq. 5-2).

### 5.3.    HYBRID POTENTIALS FOR THE SIMULATION OF PHOSPHORYL TRANSFER REACTIONS

Modelling biological systems involving macromolecules such as DNA, RNA and proteins requires an accurate description of a range of both subtle and highly specific inter- and intra-molecular interactions. One area receiving considerable attention is the study of biological processes involving phosphorylation and dephosphorylation reactions, which are important in the regulation of many cellular processes [45, 46]. Ab initio [47–50] as well as DFT methods [51–55] have been used extensively to study these phosphoryl transfer reactions using prototypical model systems. It is however important to treat these systems with explicit consideration of the enzyme environment, since there is extensive charge rearrangement during the reaction. In view of this, QM/MM methods employing different levels of QM have been quite widely used [36, 56–58]. Thus, a number of studies of phosphoryl transfer catalysed by protein kinases have been reported, and have focused on the role of the conserved aspartate residue. An early QM/MM study using PM3 with standard parameters questioned whether this residue acted as a catalytic base [57]. This study has been followed by others employing DFT, some of which have come to different conclusions concerning the role of the aspartate [59]. There is thus the need to develop a parameter set specific for this important class of reaction.

Two recent studies have used data from high level ab initio and DFT calculations on model systems to develop improved parameters to describe these reactions. Arantes and Loos developed a new parameterization of NDDO (for hydrogen, sulfur, oxygen, carbon and phosphorus) for phosphate ester hydrolysis and thiolysis by fitting to a database involving the stationary structures from eight phosphate ester reactions [calculated at the MP2/6-311+G(2df,2p)//MP2/6-31+G(d) level] [36]. The new Hamiltonian gave a notable improvement over the standard models as indicated by an RMSD of 2.8 kcal mol$^{-1}$ (reduced from 16.8 kcal mol$^{-1}$ using standard MNDO parameters) for the relative energies of species along the pathways of various phosphate ester reactions. These parameters were then employed in hybrid QM/MM calculation of a potential of mean force (PMF) for the dephosphorylation of phenyl phosphate (catalysed by the dual specificity phophatase VHR1: E $\cdots$ ROPO$_3{}^{2-}$ $\rightarrow$ E—PO$_3{}^{2-}$ + ROH, E = ester) and gave a reaction energy barrier of 16.4 kcal mol$^{-1}$, very close to the value of 15.5 kcal mol$^{-1}$ calculated using transition state theory [60] and an experimentally measured rate constant [61].

The work of Arantes and Loos involved a modified MNDO Hamiltonian in which the valence orbitals of sulfur and phosphorus were described by *s* and *p* orbitals

alone. Although, the MNDO and AM1 schemes have now been extended to include *d*-orbitals on second-row elements as in the MNDO/d scheme of Voityuk and Thiel [62], and the AM1* method of Winget et al. [39], further modification of these methods has been required to accurately model phosphoryl transfer reactions. Nam et al. [35] have developed a modified AM1/d-PhoT Hamiltonian by parameterising against a database of high-level DFT(B3LYP) calculations for RNA catalysis, including geometries and relative energies of minima, transition states and reactive intermediates, dipole moments, proton affinities, and other relevant properties. This work included a necessary modification of the core-repulsion function which involved the use of two atom specific scaling parameters to allow increased flexibility through the attenuation (or elimination) of the Gaussian core-repulsion interactions between certain atoms (refer to Eq. 5-6). The use of the new AM1/d-PhoT model was demonstrated by the computation of PMFs for a model transphosphorylation reaction in water and gave barriers which differed by only 3–5 kcal mol$^{-1}$ from DFT calculations employing an implicit solvent model.

## 5.4. QUANTUM MECHANICAL FORCE FIELDS FOR CARBOHYDRATES

An important step in the understanding of a range of fundamental biological processes is the accurate modelling of carbohydrate structure and dynamics at a molecular level [63]. The inherent flexibility and polar nature of both mono- and polysaccharides has prompted the development of a number of elaborate MM force fields in order to properly describe the subtle stereoelectronic anomeric effects in such molecules [64, 65]. An alternative to the MM approach is to use a semi-empirical method [66], which can allow the incorporation of explicit electronic polarisability, as in the fluctuating charge and dipolar force fields [67]. However, the current parameterizations of the semi-empirical methods, AM1 [20] and PM3 [21], lack sufficient accuracy [68]. For example, PM3 predicts a $^1C_4$ ring conformation as more stable than a $^4C_1$ structure (Figure 5-1) in opposition to high level QM calculations and experiment [69, 70]. Therefore we have developed a set of SRPs for use within the PM3 method by fitting the $U_{ss}$, $U_{pp}$, $\beta_s$, $\beta_p$, $\alpha$ of oxygen and $U_{ss}$, $\beta_s$, $\alpha$ of hydrogen to a set of small molecule carbohydrate analogues (10 conformers of 1,2-ethanediol and 4 conformers of methoxymethanol) calculated at the MP2/cc-pVDZ level [65, 71]. The new parameterisation is denoted PM3CARB-1 and is now available within the AMBER 9 package [72].

Compared to ab initio calculations and existing semi-empirical methods [21] the PM3CARB-1 method performs quite well for the relative energies of the different conformers used in the reference training data (Table 5-1). In fact the RMS of the relative energies (vs. MP2) for the conformers of 1,2-ethanediol is very close to that using low level ab initio calculation [0.51 (HF), 2.56 (PM3), 0.57 kcal mol$^{-1}$ (PM3CARB-1)]. The PM3CARB-1 model has also been assessed for the prediction of the relative energies and structures of a set of 14 glucose conformations

$^4C_1$ chair $\qquad$ $^1C_4$ chair

*Figure 5-1.* Definition of chair conformation in $\beta$-glucopyranose. Reproduced with permission from reference [66]. Copyright Elsevier 2004

(Table 5-2). For these, the PM3CARB-1 model has a RMS error in the relative energies of 1.60 kcal mol$^{-1}$ compared to the B3LYP/6-311+G**//HF/6-31G* calculations. The PM3 model has a corresponding error of 2.69 kcal mol$^{-1}$ (Table 5-2). For the 11 $^4C_1$ conformations, the error is 1.1 kcal mol$^{-1}$ *via* the PM3CARB-1

*Table 5-1.* Relative energies (kcal mol$^{-1}$) of 1,2-ethanediol and methoxymethanol conformations, and RMS error with respect to MP2 level of theory

| Conformer | MP2[a] | HF | AM1 | PM3 | PM3CARB-1 |
|---|---|---|---|---|---|
| *1,2-Ethanediol* | | | | | |
| g$^-$Gg$^-$ | 1.19 | 1.03 | −0.98 | −0.78 | 1.45 |
| gGg | 3.15 | 3.04 | −0.98 | −0.82 | 1.75 |
| gGg$^-$ | 0.30 | 0.69 | −0.88 | −1.39 | 1.10 |
| gTg | 3.00 | 2.59 | 0.39 | −0.43 | 2.87 |
| gTg$^-$ | 2.81 | 2.28 | 0.20 | −0.67 | 2.50 |
| tGg | 3.81 | 3.46 | 2.56 | 1.76 | 3.63 |
| tGg$^-$ | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 |
| tGt | 3.48 | 2.90 | 3.19 | −0.78[b] | 2.56 |
| tTg | 2.92 | 2.08 | 1.55 | 1.07 | 2.40 |
| tTt | 2.87 | 1.69 | 2.98 | 2.56 | 2.23 |
| RMS | | 0.51 | 1.75 | 2.56 | 0.57 |
| *Methoxymethanol* | | | | | |
| Gg | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 |
| Gg$^-$ | 2.06 | 2.89 | 2.36 | 1.50 | 3.98 |
| Tt | 6.80 | 7.30 | 9.43 | 4.09 | 8.81 |
| Tg | 2.73 | 2.66 | 2.74 | 1.24 | 1.88 |
| Total RMS | | 0.49 | 1.53 | 2.28 | 0.86 |

[a] Ab initio QM levels of theory for 1,2-ethanediol are MP2/cc-pVDZ [65] and HF/6-311+G(2d,2p) [71]; for methoxymethanol, MP2/cc-pVDZ and HF/6-311+G(2d,2p) [71].
[b] Unstable – converts to g$^-$Gg$^-$.

*Table 5-2.* Relative energies (kcal mol$^{-1}$) of glucose conformers with RMS error with respect to ab initio QM level of theory

| Conformer[a] | Ab initio[b] | PM3 | PM3CARB-1 |
|---|---|---|---|
| $^4C_1\alpha/cc/g^-g$ | 0.00 | 0.34 | 0.43 |
| $^4C_1\alpha/cc/gg^-$ | 0.01 | −0.11 | 0.35 |
| $^4C_1\alpha/cc/tg$ | 0.00 | 0.00 | 0.00 |
| $^4C_1\alpha/cl/g^-g$ | 0.91 | 0.40 | 2.07 |
| $^4C_1\alpha/cl/tt$ | 1.14 | 2.08 | 3.31 |
| $^4C_1\alpha/cl/tt$ | 1.61 | 2.08 | 3.31 |
| $^4C_1\alpha/cl/tt$ | 1.24 | 2.08 | 3.31 |
| $^4C_1\beta/cc/g^-g$ | 0.74 | 1.22 | 1.88 |
| $^4C_1\beta/cc/gg^-$ | 0.63 | 0.78 | 2.10 |
| $^4C_1\beta/cc/tg$ | 0.85 | −0.67 | 1.34 |
| $^4C_1\beta/cc/tg$ | 1.77 | −0.67 | 1.34 |
| $^1C_4\beta/cc/gg$ | 9.48 | 1.44 | 6.52 |
| $^1C_4\beta/cc/g^-g^-$ | 8.77 | 0.82 | 6.31 |
| $^1C_4\beta/cl/g^-g$ | 8.69 | −2.43 | 4.73 |
| RMS | | 2.69 | 1.60 |

[a] Adapted from reference [65], notation for glucose conformers: orientation of ring/anomer/hydrogen bond network/$CH_2OH$ group ($O_5C_5C_6O_6$ and $C_5C_6O_6H$).
[b] Ab initio QM level of theory B3LYP/6-311+G**//HF/6-31G* [65].

model, which although an improvement over PM3, is still higher than the respective errors of 0.6 kcal mol$^{-1}$ for the OPLS-AA [73], and 0.5 kcal mol$^{-1}$ for the OPLS-AA-SEI force fields [65]. Importantly, the PM3CARB-1 model predicts $^4C_1$ conformations to be more stable than $^1C_4$ conformations. For example, $^1C_4/\beta/cl/g^-g$ conformation is predicted by PM3 to be the global minimum for the 14 conformers. At the PM3CARB-1 level, this conformer is 4.7 kcal mol$^{-1}$ higher in energy than the correct lowest energy $^4C_1$ conformer, in improved agreement with the DFT calculations [65]. Moreover, the clockwise and anti-clockwise hydrogen-bonding networks in each of the 14 conformers are also shown to be in good agreement with the ab initio calculation [66]. The PM3CARB-1 method has been used to study the dynamics of the disaccharide 4-*O*-α-D-xylopyranosyl-α-D-xylopyranose. Here the disaccharide solute was treated at the QM level in the presence of 492 TIP3P water molecules. Previously, hybrid QM/MM dynamics simulations using a PM3/TIP3P potential required ring constraints to prevent unphysical transitions to $^1C_4$ conformations [74]. Use of the PM3CARB-1 potential indicates no such constraints are required (Figure 5-2).

Thus, in contrast to preceding MM approaches explicit treatment of electronic polarisability is integral to a semi-empirical QM approach and promises excellent prospects for quantitative theoretical modelling of carbohydrates across a range of condensed phase environments. The results of the PM3CARB-1 model do however indicate in line with classical force field approaches [65, 73] that perhaps greater

*Figure 5-2.* Plot of $C_5O_5C_1O_1$ angle of non-reducing saccharide residue 4-O-$\alpha$-D-xylopyranosyl-$\alpha$-D-xylopyranose from QM/MM molecular dynamics simulations using (**a**) a PM3/TIP3P potential and (**b**) a PM3CARB-1/TIP3P potential. Reproduced with permission from reference [66]. Copyright Elsevier 2004

flexibility will be required to accurately model the complete manifold of subtle carbohydrate energetics and structure.

## 5.5.    HYDROGEN TUNNELLING AND ENZYME CATALYSIS

Hydrogen motion, $H^+$, $H^-$ or H$^·$, is often involved in the rate-limiting step of many enzyme catalysed reactions. Here, QM tunnelling can be important and is reflected in the values of the measured kinetic isotope effects (KIEs) [75]. Enzyme motion

can in some cases result in local structures which may favour through-barrier tunnelling [76–80]. An accurate calculation of the potential energy surface and the degree of tunnelling is therefore important in understanding which structural features may increase vibrationally enhanced tunnelling. As examples, we studied three enzyme catalysed reactions involving hydride, proton and neutral hydrogen atom transfer using semi-empirical calculations, comparing the calculated KIEs with the available experimental data. The reactions considered are hydride transfer by liver alcohol dehydrogenase (LADH), proton transfer catalysed by methylamine dehydrogenase (MADH) and hydrogen atom transfer catalysed by soybean lipoxygenase (SLO-1) [81].

The rate of hydrogen transfer can be calculated using the direct dynamics approach of Truhlar and co-workers which combines canonical variational transition state theory (CVT) [82, 83] with semi-classical multidimensional tunnelling corrections [84]. The rate constant is calculated using [83]:

$$k(T) = \kappa (k_B T/h) \exp(-\Delta G_T^{CVT}/RT) \qquad (5\text{-}10)$$

where $\Delta G_T{}^{CVT}$ is the CVT free energy of activation and $\kappa$ is a transmission coefficient based on the semi-classical tunnelling approximation. The change in $k(T)$ upon isotopic substitution leads to the calculated KIEs. This approach has been successful in the evaluation of rate constants for small gas phase reactions [85], as well as for some enzymatic systems, and has been recently reviewed [86].

Tunnelling can occur not only through the one-dimensional PES, but also *via* a corner cutting mechanism which is included by using the centrifugal dominant small-curvature tunnelling (SCT) approximation due to Truhlar [84]. Such corner-cutting is evaluated in terms of the reaction path curvature, which gives the degree to which the reaction co-ordinate is coupled to transverse vibrational modes along the minimum energy pathway (MEP). In regions of increased reaction path curvature the effective reduced mass for motion along $s$ is reduced [87], which simulates corner-cutting. Another quantity resulting from these calculations useful for describing the degree and origin of tunnelling is the representative tunnelling energy (RTE) [84], the energy at which tunnelling is most likely to occur at a particular temperature.

The use of Eq. (5-10) to evaluate the reaction rate is characterised by the calculation of Hessians for a large number of points along the MEP which are required to locate the free energy maximum and also to evaluate the curvature required for evaluation of the transmission coefficient. In view of the associated computational expense, high-level electronic structure calculations are not feasible and alternative strategies, one of which is to use a semi-empirical method, are usually employed [81].

Our studies on the three enzymes have involved the use of semi-empirical methods, using published and also SRP parameter sets. For both LADH and MADH (Figures 5-3a and b) hybrid QM/MM models were employed [8, 9, 88–90]. In LADH the PES surface was calculated at the AM1 level [20] but scaled by data from the HF/3-21G surface [91]. The results of the CVT calculation with the SCT correction show quite modest yet contributory degrees of tunnelling, an RTE

*Figure 5-3.* Active site and calculated PES properties for the reactions studied, with the transferring hydrogen labelled as $H_T$: (**a**) hydride transfer in LADH, (**b**) proton transfer in MADH and (**c**) hydrogen atom transfer in SLO-1. (i) potential energy, (ii) vibrationally adiabatic potential energy, (iii) RTE at 300K and (iv) total reaction path curvature. Reproduced with permission from reference [81]. Copyright Elsevier 2002

*Table 5-3.* PES properties for hydrogen transfer in the LADH, MADH and SLO-1 enzymes[a]

|  | LADH | MADH | SLO-1 |
|---|---|---|---|
| Potential energy barrier | 18.6 | 19.1 | 18.2 |
| Activation energy[b] | 14.1 {15.6} | 9.9 {11.3} | 6.9 {1.2} |
| Imaginary frequency | 1229$i$ | 2218$i$ | 2913$i$ |
| RTE[c] | −0.7 | −4.8 | −6.4 |
| $\kappa$ | 3.0 (2.5) | 114 (5.7) | 780 (9.1) |
| $\kappa^{D}$ [d] | 2.7 (2.1) | 44 (3.5) | 325 (5.4) |
| KIE (CVT/SCT)[e] | 3.7 (4.4) {3.8} | 15.6 (9.9) {16.8} | 18.9 (12.5) {56} |

[a]All temperature dependent properties are for 300 K; energies in kcal mol$^{-1}$; frequencies cm$^{-1}$; all values for hydrogen unless otherwise stated. {In parenthesis} experimental results from references [76, 92, 95] for LADH, [78] for MADH and [95] for SLO-1.
[b]Calculated using the CVT/SCT method.
[c]Relative to the maximum of the adiabatic potential.
[d]Using the SCT method. (In parenthesis) Wigner values.
[e](In parenthesis) TST/W results.

within 1 kcal mol$^{-1}$ of the top of the barrier and a calculated deuterium KIE of 3.7 in excellent agreement with the experimental value of $3.78 \pm 0.07$ (Table 5-3) [76, 92, 93]. The measured deuterium KIE for the proton transfer step (rate-limiting) in MADH is considerably greater than for LADH ($16.8 \pm 0.5$, 298 K, Table 5-3) [78]. Here, a set of PM3 SRP parameters were developed using the results of high-level [CCSD(T)//B3LYP/6-311++G**] calculations for a related model system. A GA was used and only $U_{pp}$ and $\beta_p$ of the carbon and oxygen atoms involved in the proton transfer reaction were adjusted [6]. The resulting QM(PM3-SRP)/MM PES was calculated and the final PES, RTE and reaction path curvature are shown in Figure 5-3b. In contrast to LADH (Figure 5-3a), the RTE is considerably below the top of the adiabatic potential and there is enhanced curvature on both reaction and product side of the potential. We find that in contrast to the predictions using the standard PM3 parameters, the calculated KIE and activation energy (15.6, 9.9 kcal mol$^{-1}$, Table 5-3) using the SRP's are now in excellent agreement with the experimental values (16.8, 11.3 kcal mol$^{-1}$) [78].

In the enzyme lipoxygenase, where there is hydrogen atom abstraction from $C_{11}$ of linoleic acid by a hydroxyl group (Figure 5-3c), which is one of the ligands attached to a six coordinate high spin Fe(III), a maximum deuterium KIE of $56 \pm 5$ at 305 K has been reported [94–96]. In this system, we have the problem of the proper QM description of this reaction involving a transition metal. We chose to use a cluster model of the active site, involving Fe(III) in its sextet ground state [97] coordinated to three imidazoles, one acetate and one acetamide ligand (modelling His504, His690, His499, I1e839 and Asn694 in the actual enzyme) together with the OH group. As substrate, we used 1,4-pentadiene, rather than the natural substrate linoleic acid (Figure 5-3c). We used the PM3/d method employing iron parameters reported

elsewhere [98]. We note that the important energetic features of the hydrogen transfer reaction (Figure 5-3c) are predicted to be similar to those for MADH, with an RTE below the top of the barrier by $6.4\,\text{kcal}\,\text{mol}^{-1}$, with considerable curvature on both the reactant and product side of the barrier. The final results are summarised in Figures 5-3 and Table 5-3 and show important differences in the PES for the three different enzyme catalysed reactions.

## 5.6.    TRANSITION METALS

The development of semi-empirical models for the accurate description of the diverse chemistry of transition metal elements is by no means straightforward. The general absence of properly validated parameters for these elements possibly reflects the difficulties associated with firstly, the description of their complex electronic structure using a minimal basis and secondly, the problem of assembling reliable reference data for the parameterisation procedure [25, 26]. In spite of this, promising progress has been made towards developing robust semi-empirical models to treat transition metal systems. Some present-day commercial packages now include semi-empirical Hamiltonians capable of calculating transition metal complexes [99], and some parameterisations have been discussed in the literature. For example, Voityuk and Rösch have published details of an AM1/d parameterisation for molybdenum for the calculation of structures as well as heats of formation, reaction enthalpies and bond energies [25]. In addition, Winget and Clark have extended their AM1* parameterisation to include titanium and zirconium [40] and Tejero et al. have reported a set of SRPs for iron to study hydrogen abstraction catalysed by Soybean Lipoxygenase-1 (SLO-1) [44]. The parameters in this study were fitted using a reference dataset involving five model reactions (including transition states) as well as some prototypical SLO-1 model complexes.

We have used a strategy to derive a parameter set for iron [31] within the PM3 model focussing on the iron-sulfur protein rubredoxin (Rd), an important electron transfer protein [100]. The active site of Rd involves the high spin states of $[\text{Fe(SR)}_4]^{1-/2-}$ where the single iron atom is surrounded by four sulfur atoms from cysteine residues [101]. Experimental adiabatic and vertical electron detachment energies (ADE, VDE) of clusters which mimic the enzyme active site have proved valuable in providing data directly relevant to electron transfer processes in enzymes, particularly relating to reorganization energies which are central to the Marcus theory of electron transfer [102, 103]. These experimental data may be used to both test and to parameterize different computational schemes.

For transition metals compounds, excited states are often extremely important, and this should be reflected in any parameterisation scheme. We have shown that one way such information can be incorporated into a given parameterisation strategy is by fitting the one-centre parameters ($U_{\text{ss}}$, $U_{\text{pp}}$, $U_{\text{dd}}$, $G_{\text{ss}}$, $G_{\text{sd}}$, $G_{\text{dd}}$, $H_{\text{sd}}$) to experimental excitation and ionisation energies of the neutral and charged metal atom (as a starting point for future parameterisations these parameters have been reported elsewhere for

the first-row transition metals [104]). This approach combined with the SRP strategy has been used to obtain a full set of parameters for iron. The remaining two-centre parameters [$\beta_s$, $\beta_p$, $\beta_d$, $\zeta_s$, $\zeta_p$, $\zeta_d$; $\alpha$, $a_1$, $a_2$, $b_1$, $b_2$, $c_1$, $c_2$ (refer to Eq. 5-6)] were fitted using a GA to the structures and energy data obtained from DFT (B3LYP/6-31G*) calculations of the active site analogues of Rd, $[Fe(SCH_3)_4]^{1-/2-}$.

The final PM3 parameter set yields ADE, VDE and reorganisation energies ($\lambda_{oxi}$) (−2.16, −2.64, 0.48 eV) in good agreement with the DFT values (−1.90, −2.17, 0.27 eV); the PM3 Fe—S distances are also very close to the DFT values for the oxidised and reduced states of the model complex [2.31 2.43 (PM3) and 2.32, 2.43 Å (DFT)]. In Table 5-4 the structural and energetic properties are compared for eight redox couples calculated using both DFT and PM3 and also with experiment [102, 103, 105]. The level of agreement between PM3 and DFT is extremely

*Table 5-4.* Comparison of VDEs and ADEs, reorganization energies ($\lambda_{oxi}$) (eV), and Fe—X distances (Å), for Rd analogues

| Structure | | Energy | | | Fe—X | | |
|---|---|---|---|---|---|---|---|
| | | Experiment | DFT | PM3 | | DFT | PM3 |
| $[FeCl_3]^{1-/0}$ [a] | VDE | 4.36 | 4.24 | 3.49 | Reduced | 2.24 | 2.32 |
| | ADE | 4.10 | 4.00 | 3.17 | Oxidized | 2.14 | 2.23 |
| | $\lambda_{oxi}$ | 0.26 | 0.24 | 0.32 | | | |
| $[FeBr_3]^{1-/0}$ [a] | VDE | 4.42 | 4.35 | 4.46 | Reduced | 2.38 | 2.55 |
| | ADE | 4.26 | 4.17 | 4.27 | Oxidized | 2.28 | 2.46 |
| | $\lambda_{oxi}$ | 0.16 | 0.18 | 0.19 | | | |
| $[Fe(SCH_3)_3]^{1-/0}$ [b] | VDE | 3.08 | 2.94 | 2.75 | Reduced | 2.31 | 2.33 |
| | ADE | 2.80 | 2.70 | 2.40 | Oxidized | 2.21 | 2.23 |
| | $\lambda_{oxi}$ | 0.28 | 0.24 | 0.35 | | | |
| $[Fe(SCN)_3]^{1-/0}$ [b] | VDE | 4.96 | | 4.02 | Reduced | | 2.29 |
| | ADE | 4.64 | | 3.74 | Oxidized | | 2.23 |
| | $\lambda_{oxi}$ | 0.32 | | 0.28 | | | |
| $[FeCl_4]^{2-/1-}$ | VDE | | −1.33 | −1.34 | Reduced | 2.41 | 2.42 |
| | ADE | | −1.99 | −2.02 | Oxidized | 2.23 | 2.29 |
| | $\lambda_{oxi}$ | | 0.66 | 0.68 | | | |
| $Fe(SCH_3)_4]^{2-/1-}$ | VDE | | −1.90 | −2.16 | Reduced | 2.43 | 2.43 |
| | ADE | | −2.17 | −2.64 | Oxidized | 2.32 | 2.31 |
| | $\lambda_{oxi}$ | | 0.27 | 0.48 | | | |
| $[Fe(S_2$-$o$-xyl$)_2]^{2-/1-}$ [b] | VDE | 2.60 | | −1.89 | Reduced | 2.32[c] | 2.39 |
| | ADE | 2.30 | | −2.24 | Oxidized | 2.28[c] | 2.29 |
| | $\lambda_{oxi}$ | 0.30 | | 0.35 | | | |
| $[Fe(SC_6H_5)_4]^{2-/1-}$ | VDE | | | −0.42 | Reduced | | 2.38 |
| | ADE | | | −0.78 | Oxidized | | 2.31 |
| | $\lambda_{oxi}$ | | | 0.36 | | | |

[a]Reference [103].
[b]Reference [102].
[c]Experiment [105].

encouraging; the PM3 calculations predict the trend found in the DFT calculations of a significantly greater bond length change upon ionisation in the case of tetrachloride, as well as a larger reorganisation energy. The ADE of $[FeCl_4]^{2-}$ is also calculated to be smaller than that of $[Fe(SC_6H_5)_4]^{2-}$ by 1.2 eV in good agreement with an $E^0$ difference of approximately 1.0 eV, and a difference in VDE of 0.9 eV (compared to an experimental value of $1.4 \pm 0.3$ eV [106]).

A proper description of the transition metal-ligand interactions is also important when the semi-empirical method is used to calculate either the high or low-level region in an ONIOM modelling scheme [10]. To demonstrate this, we have assessed the suitability of the PM3 parameters for describing the oxidised and reduced states of the protein Rd using the ONIOM method with EE [11]. The DFT (B3LYP/6-31G*) and PM3 methods have been used to describe the high-level region (DFT:MM, PM3:MM) as well as in an ONIOM scheme which couples two different levels of QM, DFT and PM3, (DFT:PM3). The X-ray structure of the oxidised form of Rd, (PDB Code: 1IRO) was solvated with 3900 TIP3P waters and the protein relaxed using the AMBER force field [107]. In each ONIOM calculation the $[Fe(SCH_3)_4]$ cluster was taken to be the model system and was described by a QM level of theory (DFT or PM3) where the cysteine residues were terminated between $C_\beta$ and $C_\alpha$ with a hydrogen link atom. For the EE calculations, all residues within a 10 Å radius of the central iron atom were optimised whilst the remainder of the protein and solvating waters were held fixed. In all calculations, both high- and low-level regions of the protein geometry were optimised, the Fe—S bond lengths being given in Table 5-5. These structures were used to calculate ADEs and VDEs for the model region, and the associated inner sphere reorganisation energies (Table 5-6). The energy of the low-level region is not included in these quantities, since this is somewhat unstable, due in part to the lack of consideration of an ensemble of enzyme structures.

We found that the protein effects both the bond length change upon ionisation and the inner sphere reorganisation energy ($\lambda_{oxi}$). The reduction in the Fe—S bond length which is 0.1 Å for the isolated cluster, is reduced to less than 0.1 Å in the enzyme, the models which include PM3 yielding a somewhat smaller reduction than the DFT(EE) model. All models are in good agreement with estimates from EXAFS and X-ray data [108, 109] of a bond length difference between the reduced and oxidised forms of 0.04–0.07 Å. The constraints of the enzyme, which are reflected in

*Table 5-5.* Comparison of Fe—S bond lengths (Å) for oxidised and reduced states of iron-sulfur protein, Rd

|          |                      |                     | ONIOM-EE |        | ONIOM   |
|----------|----------------------|---------------------|----------|--------|---------|
|          | X-Ray[a]             | EXAFS[b]            | DFT:MM   | PM3:MM | DFT:PM3 |
| Oxidised | 2.26                 | 2.29                | 2.31     | 2.34   | 2.29    |
| Reduced  |                      | 2.33                | 2.39     | 2.39   | 2.33    |

[a]Reference [109].
[b]Reference [108].

*Table 5-6.* Comparison of ADEs and VDEs, and inner sphere reorganization energies ($\lambda_{oxi}$) (eV) for iron-sulfur protein, Rd

| | Isolated | | ONIOM-EE | | ONIOM |
|---|---|---|---|---|---|
| | DFT | PM3 | DFT:MM | PM3:MM | DFT:PM3 |
| VDE | −1.91 | −2.16 | −1.91 | −1.62 | −2.79 |
| ADE | −2.18 | −2.64 | −2.00 | −1.99 | −2.90 |
| $\lambda_{oxi}$ | 0.27 | 0.48 | 0.09 | 0.37 | 0.11 |

these bond length changes are also reflected in the reduction in the value of $\lambda_{oxi}$ for the enzyme compared with that for the isolated cluster. The value of this reduction is 0.1–0.2 eV for all three models, which is in line with the DFT:MM value of Sigfridsson et al. [110]. The results of the calculations described here are encouraging and lend support to the idea of developing SRPs for transition metals in specific enzyme environments. The success of the reparameterised PM3 model in treating the effect that different ligands and the environment have on different oxidation states of the core structure of the protein, Rd, justifies further exploration of this cost effective way of modelling complex metalloenzymes.

Since transition metal atoms are found in a diverse range of coordination environments [3], the use of a bond-specific core repulsion function (Eq. 5-9) with increased flexibility can be used to develop more general parameter sets (thus allowing a greater range of compounds to be modelled than with an SRP set), without necessarily a compromise in accuracy [26]. We refined our iron SRPs by addressing two issues, firstly, increasing the size of our reference training set from just two complexes, $[Fe(SCH_3)_4]^{1-/2-}$, to 60 iron complexes, and secondly, introducing a more flexible core-core repulsion function (as in Eq. 5-9) [25, 26]. The new reference training set included a range of high- and low-spin iron complexes (60 in total); 11 standard iron inorganic complexes (e.g. $[FeCl_4]^{1-}$), 14 iron-sulfur protein models (e.g. $[Fe_4S_4(SCH_3)_4]^{2-}$), 19 iron-heme complexes which had either one or two neutral or charged axial ligands (e.g. $[Fe(Por)(Im)_2]$), 16 non-heme iron enzyme models (e.g. $[Fe(SH)_2(\mu\text{-}S_2)Fe(SH)_2]^{2-}$). The final parameter set for iron involved 32 parameters (including 17 bond-specific terms) [26] compared to just 21 parameters in the original SRP set [31]. An impressive aspect of the new parameter set is that it performs quite well in the description of different electronic states of the neutral and charged iron-sulfur cubane models even though the reference training data only involved the complex $[Fe_4S_4(SCH_3)_4]^{2-}$ and the hydrogenase-cluster [111] (Figure 5-4). The structures and spin densities for the different spin states calculated at both DFT [112] and PM3 levels are reported in Table 5-7, and the relative energies of the different spin states in Table 5-8. As far as the structures are concerned, PM3 predicts Fe—Fe, Fe—S*, and Fe—S distances within the ranges calculated at the DFT level. The PM3 spin densities are also consistent with those calculated at the DFT level. The relative energies of these complexes calculated at the PM3 level are
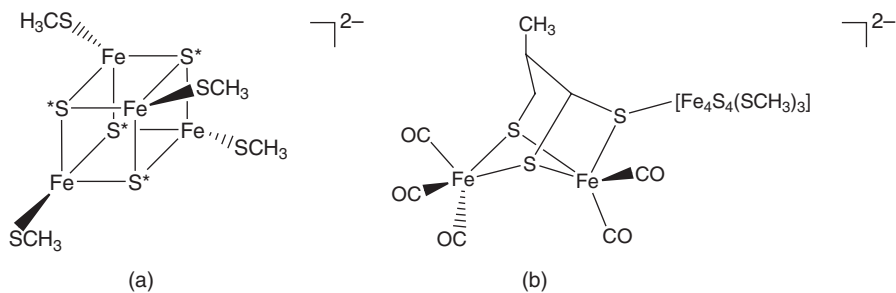
*Figure 5-4.* Structure of (**a**) [Fe$_4$S$_4$(SCH$_3$)$_4$]$^{2-}$ unit and (**b**) hydrogenase-cluster (H-cluster)

also quite promising. PM3 correctly predicts the S=1/2 state of [Fe$_4$S$_4$(SCH$_3$)$_4$]$^{1-}$ to be the lowest in energy and the S=8 state of [Fe$_4$S$_4$(SCH$_3$)$_4$]$^{4-}$ to be the least stable by 11.64 eV (DFT) and 12.10 eV (PM3). The relative energies of the other spin states are quite well reproduced, although their relative ordering, particularly of those states close in energy, does not always agree with the DFT values.

We now review the application of the modified iron parameter set [26] for modelling an actual enzyme catalysed reaction. The hydroxylation by cytochrome P450 is at present receiving considerable attention since this family of enzymes appear in all aerobic bioorganisms and perform vital bioregulatory functions [113]. One of the most important catalytic processes performed by these enzymes is the hydroxylation of C—H bonds. It is now widely accepted that this reaction proceeds by the so-called "rebound" mechanism in which the initial hydrogen abstraction from the alkane (RH) by the iron-oxo species (compound I) is followed by a radical rebound on the iron-hydroxo intermediate to yield the ferric-alcohol complex which then releases the alcohol and restores the resting state (water complex) [114]. This mechanism is evidenced by considerable work in determining kinetic isotope effects (KIEs) which suggest C—H bond breaking to be the rate determining step. In addition, the rebound mechanism has also been the subject of several DFT studies [115, 116].

We have calculated the transition state for the rate-limiting step (RLS) on the low spin (doublet) surface and also the free energy barrier and the corresponding KIE (without tunnelling). The semi-empirical calculations indicate that the PM3 transition state for the RLS differs from the DFT one [115] in that the transferring hydrogen is a little further from the oxo group than in the high-level calculation [1.16 Å (PM3) *vs.* 1.08 Å (DFT)]; in both cases the O—H bond formation is almost complete. The PM3 calculation yields a free energy barrier very close to the DFT one [24.0 (PM3) *vs.* 26.5 kcal mol$^{-1}$ (DFT)] although the associated KIE (k$_H$/k$_D$) is not sufficiently quantitative [9.0 (PM3), 5.1 (DFT), 5.9 $\pm$ 0.35 (experiment)] [115, 117]. These findings, in line with our calculations of proton transfer in MADH, emphasize the need for SRPs to accurately calculate KIEs [6, 81]. We could, most probably, achieve greater accuracy by a more focussed training set, and by some modification of the ligand atom parameters.

*Table 5-7.* Fe−X distances (Å) and Fe spin densities for different states and anions of [Fe$_4$S$_4$(SCH$_3$)$_4$]

| Charge | Total Spin | X | Fe—X | | Spin Density | |
| | | | DFT[a] | PM3 | DFT[a] | PM3 |
| --- | --- | --- | --- | --- | --- | --- |
| 0 | 0 | Fe | 2.61–2.67 | 3.32–3.34 | | 1.17, 1.17, −1.17, −1.17 |
| | | S* | 2.22 | 2.30–2.32 | | |
| | | S | 2.24 | 2.20 | | |
| 1− | 1/2 | Fe | 2.75–2.96 | 2.82–3.20 | 3.26, 3.26, −3.01, −3.01 | 2.84, 3.31, −2.60, −2.28 |
| | | S* | 2.20–2.30 | 2.26–2.54 | | |
| | | S | 2.25 | 2.26–2.32 | | |
| 1− | 9/2 | Fe | 2.77–2.83 | 2.78–3.23 | | 3.37, 2.62, 3.40, −2.66 |
| | | S* | 2.32 | 2.27–2.44 | | |
| | | S | 2.25 | 2.27–2.31 | | |
| 2− | 0 | Fe | 2.74–2.81 | 2.97–3.05 | 3.13, 3.13, −3.13, −3.13 | 3.21, 3.21, −3.21, −3.20 |
| | | S* | 2.22–2.33 | 2.36–2.41 | | |
| | | S | 2.30 | 2.35 | | |
| 3− | 1/2 | Fe | 2.78–2.81 | 2.86–3.19 | 3.29, 3.29, −3.01, −3.01 | 3.40, 3.46, −3.38, −3.38 |
| | | S* | 2.27–2.35 | 2.34–2.51 | | |
| | | S | 2.38 | 2.39–2.45 | | |
| 3− | 7/2 | Fe | 2.64–2.80 | 2.74–3.11 | 3.25, 3.23, 3.20, −3.23 | 3.46, 3.39, 1.93, −3.38 |
| | | S* | 2.32–2.34 | 2.30–2.57 | | |
| | | S | 2.37–2.41 | 2.37–2.45 | | |
| 4− | 0 | Fe | 2.81–2.85 | 2.90–3.21 | 3.19, 3.19, −3.19, −3.19 | 3.40, 3.40, −3.40, −3.40 |
| | | S* | 2.31–2.37 | 2.38–2.45 | | |
| | | S | 2.38 | 2.55 | | |
| 4− | 4 | Fe | 2.63–2.87 | 2.92–3.08 | 3.66, 3.55, 3.36, −3.10 | 3.42, 3.42, 3.41, −3.41 |
| | | S* | 2.36–2.38 | 2.37–2.51 | | |
| | | S | 2.41–2.54 | 2.51–2.56 | | |
| 4− | 8 | Fe | 2.61–2.82 | 2.90–3.16 | 3.50, 3.50, 3.51, 3.51 | 3.42, 3.42, 3.42, 3.42 |
| | | S* | 2.38–2.45 | 2.37–2.52 | | |
| | | S | 2.56 | 2.51 | | |

[a]BV86 with TZ for iron [112].

*Table 5-8.* Relative energies (eV) of spin states of $[Fe_4S_4(SCH_3)_4]$

| Charge | Spin | Relative Energy | |
|--------|------|------|------|
| | | DFT[a] | PM3 |
| 1 – | 1/2 | 0.00 | 0.00 |
| 2 – | 0 | 0.03 | 0.86 |
| 1 – | 9/2 | 0.67 | 0.08 |
| 0 | 0 | 3.05 | 3.81 |
| 3 – | 7/2 | 4.12 | 4.75 |
| 3 – | 1/2 | 4.24 | 4.26 |
| 4 – | 0 | 11.24 | 12.00 |
| 4 – | 4 | 11.26 | 11.90 |
| 4 – | 8 | 11.64 | 12.10 |

[a]BV86 with TZ for iron [112].

## 5.7.    NON-COVALENT INTERACTIONS IN BIOMOLECULES

Another important aspect of protein structure-function relationships involves the proper description of the non-covalent inter-molecular interactions responsible for protein structure and for substrate protein-binding [1]. It is now generally accepted that dispersion interactions, particularly those responsible for π–π stacking inter-action, can be of equal importance to hydrogen bonding interactions [1]. However, the accurate calculation of dispersion interactions at the QM level usually requires computationally expensive high-level ab initio methods [e.g. MP2 or CCSD(T)] and quite large basis sets, thus limiting the size of the systems that can be studied to a few tens of atoms [1]. DFT methods, although computationally cheaper, often fail to properly describe these interactions [118, 119]. In recent years progress towards addressing these issues has been made, firstly through the development of new functionals [120, 121], and secondly, through the addition of an empirical dispersive correction to the normal QM energy. The latter represents a promising solution to this important problem, and for DFT this has been achieved through the addition of an explicit $R^{-6}$ term to describe the inter-atomic dispersion interactions, giving the so-called DFT-D method [118, 122]. Self-consistent-charge density functional tight-binding methods (SCC-DFTB-D) have also been developed [123, 124]. Like many of their DFT counterparts, the most common semi-empirical methods (e.g. MNDO [19], AM1 [20] and PM3 [21]) often fail to adequately describe dispersion interactions.

As far as hydrogen-bonding is concerned, Giese et al. [34] have made important progress towards improving the description of hydrogen-bonding within the PM3 model. They introduced an exploratory semi-empirical Hamiltonian (PM3$_{BP}$) for modelling hydrogen-bonding in nucleic acid base pairs. The actual functional

form of the Hamiltonian remained unchanged and instead the authors chose to optimise the PM3 parameters for hydrogen, oxygen and nitrogen so as to refine the hydrogen bonding interactions, to reproduce experimental base pair dimer enthalpies and values from high-level density functional calculations (using the mPWPW91 exchange-correlation functional [125, 126] and the MIDI! basis set [127]). The new Hamiltonian has shown notable improvements over existing semi-empirical methods and reproduces experimental dimer interaction energies with an accuracy that rivals DFT, with a reduction in computational cost approaching three orders of magnitude. Importantly, this work demonstrated that for such systems (particularly those where dispersion interactions are important) additional ad-hoc corrections are probably required (in addition to new parameter sets) in order for these methods to have sufficient quantitative accuracy to be useful in biological applications.

In view of this and in line with the DFT-D approach described by Grimme [118], we have added an atom-atom pair-wise additive potential of the form $C_6/R^6$ to the usual semi-empirical energy [19–21] in order to account for dispersion effects [43]. Thus the dispersion corrected semi-empirical energy ($E_{PM3-D}$) is now given by;

$$E_{PM3-D} = E_{PM3} + E_{disp} \qquad (5\text{-}11)$$

where $E_{PM3}$ is the normal PM3 energy and $E_{disp}$ is an empirical term containing the dispersion correction.

$$E_{disp} = -s_6 \sum_i \sum_j \frac{C_6^{ij}}{R_{ij}^6} f_{dmp}\left(R_{ij}\right) \qquad (5\text{-}12)$$

Here, the summation is over all atom pairs, $C_6^{ij}$ is the dispersion coefficient for the pair of atoms $i$ and $j$ (calculated from the atomic $C_6$ coefficients), $s_6$ is a scaling factor which is chosen to be 1.4 in line with the value used for the BLYP functional [118], and $R_{ij}$ is the inter-atomic distance between atoms $i$ and $j$. A damping function is used in order to avoid near singularities for small distances, given by;

$$f_{dmp}\left(R_{ij}\right) = \frac{1}{1 + e^{-\alpha(R_{ij}/R_0-1)}} \qquad (5\text{-}13)$$

where $R_0$ is the sum of the atomic van der Waals radii and $\alpha$ is a parameter determining the steepness of the damping function. The atomic $C_6$ coefficients, and the $R_0$ and $\alpha$ values as well as the combination rule for the composite $C_6^{ij}$ coefficients were taken from the work of Grimme [118].

The modification of the semi-empirical Hamiltonian thus required the development of a new set of parameters for use within the PM3 method. Rather than develop a new reference database we chose to use the high-level ab initio calculations recently reported by Hobza and co-workers [1]. These calculations have been collected together in a database which can be used to judge the accuracy of less rigorous, but

also less computationally demanding methods. In all, the database contains some 165 non-covalent complexes, including 128 DNA base pairs, 19 amino acid pairs and 18 other small complexes. For each complex high-level ab initio MP2 and CCSD(T) complete basis set (CBS) interaction energies are reported. This database is sub-divided into a smaller training set containing 22 complexes labelled *S22* which is proposed for initial screening, whilst the remaining 143 complexes are labelled as the *JSCH-2005* database. The latter database involves some 38 hydrogen-bonded DNA base pairs, 32 interstrand base pairs, 54 stacked base pairs and 19 amino acid base pairs. To date these two databases (*S22* and *JSCH-2005*) have been used to assess the performance of various DFT functionals [119] and dispersion corrected DFT-D methods [128–130]. In view of the suggestion by Jurečka et al. [1] that the small molecule database (*S22*) contains all the important biological non-covalent interactions, we chose to use this set of reference molecules to parameterise the dispersion corrected PM3 method, denoted PM3-D. The larger *JSCH-2005* database was then used to extensively test the PM3-D method. Due to the absence of interactions involving sulfur containing complexes in the *S22* dataset, semi-empirical calculations of the sulfur containing complexes in the *JSCH-2005* database are not reported. Therefore our *JSCH-2005* database involves 31 hydrogen-bonded DNA base pairs, 32 interstrand pairs, 54 stacked base pairs and 17 amino acid base pairs; a total of 134 complexes.

In Table 5-9 we report the summary statistics (interaction energies and distances) of the semi-empirical calculations on the complexes in the *S22* database [1]. The deviation of the interaction energies are displayed in Figure 5-5. Clearly, for the method without the dispersive correction (PM3) the interaction energies deviate significantly from the reference values, whereas for the PM3-D method the deviations are much less (Table 5-9 and Figure 5-5). These findings are further supported with a PM3-D mean unsigned error (MUE) considerably less than for the uncorrected methods (Table 5-9).

We now compare the PM3-D method with previous uncorrected DFT calculations on the *S22* complexes [130]. For the dispersion-bonded complexes the errors in the interaction distances for the PBE, B3LYP and TPSS functionals are reported to be 0.63, 1.16 and 0.69 Å which are reduced to 0.17, 0.00 and 0.02 Å when appropriate dispersive corrections are included. We see in Table 5-9 that the PM3-D method is capable of predicting the structures of dispersion-bonded complexes with greater accuracy than some uncorrected DFT functionals and with an accuracy comparable to that for the dispersion corrected PBE functional [130].

The results of the various semi-empirical calculations on the reference structures contained within the *JSCH-2005* database (134 complexes; 31 hydrogen-bonded base-pairs, 32 interstrand base pairs, 54 stacked base pairs and 17 amino acid base pairs) are summarised in Table 5-10. The deviations of the various interaction energies from the reference values are displayed in Figure 5-5. As with the *S22* training set, the AM1 and PM3 methods generally underestimate the interactions whereas the dispersion corrected method (PM3-D) mostly over-estimates the interactions a little. Overall the PM3-D results are particularly impressive given that the method has only

*Table 5-9.* Statistics of the deviation between the *S22* database and semi-empirical interaction energies (kcal mol$^{-1}$) and distances (Å) [a]

|  | AM1 | PM3 | PM3–D |
| --- | --- | --- | --- |
| *MUE*[b] | | | |
| Hydrogen-bonded complexes (7) | 8.78(9.69) | 7.28(7.27) | 1.77(0.61) |
| Complexes with predominant dispersion contribution (8) | 3.23(6.65) | 3.00(7.52) | 1.48(1.42) |
| Mixed complexes (7) | 2.73(3.28) | 2.14(2.80) | 1.28(0.59) |
| | | | |
| *Overall Statistics – Interaction Energies* | | | |
| MSE[c] | −4.83(−6.54) | −4.09(−5.94) | 0.59(0.21) |
| MD[d] | 3.27(4.54) | 2.74(4.03) | 1.34(0.85) |
| MUE[e] | 4.83(6.54) | 4.09(5.94) | 1.51(0.90) |
| RMSE[f] | 6.25(8.47) | 5.22(7.73) | 1.65(1.18) |
| MAXE–MINE[g] | 13.16(17.41) | 9.74(19.32) | 5.13(5.42) |
| | | | |
| *Overall Statistics – Interaction Distances* | | | |
| MSE[c] | −0.83 | −0.60 | 0.08 |
| MD[d] | 0.79 | 0.79 | 0.18 |
| MUE[e] | 0.85 | 0.69 | 0.20 |
| RMSE[f] | 1.28 | 1.17 | 0.25 |
| MAXE–MINE[g] | 3.47 | 3.89 | 1.01 |

[a](In parenthesis) statistics for the interaction energies for complexes in the fixed nuclear geometries provided in the *S22* database [1].
[b](In parenthesis) number of comparisons.
[c]Mean signed error.
[d]Mean deviation.
[e]Mean unsigned error.
[f]Root mean square error.
[g]Error spread (largest positive minus largest negative error).

been parameterised against the *S22* complexes. The MUE (for the 134 complexes) is 1.26 kcal mol$^{-1}$ (PM3-D) for the stabilisation energies calculated at the fixed nuclear geometries provided in the *JSCH-2005* database. Of the different sets of complexes, the interaction energies are least well predicted for the hydrogen-bonded complexes, whereas for the interstrand and stacked base-pairs the MUEs are notably less (Table 5-10) and are comparable to those reported for the BLYP-D method [0.29 (interstrand) and 0.53 kcal mol$^{-1}$ (stacked)] [128]. Most importantly, as far as biological molecules are concerned we note that optimisation of the stacked complexes using AM1 and PM3 methods leads to different complexes upon geometry optimisation (hydrogen-bonded if possible) in line with uncorrected DFT calculations [130]. For the PM3-D method the structures of the stacked complexes are generally quite close to the reference ones [1].
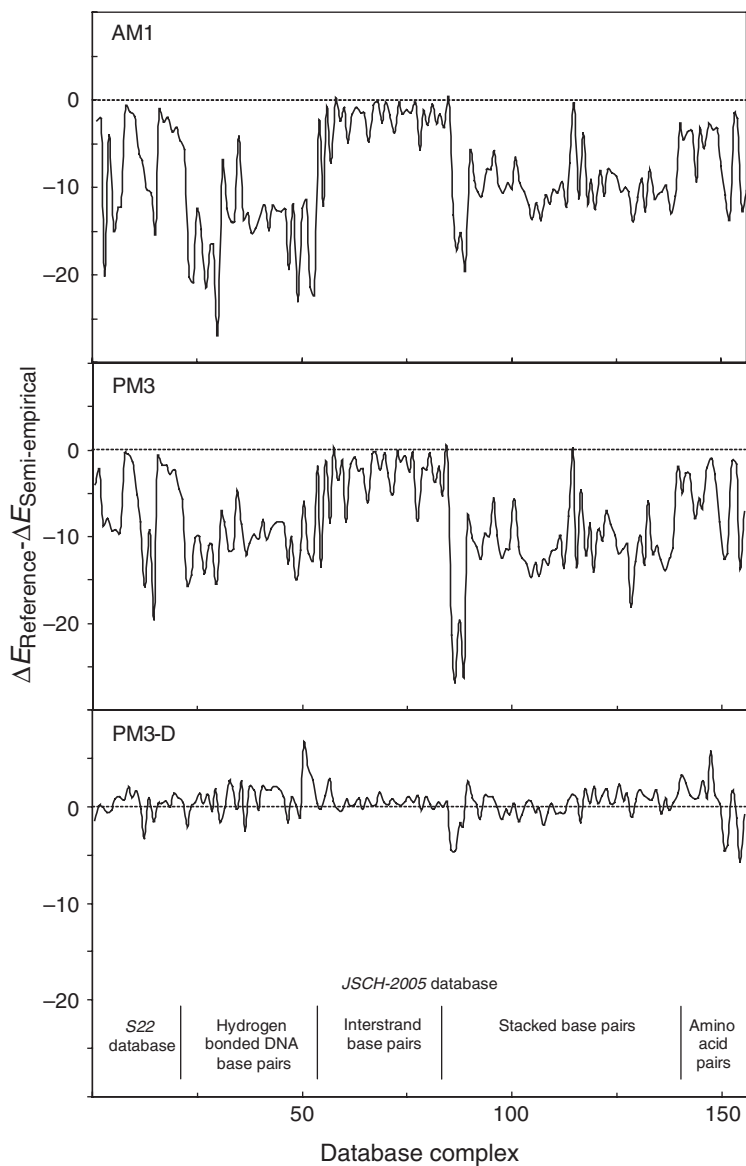
*Figure 5-5.* Deviation of semi-empirical interaction energies from *S22* and *JSCH-2005* database values
(kcal mol$^{-1}$) [1]. Reference [43] reproduced by permission of the Royal Society of Chemistry

*Table 5-10*. Statistics of the deviation between the *JSCH–2005* database and semi–empirical interaction energies (kcal mol$^{-1}$) for 134 biological complexes[a]

|  | AM1 | PM3 | PM3–D |
|---|---|---|---|
| *MUE* |  |  |  |
| Hydrogen–bonded DNA base pairs (31) | 15.17 | 10.63 | 1.66 |
| Interstrand base pairs (32) | 2.30 | 2.87 | 0.58 |
| Stacked base pairs (54) | 10.36 | 11.80 | 1.11 |
| Amino acid base pairs (17) | 5.93 | 5.60 | 2.28 |
| | | | |
| *Overall Statistics* | | | |
| MSE[b] | −8.98 | −8.60 | 0.41 |
| MD[c] | 4.74 | 4.29 | 1.17 |
| MUE[d] | 8.99 | 8.61 | 1.26 |
| RMSE[e] | 10.69 | 10.11 | 3.04 |
| MAXE–MINE[f] | 26.97 | 27.07 | 12.40 |

[a] Calculated at the *JSCH-2005* database geometries [1]. (In parenthesis) number of comparisons.
[b] Mean signed error.
[c] Mean deviation.
[d] Mean unsigned error.
[e] Root mean square error.
[f] Error spread (largest positive minus largest negative error).

We have also calculated thirteen neutral amino acid pairs containing phenylalanine taken from the hydrophobic core the protein Rd (Table 5-11). Overall, the MUE for the interaction energies at the PM3-D level is less than half of the values at AM1 and PM3 levels. Comparing our corrected semi-empirical method with MM calculations on the F30 cluster of Rd (PDB Code: 1RB9) we see that in general the MM calculations predict attractive interactions [131]. These findings are in contrast to B3LYP/6-31G* calculations on the same cluster in which all the pair interactions are predicted to be repulsive [131]. Overall for the five amino acid pairs considered, the interaction energy MUE at the PM3-D level (1.82 kcal mol$^{-1}$) is very close to those calculated using the CHARM22 and MMFF94 potentials (1.82 and 1.88 kcal mol$^{-1}$).

The work has shown that the addition of a dispersive correction to the normal PM3 energy, with appropriate optimisation of the parameters, can yield interaction energies between biologically important groups to within on average, 1–1.5 kcal mol$^{-1}$ of the results of high-level ab initio calculations [43]. The dispersion corrected method has been extensively tested on two databases involving a total of 156 biologically relevant molecules. The results are particularly pleasing since the semi-empirical parameters have been optimised using only a small training set which has been suggested to include all the typical biological interactions. Our work does however indicate that to increase the accuracy of current semi-empirical methods will in fact require further modifications of the functional form of the Hamiltonian in addition to the development of new parameter sets.

*Table 5-11.* Molecular mechanics and semi–empirical interaction energies (kcal mol$^{-1}$) of F30 cluster in Rd[a]

| Amino Acid Base Pair | F30 | | | | | Sum | MUE |
|---|---|---|---|---|---|---|---|
| | F49 | K46 | L33 | Y13 | Y4 | | |
| Ab initio reference[b] | −3.30 | −3.10 | −5.00 | −3.90 | −7.00 | −22.30 | |
| *Semi–empirical* | | | | | | | |
| AM1 | 0.30 | −0.34 | −0.47 | −0.20 | 2.38 | 1.67 | 4.79 |
| PM3 | −0.54 | −1.08 | 0.02 | −1.27 | 0.97 | −1.90 | 4.08 |
| PM3–D | −4.44 | −4.50 | −8.29 | −6.37 | −7.81 | −31.41 | 1.82 |
| *DFT–D*[c] | | | | | | | |
| BLYP–D/TZV(2d,2p) | −2.60 | −3.05 | −5.14 | −4.14 | −5.10 | −20.03 | 0.61 |
| *Molecular Mechanics*[d] | | | | | | | |
| AMBER parm94, parm99 | −19.30 | −13.00 | −6.60 | −3.10 | −5.60 | −47.60 | 5.94 |
| CHARM22 | −2.10 | −1.30 | −2.70 | −2.40 | −4.70 | −13.20 | 1.82 |
| MMFF94 | −1.50 | −4.10 | −2.20 | −2.10 | −5.00 | −14.90 | 1.88 |

[a]PDB Code: 1RB9.
[b]CCSD(T)/CBS values from *JSCH-2005* database [1].
[c]Morgado et al. [128].
[d]Vondrášek et al. [131].

## 5.8.    CONCLUSIONS

The most practical way of including a QM description within a model of protein structure and reactivity is an ever-present problem. Despite increases in computational power and advances in method development, the size and complexity of biological systems generally limits the use of highly-correlated QM methods to systems containing just a few tens of atoms. However, such calculations are of great value in providing benchmarks against which less computationally demanding methods can be judged and parameterised, as illustrated by the work of Hobza and co-workers [1]. Progress is currently being made in the application of DFT methods as part of hybrid QM/MM embedding schemes to quite large molecular clusters [132, 133] although such methods are generally too demanding computationally when many conformational states need to be sampled. One way to tackle this problem involves the evaluation of a free energy correction between a full MM and a QM/MM model [134].

Whatever scheme is chosen, a method is needed to rapidly evaluate the appropriate QM energy. If this could be done with confidence using a semi-empirical method the solution of many previously inaccessible problems could well become possible. This has provided the catalyst for the present renaissance in the development of such

methods. However, just how valuable these "QM force-fields" will turn out to be for solving important problems will only become clear as they are more widely used.

## ACKNOWLEDGMENTS

## REFERENCES

1. Jurečka P, Šponer J, Černy J, Hobza P (2006) Phys Chem Chem Phys 8:1985
2. Pu J, Gao J, Truhlar DG (2006) Chem Rev 106:3140
3. Holm RH, Kennepohl P, Solomon EI (1996) Chem Rev 96:2239
4. Bolhuis PG, Dellago C, Geissler PL, Chandler D (2000) J Phys Condens Matter 12:A147
5. Dimelow R, Bryce RA, Masters AJ, Hillier IH, Burton NA (2006) J Chem Phys 124:114113
6. Tresadern G, Wang H, Faulder PF, Burton NA, Hillier IH (2003) Mol Phys 101:2775
7. Alhambra C, Luz Sanchez M, Corchado J, Gao J, Truhlar DG (2001) Chem Phys Lett 347:512
8. Field MJ, Basch M, Karplus M (1990) J Comput Chem 11:700
9. Gao J, Xia X (1992) Science 258:631
10. Maseras F, Morokuma K (1995) J Comput Chem 16:1170
11. Vreven T, Morokuma K, Farkas O, Schlegel HB, Frisch MJ (2003) J Comput Chem 24:760
12. Senn MH, Thiel W (2007) Curr Opin Chem Biol 11:182
13. Werner HJ, Manby FR, Knowles PJ (2003) J Chem Phys 118:8149
14. Morgado C, McNamara JP, Hillier IH, Sundararajan M (2005) Mol Phys 103:905
15. Clark T (2000) J Mol Struct (Theochem) 530:1
16. Pople JA, Segal GA (1965) J Chem Phys 43:S136
17. Pople JA, Beveridge DL, Dobosh PA (1967) J Chem Phys 47:2026
18. Dewar MJS, Thiel W (1977) Theor Chim Acta 46:89
19. Dewar MJS, Thiel W (1977) J Am Chem Soc 99:4899
20. Dewar MJS, Zoebisch EG, Healy EF, Stewart JJP (1993) J Am Chem Soc 115:5348
21. Stewart JJP (1989) J Comput Chem 10:209
22. Stewart JJP (2004) J Mol Model 10:155
23. Thiel W, Voityuk AA (1992) Theor Chim Acta 81:391
24. Thiel W, Voityuk AA (1996) Theor Chim Acta 93:315
25. Voityuk AA, Rösch N (2000) J Phys Chem A 104:4089
26. McNamara JP, Sundararajan M, Hillier IH, Ge J, Campbell A, Morgado C (2006) J Comput Chem 27:1307
27. Rossi I, Truhlar DG (1995) Chem Phys Lett 233:231
28. Brothers EN, Merz Jr KM (2002) J Phys Chem B 106:2779
29. Hutter MC, Reimers JR, Hush NS (1998) J Phys Chem B 102:8080
30. Cundari TR, Deng J, Fu W (2000) Int J Quantum Chem 77:421
31. Sundararajan M, McNamara JP, Hillier IH, Wang H, Burton NA (2005) Chem Phys Lett 404:9
32. McNamara JP, Sundararajan M, Hillier IH (2005) J Mol Graph Model 24:128

33. McNamara JP, Berrigan SD, Hillier IH (2007) J Chem Theory Comput 3:1014
34. Giese TJ, Sherer EC, Cramer CJ, York DM (2005) J Chem Theory Comput 1:1275
35. Nam K, Cui Q, Gao J, York DM (2007) J Chem Theory Comput 3:486
36. Menegon Arantes G, Loos M (2006) Phys Chem Chem Phys 8:347
37. Csonka GI, Ángyán JG (1997) J Mol Struct (Theochem) 393:31
38. Repasky MP, Chandrasekhar J, Jorgensen WL (2002) J Comput Chem 23:1601
39. Winget P, Horn AHC, Selçuki C, Martin B, Clark T (2003) J Mol Model 9:408
40. Winget P, Clark T (2005) J Mol Model 11:439
41. Equation 9 differs from the one given in reference [25]; in the original publication, the factor of "2" was inadvertently omitted.
42. Dewar MJS, Jie C, Zoebisch EG (1988) Organometallics 7:513
43. McNamara JP, Hillier IH (2007) Phys Chem Chem Phys 9:2362
44. Tejero I, González-Lafont Á, Lluch JM (2007) J Comput Chem 28:997
45. Jackson MD, Denu JM (2001) Chem Rev 101:2313
46. Westheimer FH (1987) Science 235:1173
47. Xu D, Guo H, Liu Y, York DM (2005) J Phys Chem B 109:13827
48. Chen X, Zhan CG (2004) J Phys Chem A 108:6407
49. Menegon Arantes G, Chaimovich H (2005) J Phys Chem A 109:5625
50. Menegon G, Loos M, Chaimovich H (2002) J Phys Chem A 106:9078
51. Liu Y, Gregersen BA, Hengge A, York DM (2006) Biochemistry 45:10043
52. Liu Y, Gregersen BA, Lopez X, York DM (2005) J Phys Chem B 109:19987
53. López CS, Faza ON, de Lera AR, York DM (2005) Chem Eur J 11:2081
54. Liu Y, Lopez X, York DM (2005) Chem Commun 31:3909
55. Lopez X, Dejaegere A, Leclerc F, York DM, Karplus M (2006) J Phys Chem B 110:11525
56. Hart JC, Burton NA, Hillier IH, Harrison MJ, Jewsbury P (1997) Chem Commun 15:1431
57. Hart JC, Hillier IH, Burton NA, Sheppard DW (1998) J Am Chem Soc 120:13535
58. Hart JC, Sheppard DW, Hillier IH, Burton NA (1999) Chem Commun 1:79
59. Valiev M, Kawai R, Adams JA, Weare JH (2003) J Am Chem Soc 125:9926
60. Truhlar DG, Garrett BC, Klippenstein SJ (1996) J Phys Chem 100:12771
61. Zhang ZY, Wu L, Chen L (1995) Biochemistry 34:16088
62. Thiel W, Voityuk AA (1996) J Phys Chem 100:616
63. Weis WI (1997) Curr Opin Struct Biol 7:624
64. Norskov-Lauritsen L, Allinger NL (1984) J Comput Chem 5:326
65. Kony D, Damm W, Stoll S, van Gunsteren WF (2002) J Comput Chem 23:1416
66. McNamara JP, Muslim AM, Abdel-Aal H, Wang H, Mohr M, Hillier IH, Bryce RA (2004) Chem Phys Lett 394:429
67. Stern HA, Kaminski GA, Banks JL, Zhou R, Berne BJ, Friesner RA (1999) J Phys Chem B 103:4730
68. Woods RJ, Szarek WA, Smith VH (1991) Chem Commun 5:334
69. Barrows SE, Dulles FJ, Cramer CJ, French AD, Truhlar DG (1995) Carbohydr Res 276:219
70. Apprell M, Strati G, Willett JL, Momany FA (2004) Carbohydr Res 339:537
71. Reiling S, Schlenkrich M, Brickmann J (1996) J Comput Chem 17:450
72. Case DA, Darden TA, Cheatham III TE, Simmerling CL, Wang J, Duke RE, Luo R, Merz KM, Pearlman DA, Crowley M, Walker RC, Zhang W, Wang B, Hayik S, Roitberg A, Seabra G, Wong KF, Paesani F, Wu X, Brozell S, Tsui V, Gohlke H, Yang L, Tan C, Mongan J, Hornak V, Cui G, Beroza P, Mathews DH, Schafmeister C, Ross WC, Kollman PA (2006) AMBER 9, University of California, San Francisco
73. Damm W, Frontera A, Tirado-Rives J, Jorgensen WL (1997) J Comput Chem 18:1955
74. Muslim AM, Bryce RA (2004) Chem Phys Lett 388:473

75. Bell RP (1980) The tunnel effect in chemistry. Chapman and Hall, London
76. Billeter SR, Webb SP, Agarwal PK, Iordanov T, Hammes-Schiffer S (2001) J Am Chem Soc 123:11262
77. Kohen A, Klinman JP (1998) Acc Chem Res 31:397
78. Basran J, Sutcliffe MJ, Scrutton NS (1999) Biochemistry 38:3218
79. Basran J, Patel S, Sutcliffe MJ, Scrutton NS (2001) J Biol Chem 276:6234
80. Northrop DB, Cho YK (2000) Biochemistry 39:2406
81. Tresadern G, McNamara JP, Mohr M, Wang H, Burton NA, Hillier IH (2002) Chem Phys Lett 358:489
82. Truhlar DG, Garrett BC (1980) Acc Chem Res 13:440
83. Truhlar DG, Isaacson AD, Garrett BC (1985) Generalised transition state theory. In: Baer M (ed) Theory of Chemical Reaction Dynamics, Vol 4. CRC, Boca Raton, p 65
84. Liu YP, Lynch GC, Truong TN, Lu DH, Truhlar DG, Garrett BC (1993) J Am Chem Soc 115:2408
85. Allison TC, Truhlar DG (1998) Testing the accuracy of practical semiclassical methods: variational transition state theory with optimized multidimensional tunnelling. In: Thompson DL (ed) Modern Methods for Multidimensional Dynamics Computations in Chemistry. World Scientific, Singapore, p 618
86. Fernandez-Ramos A, Ellingson BA, Garrett BC, Truhlar DG (2007) Variational transition state theory with multidimensional tunneling. In: Lipkowitz KB, Cundari TR, Boyd DB (eds) Reviews in Computational Chemistry, Vol 23. Wiley-VCH, New York, p 125
87. Skodje RT, Truhlar DG, Garrett BC (1981) J Phys Chem 85:3019
88. Harrison MJ, Burton NA, Hillier IH (1997) J Am Chem Soc 119:12285
89. Warshell A, Levitt M (1976) J Mol Biol 103:227
90. Singh UC, Kollman PA (1986) J Comp Chem 7:718
91. Tresadern G, Faulder PF, Gleeson P, Tai Z, MacKenzie G, Burton NA, Hillier IH (2003) Theor Chem Acc 109:108
92. Shearer GL, Kim K, Lee KM, Wang CK, Plapp BV (1993) Biochemistry 32:11186
93. Bahnson BJ, Klinman JP (1995) Methods Enzymol 249:373
94. Jonsson T, Glickman MH, Sun S, Klinman JP (1996) J Am Chem Soc 118:10319
95. Glickman MH, Klinman JP (1996) Biochemistry 35:12882
96. Minor W, Steczko J, Stec B, Otwinowski Z, Bolin JT, Walter R, Axelrod B (1996) Biochemistry 35:10687
97. Finnen DC, Pinkerton AA, Dunham WR, Sands RH, Funk MO (1991) Inorg Chem 30:3960
98. HYPERCHEM (TM), Hypercube Inc., Gainesville, Florida, USA
99. Stewart JJP (2001) MOPAC 2002, Fujitsu Limited, Tokyo, Japan
100. Lovenberg W (1977) Iron-sulfur proteins. Academic, New York
101. Niu S, Wang XB, Nicolas JA, Wang LS, Ichiye T (2003) J Phys Chem A 107:2898
102. Wang XB, Wang LS (2000) J Chem Phys 112:6959
103. Yang X, Wang XB, Fu YJ, Wang LS (2003) J Phys Chem 107:1703
104. Mohr M, McNamara JP, Wang H, Rajeev SA, Ge J, Morgado CA, Hillier IH (2003) Faraday Discuss 124:413
105. Lane RW, Ibers JA, Frankel RB, Papefthymiou GC, Holm RH (1977) J Am Chem Soc 91:84
106. Kennopohl P, Solomon EI (2003) Inorg Chem 42:689
107. Case DA, Pearlman DA, Caldwell JW, Cheatham III TE, Wang J, Ross WS, Simmerling CL, Darden TA, Merz KM, Stanton RV, Cheng AL, Vincent JJ, Crowley M, Tsui V, Gohlke H, Radmer RJ, Duan Y, Pitera J, Massova I, Seibel GL, Singh UC, Winer PK, Kollman PA (2002) AMBER 7, University of California, San Francisco

108. Xiao Z, Lavery MJ, Ayhan M, Scrofani SDB, Wilce MCJ, Guss JM, Tregloan PA, George GN, Wedd AG (1998) J Am Chem Soc 120:4135

109. Dauter Z, Wilson KS, Sieker LC, Moulis JM, Meyer J (1996) Proc Natl Acad Sci USA 93:8836

110. Sigfridsson E, Olsson MHH, Ryde U (2001) Inorg Chem 40:2509

111. Tard C, Liu X, Ibrahim SK, Bruschi M, De Gioia L, Davies SC, Yang X, Wang LS, Sawers G, Pickett CJ (2005) Nature 433:610

112. Torres RA, Lovell T, Noodleman L, Case DA (2003) J Am Chem Soc 125:1923

113. Meunier B, de Visser SP, Shaik S (2004) Chem Rev 104:3947

114. Groves JT, Hang YZ (1995) Models and mechanics of cytochrome p-450 action. In: Ortiz de Montellano PR (ed) Cytochrome p-450: structure, mechanism and biochemistry, 2nd edn Plenum, New York, p 3

115. Shaik S, de Visser SP, Oligiaro F, Schwarz H, Schröder D (2002) Curr Opin Chem Biol 6:556

116. de Visser SP, Ogliaro F, Sharma PK, Shaik S (2002) J Am Chem Soc 124:11809

117. Manchester JI, Dinnocenzo JP, Higgins LA, Jones JP (1997) J Am Chem Soc 119:5069

118. Grimme S (2004) J Comput Chem 25:1463

119. Zhao Y, Truhlar DG (2007) J Chem Theory Comput 3:289

120. Zhao Y, Truhlar DG (2005) J Phys Chem A 109:5656

121. Zhao Y, Truhlar DG (2005) Phys Chem Chem Phys 7:2701

122. Grimme S (2006) J Comput Chem 27:1787

123. Valdés H, Řeha D, Hobza P (2006) J Phys Chem B 110:6385

124. Dobés P, Otyepka M, Strnad M, Hobza P (2006) Chem Eur J 12:4297

125. Adamo C, Barone V (1998) J Chem Phys 108:664

126. Perdew J, Burke K, Wang Y (1996) Phys Rev B: Condens Matter Mater Phys 54:16533

127. Easton RE, Giesen DJ, Welch A, Cramer CJ, Truhlar DG (1996) Theor Chem Acc 93:281

128. Morgado C, Vincent MA, Hillier IH, Shan X (2007) Phys Chem Chem Phys 9:448

129. Anthony J, Grimme S (2006) Phys Chem Chem Phys 8:5287

130. Jurečka P, Černy J, Hobza P, Salahub DR (2007) J Comput Chem 28:555

131. Vondrášek J, Bendová L, Klusák V, Hobza P (2005) J Am Chem Soc 127:2615

132. Shaik S, Kumar D, de Visser SP, Altun A, Thiel W (2005) Chem Rev 105:2279

133. Sundararajan M, Surendran R, Hillier IH (2005) Chem Phys Lett 418:92

134. Rosta E, Klähn M, Warshel A (2006) J Phys Chem B 110:2934

CHAPTER 6

# DESIGN OF NEXT GENERATION FORCE FIELDS FROM AB INITIO COMPUTATIONS: BEYOND POINT CHARGES ELECTROSTATICS

G.A. CISNEROS[1], T.A. DARDEN[1], N. GRESH[2], J. PILMÉ[3,4,5], P. REINHARDT[4,5], O. PARISEL[4,5], AND J.-P. PIQUEMAL[4,5]

[1] *Laboratory of Structural Biology, National Institute of Environmental Health Sciences, MD F0-08, 111 TW. Alexander Dr., Research Triangle Park, NC 27709, USA*
[2] *Laboratoire de Pharmacochimie Moléculaire et Cellulaire, U648 INSERM, UFR Biomédicale, Université René-Descartes, 45, rue des Saints-Pères, 75006 Paris, France*
[3] *Université de Lyon, Université Lyon 1, Faculté de pharmacie, F-69373 Lyon, Cedex 08, France*
[4] *UPMC Univ Paris 06, UMR 7616, Laboratoire de Chimie Théorique, case courrier 137, 4 place Jussieu, F-75005, Paris, France, e-mail: jpp@lct.jussieu.fr*
[5] *CNRS, UMR 7616, Laboratoire de Chimie Théorique, case courrier 137, 4 place Jussieu, F-75005 Paris, France*
Correspondance to J-P Piquemal, jpp@lct.jussieu.fr

**Abstract:**     We present an overview of the energy functions used in two Anisotropic Polarizable Molecular Mechanics (APMM) procedures namely SIBFA (Sum of Interactions Between Fragments Ab initio computed) and GEM (Gaussian Electrostatic Model). As SIBFA is a second generation APMM scheme based on distributed multipoles, GEM is the first third generation APMM as it uses distributed hermite densities obtained from density fitting. The two approaches are formulated and calibrated on the basis of quantum chemistry. They embody nonclassical effects such as electrostatic penetration, exchange-polarization, and charge transfer. We address here the technical issues of anisotropy, nonadditivity, transferability and computational speedup of methods. In addition, we review the several ab initio intermolecular energy decomposition techniques that can be used to refine polarisable force fields. As we summarize their differences and similarities, we present our own scheme based on Fragment Localized Kohn-Sham orbitals through a Singles-Configuration Interaction (CI) procedure. We also present a chemically intuitive method based on the Electron Localization Function (ELF) which allows to unravel the local electrostatic properties beyond atomic centers: i.e., on bonds, lone pairs and $\pi$ system, an useful asset to understand bonding in molecules in order to build models

## 6.1.    INTRODUCTION

Nowadays, modern molecular modelling techniques propose numerous potential applications, from material sciences to protein structure prediction and drug design. Indeed, classical Molecular dynamics (MD) is now able to provide useful information to experimentalist as simulations are getting closer and closer to relevant biological timescales. Nevertheless, if MD is now able to produce microsecond trajectories, one should ask about the possible improvements of such simulations. At this point, two directions can be taken. The first consists in increasing the speed of MD softwares by coupling improved sampling methodology to massively parallel computers, having as goal to reach the second timescale. However, if this strategy will probably offer some interesting insights about biophysical process, there is no doubt that the question of the accuracy of the used empirical energy functions, the so-called force fields, should be raised. Indeed, current simulations are mainly aimed to compute free energies and despite success, actual data are already sufficient to demonstrate that current molecular mechanics (MM) potentials have serious shortcomings [1]. This can be easily understood when considering that free energy required an accurate evaluation of both enthalpic and entropic contributions. If entropy can be recovered through sampling efforts, enthalpy needs to be approximate from their quantum mechanical expression. In a way, classical MD can be simply seen as an approximate quantum Born-Oppenheimer MD approach treating the atomic nuclei as classical particles subject to interatomic forces. Presently, these latter remain obtained from empirical potentials far to reproduce first principles results. Therefore, MD should not be able to quantitatively describe vast numbers of systems dominated by difficult weak interactions such as H-bonds networks, metalloproteins and metal clusters, highly charges systems etc., where Chemistry and electron correlation/relativity dominate. For these systems, the right tools are required. In this context, Anisotropic Polarizable Molecular Mechanics (APMM) procedures have been developed (see Reference [2] and references therein). These approaches share the common characteristic of including a more evolved representation of the electrostatic contribution to the interaction energy compared to the usual point charge approximation, allowing a close reproduction of the anisotropic features of the ab initio Coulomb contribution. As we will see, some of them use distributed multipoles (sometime damped in order to include short-range penetration effects) or electronic Hermite densities for the latest generation. The philosophy of such approaches relies on an extensive use of quantum mechanics defining the so-called: "bottom-up strategy" [2]. First, the electrostatic moments or hermites densities are directly obtained from an ab initio calculation of the considered gas phase isolated molecule and stored in a library. Second, all intermolecular components of the force field should faithfully reproduce their ab initio counterpart as obtained from energy-decomposition procedures at the

Hartree-Fock, DFT or CCSD levels. Because they can reproduce such quantities, APMM procedures should account for an accurate description of the interactions including polarization cooperative effects and charge transfer. They should also enable the reproduction of local electrostatic properties such as dipole moments an also facilitate hybrid Quantum Mechanical/Molecular Mechanical (QM/MM) embeddings.

In this contribution, we will review some aspects of this strategy. First, we will explore some recent ab initio techniques that can be used for the refinement of APMM approaches. Among them, we will discuss the differences between energy decomposition approaches namely: Kitaura-Morokuma [3], Constrained Space Orbital Variations (CSOV) [4–6], Reduced Variational Space (RVS) [7], Ziegler-Baerends [8, 9] and Symmetry Adapted Perturbation Theory [10] procedures. Moreover, we will focus on a newly developed energy decomposition approach using Fragment-localized Kohn-Sham orbitals through a Singles-Configuration Interaction (CI) procedure [11]; and on a general approach to unravel local electrostatic properties, the so-called DEMEP [12] (Distributed Electrostatic Moments based on the Electron localization function Partition). In a second part, we will detail two APMM approaches in development in our labs. The first, called SIBFA (Sum of Interactions Between Fragments Ab initio computed) [2, 13] is a second APMM generation based on distributed multipoles. The second, named GEM (Gaussian Electrostatic Model) [2, 14–16] is the first APPM of the third generation based on electron density. Focusing on methodology, we will put in perspective the physical basis underlying the development of such MM energy functions and the possibility for a computation speedup, a key step to perform simulations.

## 6.2.    AB INITIO TECHNIQUES: FROM INTERMOLECULAR INTERACTIONS TO LOCAL ELECTROSTATIC PROPERTIES

Intermolecular Energy decomposition analyses (EDA) are very useful approaches to calibrate force fields. Indeed, an evaluation of the different physical components of the interaction energy, especially of the many-body induction, is a key issue for the development of polarisable models.

### 6.2.1.    Intermolecular Energy Decomposition Schemes: Equivalence Between Terms

However, due to the availability of numerous techniques, it is important to point out here the differences and equivalence between schemes. To summarize, two EDA families can be applied to force field parametrization. The first EDA type of approach is labelled SAPT (Symmetry Adapted Perturbation Theory). It uses non orthogonal orbitals and "recomputes" the total interaction upon perturbation theory. As computations can be performed up to the Coupled-Cluster Singles Doubles (CCSD) level, SAPT can be seen as a reference method. However, due to the cost of the use of non-orthogonal molecular orbitals, pure SAPT approaches remain limited

to small systems, even if Kohn-Sham orbitals based or local SAPT approaches tend to overcome such difficulties. The second family of methods is variational and based on the supermolecule approach ($\Delta E = E_{AB} - E_A - E_B$) following the early Kitaura-Morokuma (KM) and Ziegler schemes. It includes also the Constrained Space Orbital Variations (CSOV) (and the Reduced Variational Space (RVS), essentially similar to CSOV) approach. These methods are limited to the HF or DFT levels. Following a perturbation terminology, all EDA schemes can be partitioned between first, second and higher order terms:

$$\Delta E = E_1 - E_2 - \delta E_{higher-orders} + BSSE \qquad (6\text{-}1)$$

All schemes furnish two first-order terms. The first is the electrostatic interaction of the frozen monomers denoted $E_{es}$ in the variational approaches and which is strictly equivalent to the $E_{pol}^{10}$ in SAPT. The second is an exchange-repulsion term $E_{exch\text{-}rep}$ (denoted $E_{exch}^{10}$ in SAPT).The sum of them is sometime called frozen-core contribution ($E_{FC}$) like in the variational CSOV scheme of Bagus et al. This $E_{FC}$ is also equivalent to the Heitler-London energy, employing the unperturbed monomers orbitals. At the HF level, despite a different use of operators (**V** for SAPT, vs. **H** for the variational methods), these terms should be equivalent for all approaches if a reasonable basis set is used. Second-Order terms are more problematic and can be divided into a so-called induction term and a dispersion component, each one of these terms being associated to a repulsive second-order exchange term. At the HF level, the SAPT induction term ($E_{pol}^{20}$) should be equivalent to the BSSE corrected Orbital relaxation term of the Ziegler scheme also called Orbital Interaction. This latter Orbital interaction term corresponds itself to the sum of polarization ($E_{pol}$), charge transfer ($E_{CT}$) and BSSE term in the CSOV or RVS approaches. The subtle question of the evaluation of the sole ab initio polarization energy (without charge transfer), so important for the evaluation of the accuracy of polarizable models, is important as its evaluation requires to conserve the antisymmetry of the wavefunction through relaxation of the monomers. Such computation remains limited to the CSOV and RVS scheme as the Morokuma scheme violates the antisymmetry leading to an overestimation of the polarization (and of the charge transfer term) (see References [14, 17–19] and reference therein). That way the CSOV and RVS $E_{pol}$ (and $E_{CT}$) term embodies the $E_{pol\text{-}exch}$ term through conserved MOs orthogonality. It is also important to note that higher order coupling are not included in CSOV and RVS, that way, such polarization term can be seen as a lower bound for the evaluation of polarization. Strategies to use these schemes have been previously reported.

SAPT methods remain the only approaches allowing the evaluation of dispersion (Figure 6-1).

At this point, it is important to notice that in general, the sum of the contributions do not match exactly $\Delta E$ as higher order terms are present. The difference between the sum of contributions and $\Delta E$ is denoted $\delta E$. Concerning the variational schemes, $\delta E$ is generally small in the CSOV (or RVS) approach thanks to the antisymmetry conservation and not present in the Ziegler scheme as the $E_{OI}$ term is taking into account a fully relaxed wavefunction. It is not the case for the KM scheme which

$$\Delta E_{Ziegler-Baerends}^{HF,DFT} = E_{es} + E_{exch-rep} + E_{OI} = E_1 + E_2 + BSSE$$

$$\Delta E_{CSOV}^{HF,DFT,MCSCF} = E_{FC} + E_{pol/CSOV\ (A)} + E_{pol/CSOV\ (B)} + E_{CT/CSOV_{A\to B}} + E_{CT/CSOV_{B\to A}} + BSSE + \delta E$$

$$with: E_{FC} = E_{es} + E_{exch-rep}; E_{pol/CSOV} = E_{pol} + E_{Pol-exch}$$

$$\Delta E_{SAPT}^{HF} = E_{es}^{10} + E_{exch-rep}^{10} + E_{ind,resp}^{20} + E_{exch-ind,resp}^{20} + \delta E_{resp}^{HF}$$

$$\Delta E_{SAPT}^{CCSD} = \Delta E_{SAPT}^{HF} + correlation-corrections$$

*Figure 6-1.* Notations for usual energy decomposition schemes

can embody very large $\delta E$ (sometime denoted $E_{mix}$) in presence of charge species. For other reasons (especially due to some convergence difficulties of the perturbation series for induction, see Ref. [20–22] for details) the same problem can occur for SAPT. Table 6-1 summarizes these informations.

### 6.2.2. Beyond Two-Body Interaction: Fragment-Localized Kohn-Sham Orbitals via a Singles-CI Procedure

As discussed below, EDA schemes are generally limited to dimer interactions (up to small trimer for SAPT). If the RVS scheme allows an evaluation of contribution for more than two molecules at the HF level, EDA methods allowing the inclusion of electron correlation did not exist up to a very recent time (see Head-Gordon's Scheme [23]) for the computation of large assemblies of molecules. We present here the methodology at the basis of a new potentially linear scaling local approach based on Fragment-localized Kohn-Sham orbitals via a Singles-CI procedure [11].

#### 6.2.2.1. Method: Fragment-Localized Kohn-Sham Orbitals

In the literature we may find the procedure for creating localized Hartree-Fock orbitals via an energy minimization based on a CI procedure employing mono-excitations (see for instance Reference [24]). The scheme starts from a set of given (guess) orbitals and solves iteratively the Hartree-Fock equations via the steps:

1. Symmetric (Löwdin) orthogonalisation of the orbitals via $\mathbf{S}^{-1/2}$
2. Construction of the Fock matrix
3. Calculation of the total energy
4. Construction and diagonalisation of an approximate Singles-CI matrix
5. Use in first order of the CI coefficients to correct the occupied and virtual molecular orbitals
6. Return to step 1

In step 3, a criterion of convergence may be introduced to terminate the iterations. Two other points should be mentioned: instead of taking the correct Singles-CI matrix, we may resort to a simpler one, omitting single bi-electronic integrals and using only Fock-matrix elements as:

$$\left\langle \Phi_i^a \left| \mathrm{H} \right| \Phi_i^b \right\rangle \approx F_{ab}\delta_{ij} - F_{ij}\delta_{ab} \tag{6-2}$$

*Table 6-1.* Contribution to the total interaction energy from different energy decomposition schemes

| ΔE Contributions | First order (or Frozen Core) | | | Second Order | | | | Higher orders (δE) |
|---|---|---|---|---|---|---|---|---|
| Methods | Electrostatics | Exch.-rep. | Correlation corrections | Induction | Exchange-induction | Ind. and exch-Ind. Correlation corrections | Dispersion | |
| SAPT | Yes | Yes | Yes (up to CCSD) | Yes | Yes | Yes Yes (up to CCSD) | Yes (+exchange dispersion; up to CCSD) | Yes (up to third order) |
| Kitaura-Morokuma | Yes | Yes | No | Yes ($E_{ind}=E_{pol}+E_{ct}$) | No | No | No | Yes (by difference from ΔE) |
| CSOV | Yes | Yes | Through DFT or MCSCF | Yes ($E_{ind}=E_{pol}+E_{ct}$) | Yes (included in $E_{pol}$ and $E_{ct}$) | Through DFT or MCSCF | No | Yes (by difference from ΔE) |
| Ziegler/Baerends | Yes | Yes | Through DFT | Yes ($E_{ind}=E_{OI}$) | Yes (included in $E_{OI}$) | Through DFT | No | No (included in $E_{OI}$) |

From the obtained wavefunction:

$$\Psi = \Phi_0 + \sum_i^a c_i^a \Phi_i^a \tag{6-3}$$

we use the coefficients for correcting the orbitals as:

$$\varphi_i' = \varphi_i + \sum_a c_i^a \varphi_a \text{ (occupied orbitals)} \tag{6-4}$$

$$\varphi_a' = \varphi_a - \sum_i c_i^a \varphi_i \text{ (virtual orbitals)} \tag{6-5}$$

Including the correction for the virtual orbitals ensures the orthogonality between occupied and virtual orbitals. Nevertheless, within the two separate orbital spaces, the orbitals must be re-orthogonalized in each iteration.

The advantage of the scheme lies in possibility to cut indices with a distant dependent selection criterion, rendering the method potentially linear scaling. As a consequence, orbitals for periodic structures may be created in this way (see References [25, 26]).

We may ask now, whether the same procedure may be applied to density-functional theory, just by replacing the Fock operator by the corresponding Kohn-Sham operator. To this end we have to look at the minimization of the total energy with respect to the density of a multi-determinantal wavefunction $\Psi$. We write the density as:

$$\Psi = c_0 \Phi_0 + \sum_i c_I \Phi_I$$

$$\rho(\vec{r}) = N \int \ldots \int d^3 r_1 \ldots d^3 r_{N-1} \left| \Psi^2(\vec{r}_1, \ldots, \vec{r}_{N-1}, \vec{r}) \right| \tag{6-6}$$

$$= c_0^2 \rho_{\Phi_0}(\vec{r}) + \sum_i c_I^2 \rho_{\Phi_I}(\vec{r}) + 2 \sum_{I<J} c_I c_J \varphi_k^I(\vec{r}) \varphi_l^J(\vec{r})$$

Following Reference [27], we may write the variation of the exchange-correlation energy as:

$$\int \upsilon^{XC}(\vec{r}) \rho_I(\vec{r}) d^3 r = \langle \Phi_I | V^{XC} | \Phi_I \rangle$$

$$\int \upsilon^{XC}(\vec{r}) \varphi_k^I(\vec{r}) \varphi_j^I(\vec{r}) d^3 r = \langle \Phi_I | V^{XC} | \Phi_J \rangle \tag{6-7}$$

$$\frac{\partial \rho(\vec{r})}{\partial c_I} = 2 c_I \rho_{\Phi_I}(\vec{r}) + 2 \sum_{J \neq I} c_J \varphi_k^I(\vec{r}) \varphi_l^J(\vec{r}) \tag{6-8}$$

$$\frac{\delta E^{XC}[\rho]}{\delta c_I} = 2c_I \langle \Phi_I| V^{XC} |\Phi_I\rangle + 2 \sum_{I \neq J} c_J \langle \Phi_I| V^{XC} |\Phi_J\rangle \qquad (6\text{-}9)$$

As the same construction holds for the Coulomb energy and the mono-electronic part, we obtain equations completely analogous to the system of linear equations for the Singles-CI:

$$E = \langle \Phi_0| K |\Phi_0\rangle + \sum_I c_I \langle \Phi_0| K |\Phi_I\rangle$$

$$c_I E = \langle \Phi_0| K |\Phi_I\rangle + c_I \langle \Phi_I| K |\Phi_I\rangle + \sum_{J \neq 0, I} c_J \langle \Phi_I| K |\Phi_J\rangle \qquad (6\text{-}10)$$

which we have to solve for the coefficients $c_I$ at each SCF iteration. Indeed, the only difference to Hartree-Fock theory lies in the use of the Kohn-Sham operator $K = T + Z + J + V^{XC}$ instead of the usual Hamiltonian $H = T + Z + 1/r_{12}$, reduced in the CI matrix to Fock-matrix elements.

### 6.2.2.2.    *Usefulness: From Energy Decomposition to Local Properties*

Apart from the question of linear scaling methods, we may employ the so-constructed orbitals for studying weakly interacting complexes. Of course, usual functionals do not include the important dispersion terms, but such an approach remains effective to study induction in large assemblies of molecules and, as we will see, for extracting monomer properties and interaction-induced changes of these.

(a) **Application to energy decomposition:** We first tested the accuracy our approach by implementing a Ziegler-Baerends type scheme by separately computing the electrostatic, exchange-repulsion and Orbital Interaction components of the interaction energy.

As expected we observed the invariance to orbitals localization of our decompositions scheme on a previously investigated linear water dimer configuration: localized and canonical orbitals lead to rigorously the same energy contributions, which is not the case for all decomposition schemes (due to projections or approximative orbital rotations). Concerning force field parametrization, it is interesting to observe the in influence of the addition of the second monomer basis functions. It clearly affects the energy components by diminishing the value of electrostatic and increasing the value of exchange-repulsion energy values. This "BSSE-like" (BSSE stands for Basis Set Superposition Error) effect is then clearly pronounced for Frozen Core (or first order in the SAPT terminology) as BSSE clearly acts on the other components. Table 6-2 displays such effect on the canonical water dimer. The decompositions have been easily extended beyond dimer systems [11], allowing the calculation of many-body contributions in contrast to SAPT or CSOV calculations, often restricted to the implementation of 2-body terms [14]. The scheme is also interesting to compute local fragment properties such as dipoles moments. As an example we may look at the dipole moment of two interacting $NH_3$ molecules. For each molecule we may

*Table 6-2.* Effect of dimer basis set on the components of the interaction energy

| Method | $E_S$ | $E_{exch-rep}$ | $E_{FC}$ | $E_{OI}$ + BSSE | $\Delta E$ |
|---|---|---|---|---|---|
| Monomers in the respective monomer basis | | | | | |
| HF(6d) | −8.27 | 6.91 | −1.36 | −2.18 | −3.55 |
| B3LYP(6d) | −8.03 | 7.48 | −0.55 | −3.30 | −3.86 |
| Monomers in the dimer basis | | | | | |
| HF(6d) | −8.30 | 7.02 | −1.28 | −2.27 | −3.55 |
| B3LYP(6d) | −8.18 | 6.87 | −1.31 | −3.04 | −4.35 |

*Table 6-3.* Comparison of the dipoles of the isolated individual monomers (dipole M) compare to the dipole moments of molecules within the dimer (dipole D) via the interaction, calculated with different functionals. Units are atomic units, and we give as well the difference in length and orientation

| NH3-NH3 | Dipole(M) | Dipole(D) | Difference | angle |
|---|---|---|---|---|
| HF | 0.613 | 0.664 | 0.051 | 8.8 |
| | 0.613 | 0.690 | 0.077 | 0.1 |
| BLYP | 0.552 | 0.635 | 0.0803 | 15.5 |
| | 0.535 | 0.677 | 0.142 | 1.2 |

calculate a dipole moment separately, and look for the deformation of the monomer orbitals (M) when constructing the dimer orbitals (D) via the described Singles-CI procedure. We have a good trace for the deformation, as the iterations only deform in a minimal sense the starting guess orbitals (Table 6-3).

## 6.2.3. Distributed Electrostatic Moments Based on the Electron Localization Function Partition

(**a**) **Theory:** In addition, as a fine understanding of cooperative effects is required in order to test the validity and the transferability of force fields parameters, some of us have been developing methodologies enabling the evaluation of local chemically intuitive distributed electrostatic moments using the topological analysis of the Electron Localization Function (ELF) [12].

For over a decade, the topological analysis of the ELF has been extensively used for the analysis of chemical bonding and chemical reactivity. Indeed, the Lewis' pair concept can be interpreted using the Pauli Exclusion Principle which introduces an effective repulsion between same spin electrons in the wavefunction. Consequently, bonds and lone pairs correspond to area of space where the electron density generated by valence electrons is associated to a weak Pauli repulsion. Such a property was noticed by Becke and Edgecombe [28] who proposed an expression of ELF based on the laplacian of conditional probability of finding one electron of spin σ at $r_2$, knowing that another reference same spin electron is present at $r_1$. Such a function

was later linked by Savin [29] to a local excess of kinetic energy due to the Pauli repulsion and reformulated by taking the homogenous electron gas as reference. That way, the ELF function (denoted $\eta$) can be interpreted as a measure of the Pauli repulsion in the atomic or molecular space and allows an access to the probability of finding two same spin electrons:

$$\eta(\mathbf{r}) = \frac{1}{1 + (\frac{D}{D_0})^2} \tag{6-11}$$

where D is a measure of kinetic energy excess and $D_0$ is the kinetic energy of a same density homogenous electron gas. ELF is defined to have values restricted between 0 and 1 in order to tend to 1 where parallel spins are highly improbable (there is therefore a high probability of opposite-spin pairs), and to zero in regions where there is a high probability of same-spin pairs. The ELF function can be interpreted as a signature of the electronic-pair distribution but, in contrast to pair functions, it can be more easily calculated and interpreted.

Once computed on a 3D grid from a given ab initio wave function, the ELF function can be partitioned into an intuitive chemical scheme [30]. Indeed, core regions, denoted C(X), can be determined for any atom, as well as valence regions associated to lone pairs, denoted V(X), and to chemical bonds (V(X,Y)). These ELF regions, the so-called basins (denoted $\Omega$), match closely the domains of Gillespie's VSEPR (Valence Shell Electron Pair Repulsion) model. Details about the ELF function and its applications can be found in a recent review paper [31].

It has been recently shown [12] that the ELF topological analysis can also be used in the framework of a distributed moments analysis as was done for Atoms in Molecules (AIM) by Popelier and Bader [32, 33]. That way, the $M_0(\Omega)$ monopole term corresponds to the opposite of the population (denoted N):

$$M_0(\Omega) = -\int_\Omega \rho(\mathbf{r})d\tau = -N(\Omega) \tag{6-12}$$

The first moments or dipolar polarization components of the charge distribution are defined by three-dimensional integrals for a given basin $\Omega$ according to:

$$M_{1,x}(\Omega) = -\int_\Omega (x - X_c)\rho(\mathbf{r})\,d\tau$$

$$M_{1,y}(\Omega) = -\int_\Omega (y - Y_c)\rho(\mathbf{r})\,d\tau \tag{6-13}$$

$$M_{1,z}(\Omega) = -\int_\Omega (z - Z_c)\rho(\mathbf{r})\,d\tau$$

where Xc, Yc, and Zc are the Cartesian coordinates of the basin centres.

The five second-moment spherical tensor components can also be calculated and are defined as the quadrupolar polarization terms. They can be seen as the ELF basin equivalents to the atomic quadrupole moments introduced by Popelier [32] in the case of an AIM analysis:

$$M_{2,zz}(\Omega) = -\frac{1}{2}\int_\Omega (3(z-Z_c)^2 - \mathbf{r}^2)\rho(\mathbf{r})\,d\tau$$

$$M_{2,x^2-y^2}(\Omega) = -\frac{\sqrt{3}}{2}\int_\Omega [(x-X_c)^2 - (y-Y_c^2)]\rho(\mathbf{r})\,d\tau$$

$$M_{2,xy}(\Omega) = -\sqrt{3}\int_\Omega (x-X_c)(y-Y_c)\rho(\mathbf{r})\,d\tau \qquad (6\text{-}14)$$

$$M_{2,xz}(\Omega) = -\sqrt{3}\int_\Omega (x-X_c)(z-Z_c)\rho(\mathbf{r})\,d\tau$$

$$M_{2,yz}(\Omega) = -\sqrt{3}\int_\Omega (y-Y_c)(z-Z_c)\rho(\mathbf{r})\,d\tau$$

The first- or second-moment basin magnitude is then defined as the square root of the sum of squared components:

$$|\mathbf{M}(\Omega)| = \sqrt{\sum_i M_i(\Omega)^2} \qquad (6\text{-}15)$$

Thanks to the invariance of the magnitude of any multipole rank ($|M1|$ or $|M2|$) with respect to the axis for a given bond or lone pair, the approach allows us to compare the dipolar or quadrupole polarization of a given basin in different chemical environments.

That way, the Distributed Electrostatic Moments based on the ELF Partition (DE-MEP) allows computing of local moments located at non-atomic centres such as lone pairs, $\sigma$ bonds and $\pi$ systems. Local dipole contributions have been shown to be useful to rationalize inductive polarization effects and typical hydrogen bond interactions. Moreover, bond quadrupole polarization moments being related to a $\pi$ character enable to discuss bond multiplicities, and to sort families of molecules according to their bond order.

**(b) Applications:** It is then possible to compute a chemically intuitive distributed analysis of electrostatic moments based on ELF basins. As this partition of the to-tal charge density provides an accurate representation of the molecular moments (dipole, quadrupole etc. . .), the distributed ELF electrostatic moments allows the computation of local moments located at non-atomic centers such as lone pairs, bonds and $\pi$ systems. It has been recently shown [12] that local dipole contributions

*Figure 6-2.* ELF moments for the canonical water dimer (dipole moments in Debye, $M_1$ in au.)

V(O) N = 2.37 e
|M₁| = 0.951

Acceptor
μ = 2.06 D

Donor
μ = 2.09 D

V(O) N = 2.50
|M₁| = 0.998

V(O) N = 2.15
|M₁| = 0.928

*Figure 6-3.* Correlation between the quadrupolar polarization (|**M**2|) of CO bonds in selected molecule and of the bond multiplicity

*Figure 6-4.* Electron localization function domains (concentration of electrons) in glycine. Lone pair domains are displayed in red

can be used to rationalize inductive polarization effects and so should be able to give some new insight towards a better understanding of local density modifications due to the cooperative effects. Figure 6-2 shows the difference of local dipole moment within the canonical water dimer. We can clearly distinguish the acceptor molecule from the donor one. Moreover, the local M1 value of the lone pair involved in the hydrogen bond (i.e. Acceptor molecule) is clearly higher than for the lone pairs of the donor molecule.

Following the same idea, it has been shown that the quadrupole moment of a bond was related to its $\pi$ character, allowing the discussion on its multiplicity. Then it becomes possible to discuss the influence of the intra- or inter-molecular environment on a given constituent of a molecule. Figure 6-3 displays such influence on a $C{=}O$ bond through a large set of molecules.

In a recent study of the transferability of moments, it has shown that stable trends are actually observed for the chemical bond features along investigated test peptide chains (Figure 6-4 and Table 6-4).

Such results are interesting for force field development as they clearly establish the existence of conserved "electrostatic blocks" within amino acids, an encouraging step for transferability of force field parameters.

## 6.3. DEVELOPMENT OF NEXT GENERATION POLARIZABLE FORCE FIELDS: FROM SIBFA TO GEM

### 6.3.1. Sum of Interaction Between Fragments Ab Initio (SIBFA)

SIBFA [2, 13] is a polarizable molecular mechanics procedure, formulated as a sum of five energy contributions, each of which is destined to reproduce its counterpart from reference EDA ab initio computations. The intermolecular interaction energy is formulated as:

$$\Delta E_{int} = E_{MTP^*} + E_{rep^*} + E_{pol} + E_{ct} + E_{disp}(+E_{LF}) \qquad (6\text{-}16)$$

*Table 6-4.* Local dipole contributions and magnitude of the first and second moments of some basins of typical basins involved in the main chain $C^2H(NH_2)C^1O^1O^2H$ of amino acid

| Amino acid | V($C^1$, $O^1$) | | V($C^1$, $O^2$) | | V($C^2$, N) | | V(N) | | | V($O^2$) | | | Total |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | $|M_1|$ | $|M_2|$ | $|M_1|$ | $|M_2|$ | $|M_1|$ | $|M_2|$ | $|M_1|$ | $|M_2|$ | $|\mu|$ | $|M_1|$ | $|M_2|$ | $|\mu|$ | $|\mu|^{a,b}$ |
| Glycine[c,d] | 0.255 (0.247) | 1.249 (1.221) | 0.052 (0.050) | 0.257 (0.270) | 0.183 (0.207) | 0.166 (0.177) | 0.952 (1.055) | 0.125 (0.100) | 11.6 (12.0) | 3.308 (2.955) | 2.446 (2.136) | 11.6 (11.5) | 0.50 (0.48) 0.48 (0.46) |
| Valine[c] | 0.281 | 1.241 | 0.059 | 0.256 | 0.177 | 0.173 | 0.931 | 0.167 | 11.5 | 3.296 | 2.417 | 12.2 | 0.55 (0.55) |
| Tyrosine[c] | 0.270 | 1.210 | 0.059 | 0.268 | 0.177 | 0.172 | 0.929 | 0.161 | 11,5 | 3.286 | 2.454 | 12,0 | 0.95 (0.92) |

[a] Total molecular dipole of the amino acid in a.u.
[b] The values given in parentheses are obtained from the SCF calculation provided by the GAUSSIAN03 software.
[c] Optimized at B3LYP/6-31+G(d,p) level of computation.
[d] The values given in parentheses correspond to a single point calculation at the B3LYP/Aug-cc-pVTZ level of computation.

Which denotes respectively the short-range penetration corrected electrostatic multipolar ($E_{MTP^*}$) energy, short-range repulsion ($E_{rep^*}$), polarization ($E_{pol}$), charge-transfer ($E_{ct}$), and dispersion ($E_{disp}$) contributions. In presence of an open-shell cation, a ligand field correction is introduced ($E_{LF}$).

The connectedness of SIBFA to quantum chemistry stems from the use of distributed multipoles and polarizabilities. They are derived from the molecular orbitals of any given molecular fragment using procedures due to Claverie and coworkers [34] concerning the multipoles and by Garmer and Stevens concerning the polarizabilities [35, 36]. They are then stored in the SIBFA library of fragments along with the fragment internal geometry and types of successive atoms and used in subsequent inter- or intramolecular interactions that involve that fragment. SIBFA can be seen as a set of parametric equations aiming to reproduce the required integrals produced by Localized Molecular Orbitals Theory.

We have previously [2] emphasized the features that an MM methodology should have in view of a meaningful reproduction of QC, namely separability, anisotropy, non-additivity and transferability.

Non-additivity and anisotropy of the interaction potential are critical features in molecular recognition and docking. Non-additivity in SIBFA stems from both second-order contributions, $E_{pol}$ and $E_{ct}$. That of $E_{pol}$ stems from the vector addition of the polarizing fields on a given centre and the use of the square of its norm. Iteratively accounting for the effects of the induced dipoles further enhances non-additivity. That of $E_{ct}$ is conferred by the modulation of the ionization potential of the electron donor on the one hand, and of the electronic affinity of the electron acceptor on the other hand, by the electrostatic potential that each undergoes in a multimolecular complex. Moreover, such potentials embody components due to the induced dipoles, whose amplitudes themselves depend non-additively upon the fields. An additional coupling to nonadditive polarization effects stems from the increase of the effective radius of the electron donor, intervening in the exponential of $E_{ct}$, by a term proportional to the magnitude of the field undergone by the electron donor.

The anisotropy of $E_{MTP^*}$ stems for the use of distributed multipoles on atoms and on the barycentres of the chemical bonds, thus advancing beyond the assumption of spherical symmetry incurred by the use of atom-centred point-charges. That of $E_{pol}$ stems from: (i) the multipolar nature of the polarizing field; (ii) the use of lone-pair polarizabilities that are off-centred, being located on the barycentres of the Boys localized lone-pair orbitals; (iii) and the use of polarizability tensors instead of scalars. The anisotropies of both short-range contributions, $E_{rep}$ and $E_{ct}$, which are overlap-dependent terms, is conferred by the use of localized lone-pairs accounting for hybridization. The anisotropy of $E_{disp}$ is, similarly, conferred by the introduction of fictitious atoms on the localized lone pairs. We detail here the methodology used for each one of the components of the SIBFA intermolecular interaction energy. Such equations have been shown to be transferable for intermolecular interactions, see Reference [2] and references therein.

- **Multipolar Electrostatic contribution: penetration corrected EMTP***

In SIBFA, electrostatics is computed upon using distributed multipoles (monopoles, dipoles, quadupole) located on atoms and bond midpoints as:

$$E_{MTP} = E_{mono-mono} + E_{mono-dip} + E_{mono-quad} + E_{dip-dip} + E_{dip-quad} + E_{quad-quad} \tag{6-17}$$

If we review the separated components of EMTP, the electrostatic energy appears mainly dominate by the terms involving the charge (Table 6-5).

However, if we analyse the functional form used to compute the charge–charge interaction,

$$E_{mono-mono} = q_i \cdot q_j / \mathbf{r} \tag{6-18}$$

where $\mathbf{r}$ is the distance between $q_i$ and $q_j$, we can easily see that it remains very different from the quantum chemistry formulation [6]:

$$E_c = -2 \sum_i \sum_{\nu} Z_{\nu} \int (|\varphi_i(1)|^2)/(r_{1_{\nu}}) d\iota_1 - 2 \sum_j \sum_{\mu} Z_{\mu} \int (|\varphi_i(1)|^2)/(r_{2_{\mu}}) d\iota_2$$

$$+ 4 \sum_i \sum_j \int (|\varphi_i(1)|^2 |\varphi_j(2)|^2)/(r_{12}) d\iota_1 d\iota_2 + \sum_{\mu} \sum_{\nu} Z_{\mu} \cdot Z_{\nu}/r_{\mu\nu} \tag{6-19}$$

where $\mu$ and $\varphi_i$ are respectively the nucleus and the unperturbed MOs of monomer A ; and $\nu$ and $\varphi_j$, those of monomere B.

Indeed, the ab initio integrals exhibit an exponential decay at short-range which is not present in any of the EMPT energy terms. This comportment of integrals is at the origin of the so-called penetration energy, an overlap dependant term which is, by definition, absent of the long-range multipolar approximation.

In our approach [18, 37], we have modified the formulation of the terms involving the charges (mono-mono, mono-dip and mono quad term) to screen the electrostatic interaction.

*Table 6-5.* Contributions to the multipolar electrostatic energy (kcal/mol) for various complexes at their equilibrium geometry

| Complexes | mono-mono. | mono-dip. | mono-quad. | dip-dip. | dip-quad. | Quad-quad. |
|---|---|---|---|---|---|---|
| $(H_2O)_2$ linear | −3.3 | −2.6 | −1.1 | −0.5 | −0.1 | −0.3 |
| $(HCONH_2)_2$ linear | −7.3 | −2.8 | 1.6 | 0.1 | −0.7 | 0.6 |
| $Cu^{2+} - H_2O$ | −46.2 | −27.9 | 0.9 | 0.0 | 0.0 | 0.0 |
| $HCOO^- - H_2O$ monodentate | −11.3 | −4.9 | −0.3 | −0.1 | −0.7 | 0.6 |
| $H_3CNH_3^+ - H_2O$ | −12.5 | −7.6 | 0.1 | −0.2 | −0.1 | 0.0 |

First, we have modified the mono-mono term (now denoted $\mathbf{E_{mono\text{-}mono*}}$) to propose a functional form mimicking the three terms present in ab initio, namely the nucleus–nucleus repulsion, the electron–nucleus attraction and the electron–electron repulsion. For two interacting centers i and j, the modified mono-mono term is:

$$E_{mono\text{-}mono*} = [Z_iZ_j - \{Z_i(Z_j - q_j)(1 - \exp(-\alpha_i.r)) + Z_j(Z_i - q_i)$$
$$(1 - \exp(-\alpha_j.r))\}$$
$$+ (Z_i - q_i)(Z_j - q_j)(1 - \exp(-\beta_i.r))(1 - \exp(-\beta_j.r))]^*(1/r) \qquad (6\text{-}20)$$

Where $Z_i$ and $Z_j$ are the valence electrons for the i and j atoms. This number is set to 0 for sites located on bonds. $\alpha_i$ and $\beta_i$ are parameters depending on effective van der Waals radii (denoted $r_{vdw}$) and given by:

$$\boldsymbol{\alpha_i = \gamma/r_{vdw\ i}\ \textbf{and}\ \beta_i = \delta/r_{vdw\ i}}$$

$\gamma$ and $\delta$ are fixed parameters depending on the reference ab initio level (methodology and chosen basis set). They are transferable to any atom and are evaluated once and for all upon fitting on a set of $H_2$ or $H_2O$ dimers geometries. For bonds monopoles, the $r_{vdw}$ values are given by the arithmetic mean of the radii forming the bond.

From a physical point of view, this new formulation includes exponential terms that are in agreement with the observed ab initio and experimental results. Moreover, it is easy to verify that the new expression converges to the classical one when r increases. That way, at long range, where the multipolar approximation is valid, the exponential part dies whereas, at short distances, the monopole–monopole interaction embodies a part of the penetration energy. Consequently, $\mathbf{E_{mono\text{-}mono*}}$ has the correct dependence at any range.

The second modification acts on the monopole-dipole term.

$$E_{mono\text{-}dip} = -\mu_j.\boldsymbol{\xi} \qquad (6\text{-}21)$$

Where $\xi$, the electric field due to a charge $q_i$ located at a point j is equal to:

$$\boldsymbol{\xi} = q_i\mathbf{r_{ij}}/r_{ij}^3 \qquad (6\text{-}22)$$

where $\mathbf{r_{ij}}$ is the vector oriented along r from i towards j.

We chose here to only modify $\xi$ to obtain the $E_{mono\text{-}dip}^*$ component:

$$E_{mono\text{-}dip*} = -\mu_j.\boldsymbol{\xi}^* \qquad (6\text{-}23)$$

where

$$\boldsymbol{\xi}^* = \{Z_i - (Z_i - q_i)(1 - \exp(-\eta r))\}.\mathbf{r_{ij}}/r_{ij}^3 \qquad (6\text{-}24)$$

$\eta$ is given by:

$$\eta = \chi/((r_{vdw\ i} + r_{vdw\ j}))/2$$

As for the charge-charge term, $\chi$ is a constant depending on the chosen level of reference ab initio computation and converges to the classical form as r increases.

At this level, our formulation includes a penetration correction for terms varying like $R^{-1}$ (monopole-monopole), like $R^{-2}$ (monopole-dipole) but does not include any correction for terms varying like $R^{-3}$ (dipole-dipole and monopole-quadupole). We then added a correction for the monopole-quadupole interaction.

For the evaluation of this term, we used a non usual formalism (see Ref. [20] and references therein), the so-called axial quadupole. Indeed, it is possible to define any quadupole as the sum of 3 axial quadupoles oriented towards the main axis $(e_1, e_2, e_3)$ of a local frame. We then obtain:

$$Q = \sum_{i=1}^{3} Qij(ei \otimes ej) \qquad (6\text{-}25)$$

Thanks to tensors mathematic properties, it is possible to add the same constant to each one of the diagonal terms, which allows the elimination of one of the axial qaudrupoles.

That way, the mono-quad interaction energy is given by:

$$E_{mono\text{-}quad} = E_{mq1} + E_{mq2} \qquad (6\text{-}26)$$

where $E_{mq1}$ and $E_{mq2}$ are respectively the monopole-axial quadupole interaction ($E_{mq1}$ and $E_{mq2}$ are different and represent the true quadupole) given by:

$$E_{mq1} = q^*(Q_1/2r^3)[3(\mathbf{a}.\mathbf{r}/r)^2 - 1] \qquad (6\text{-}27)$$

Where $\mathbf{a}$ is the unit vector defined by the local frame defining the axial quadupole, $\mathbf{r}$ is the vector oriented towards r from the monopole to the axial quadupole. $Q_a$ is the corresponding matrix element. Following the modification of the charge-dipole interaction, we introduced a modified mono-quad interaction, namely $E_{mono\text{-}quad^*}$:

$$E_{mono\text{-}quad^*} = E_{mq1^*} + E_{mq2^*} \qquad (6\text{-}28)$$

with:

$$E_{mq1^*} = \{Z_i - (Z_i - q_i)(1 - \exp(-\varphi r))\{^*\varphi(Q_a/2r^3)[3(\mathbf{a}.\mathbf{u})^2 - 1] \qquad (6\text{-}29)$$

$\varphi$ is given by:

$$\varphi = \Omega/(r_{vdw\,i} + r_{vdw\,j})/2$$

$\Omega$ is a constant dependant on the chosen level of reference ab initio computation.

To conclude, the penetration corrected $E_{MTP^*}$ interaction energy is then computed as:

$$E_{MTP^*} = E_{\text{mono-mono}^*} + E_{\text{mono-dip}^*} + E_{\text{mono-quad}^*} + E_{\text{dip-dip}} + E_{\text{dip-quad}} + E_{\text{quad-quad}}$$

- **Short Range exchange-repulsion: $\mathbf{E_{rep}^*}$**

To determine an expression of $E_{rep}^*$, we chose to follow the theoretical results of Murrell [38, 39] who proposed a simplified ab initio perturbation scheme to represent the exchange-repulsion energy based on an overlap expansion of localized Molecular Orbitals (MOs). Studying the interaction between hydrogen atoms, they shown that a $KS^2/R$ relation, S being the overlap between MOs, was not able to accurately reproduce the exchange-repulsion energy. They then proposed an extended $S^2(AR^{-1}+BR^{-2})$ expression that we have used to formulate $E_{rep}^*$. Following early expression (see References [2, 13] and references therein) based on LMOs approximation derived by Claverie, we have expressed $E_{rep}^*$ [18, 40, 41] as a sum of bond-bond, bond-lone pair and lone pair-lone pair repulsion:

$$\mathbf{E_{rep}^*} = C_1\left(\sum_{AB}\sum_{CD}\mathbf{rep}^*(AB, CD) + \sum_{AB}\sum_{L\gamma}\mathbf{rep}^*(AB, L\gamma) + \right.$$
$$\left. \sum_{L\alpha}\sum_{CD}\mathbf{rep}^*(L_\alpha, CD) + \sum_{L\alpha}\sum_{L\gamma}\mathbf{rep}^*(L_\alpha, L_\gamma)\right) \tag{6-30}$$

Where each one of the repulsion term includes two components: one varying like $1/R$, the other like $1/R^2$:

$$\mathbf{Rep}^*(AB, CD) = N_{occ}(AB)N_{occ}(CD)S^2(AB, CD)/R_{AB,CD} + N_{occ}(AB)N_{occ}$$
$$(CD)S^2(AB, CD)/(R_{AB,CD})^2 \tag{6-31}$$

AB and CD denoted the center positions of the bonds formed by atoms A and B; and C and D respectively. $L_\alpha$ and $L\gamma$ represent the lone pair positions. As we will see, this formulation takes into account bonds and lone pairs hybridation, each one of the term depending of an overlap functional. $N_{occ}(AB)$ and $N_{occ}(CD)$ are the electron occupation numbers of the AB and BC bonds. Therefore, $N_{occ}$ is equal to 2 for usual bonds and lone pairs. $R_{AB,CD}$ is the distance between the barycenters of the AB and CD bonds.

The S overlap expression relies on the general situation where 4 atoms form 2 bonds having their valence electrons involved in $sp^n$ hybrid Mos. In the context of Slater orbitals, $c_s$ et $c_p$ are the hybridation coefficients and, for example, S is then formulated for the bond-bond term as:

$$\mathbf{S(AB, CD)} = (c_{sI}c_{sk}\langle 2_{sI}2_{sk}\rangle + c_{pI}c_{sK}\langle 2p_{\sigma I}2_{sK}\rangle\mathbf{cos(IJ, IK)}$$
$$+ c_{pK}c_{sI}\langle 2p_{\sigma K}2_{sI}\rangle\mathbf{cos(KL, KI)} + c_{pI}c_{pK}\langle 2p_{\sigma I}2_{\sigma K}\rangle\mathbf{cos(IJ, IK)cos(KL, KI)}) \tag{6-32}$$

To simplify the problem, we introduce the following approximation:

$$\langle 2\mathbf{p}_{\sigma\mathbf{A}}2_{\mathbf{sC}}\rangle = \mathbf{m}_{\mathbf{AC}}\langle 2_{\mathbf{sA}}2_{\mathbf{sC}}\rangle \tag{6-33}$$

$\mathbf{m}_{\mathbf{AC}}$ is a parameter obtained from computations of overlap integrals between atoms A and B obtained from Mulliken [42] and Roothaan [43] approximations using Slater orbitals. Their values are tabulated and depend on a given atoms couple.

$\langle 2_{\mathbf{sA}}2_{\mathbf{sC}}\rangle$ can be approximate by an exponential depending on the distance separating atoms A and C and modulated on effective van der Waals radii:

$$\langle 2_{\mathbf{sA}}2_{\mathbf{sC}}\rangle = \mathbf{M}_{\mathbf{AC}}\exp(-\alpha\rho_{\mathbf{AC}}) \tag{6-34}$$

with

$$\rho_{\mathbf{AC}} = \mathbf{R}_{\mathbf{AC}}/4\sqrt{\mathbf{W}_A\mathbf{W}_C}$$

and

$$\mathbf{M}_{\mathbf{AC}} = \sqrt{\mathbf{K}_{AC}(1 - \mathbf{Q}A/N_{VAL}^A)(1 - Qc/N_{VAL}^C)}$$

$\mathbf{Q}_A$ and $\mathbf{Q}_C$ are the charges obtained from the multipolar expansion of the interacting A and C molecular charge distributions, $N_{VAL}^A$ and $N_{VAL}^C$ being their respective number of valence electrons. $W_A$ and $W_C$ are the A and C atoms effective van der Waals radii. $K_{AC}$ is a proportionality factor tabulated upon the atomic numbers of the A and C atoms. $\alpha$ is a constant fixed to 12.35. The same treatment is applied to the others terms of the repulsion energy.

It is important to point out that recent results on density based overlap integrals [16] confirm the interest of the formulation of $E_{rep}{}^*$ as a sum of bond-bond, bond-lone pair and lone pair-lone pair repulsion: indeed, core electrons do not contribute to the value of the overlap integrals.

• **Polarization contribution**

$\mathbf{E_{pol}}$ also relies on a local picture as it uses polarizabilities distributed at the Boys LMOs centroids [44] on bonds and lone pairs using a method due to Garmer et al. [35]. In this framework, polarizabilities are distributed within a molecular fragment an therefore, the induced dipoles do not need to interact together (like in the Applequist model) within a molecule as their value is only influenced by the electric fields from the others interacting molecules.

The general expression of the polarization energy at center I located at the centroid of an LMO of a A molecule is:

$$E_{pol}(i) = -0.5 \sum_j \Delta\mu(i) E_0(j)$$

$$\Delta\mu(i) = \alpha(i) \sum_j^{xyz} E(\Delta\mu(i)) + E_0(j)$$

(6-35)

$E_0$ and E $(\Delta\mu(i))$ are respectively the electric fields generated by the permanent and induced multipoles moments. $\alpha(i)$ represents the polarisability tensor and $\Delta\mu(i)$ is the induced dipole at a center i. This computation is performed iteratively, as $E_{pol}$ generally converges in 5–6 iterations. It is important to note that in order to avoid problems at the short-range, the so-called polarization catastrophe, it is necessary to reduce the polarization energy when two centers are at close contact distance. In SIBFA, the electric fields equations are "dressed" by a Gaussian function reducing their value to avoid such problems.

The initial electric field generated by a centre i of molecule A on a center j of molecule B is denoted $E_{i\rightarrow j}^{init}$ is modified by a Gaussian function denoted S [45] as:

$$E_{i\rightarrow j}^{final} = (1 - S(i, j)) E_{i\rightarrow j}^{init}$$

$$S(i, j) = \Theta_i \text{E} \exp(-\Gamma(R_{ij}^2)/(r_{vdw}(i) + r_{vdw}(j))$$

(6-36)

$R_{ij}$ is the distance between centers i and j. $\Theta_i$ is the monopole associated to center i. E and $\Gamma$ are empirical parameters associated to each atom types as $r_{vdw}$ are the atom effective radii.

It is important to point that parametrization procedure of the short-range damping is really important. In SIBFA, in order to embody short-range penetration and exchange-polarization effects, the fit is performed upon CSOV (or RVS) polarization energy which embodies exchange effects (see Section 6.1). To do so, the SIBFA polarization energy "prior iterating" is adjusted to the CSOV value which corresponds to the relaxation of a molecule A in the field of a frozen B molecule (before the compution of higher-orders induction terms δE). Details can be found in Refs. [18, 19].

- **Charge transfer contribution**

As for $E_{rep}$*, **$E_{ct}$** is derived from an early simplified perturbation theory due to Murrel [46]. Its formulation [47, 48] also takes into account the $L_\alpha$ lone pairs of the electron donor molecule (denoted molecule A). Indeed, they are the most exposed in this case of interaction (see Section 6.2.3) and have, with the π orbital, the lowest ionization potentials. The acceptor molecule is represented by bond involving an hydrogen (denoted BH) mimicking the set, denoted $\phi_{*BH}$, of virtual bond orbitals involved in the interaction.

$E_{ct}$ is expressed as:

$$\mathbf{E_{tc}} = -2\mathbf{C} \sum_{\mathbf{L\alpha}} \mathbf{N_{occ}(\alpha)(T_{\alpha\beta*})^2/\Delta E_{\alpha\beta*}}$$

(6-37)

C is a constant which has been calibrated in order to reproduce $E_{ct}$ at the equilibrium geometry of the water dimer. This value is transferable for all non metallic acceptors. $N_{occ}(\alpha)$ is the occupation number $L_\alpha$ lone pair.

$T_{\alpha\beta*}$ is a function of:

i) The transition density overlap of the $L\alpha$ donor lone pairs and the bond BH centroid, expressed with the same approximations as $E_{rep*}$.

ii) The electrostatic potential applied on A by all the other molecules ($\sum_C \mathbf{V_{C\rightarrow A}}$). It is important to point out that the fields due to the polarization converged induced dipoles are taken into account in order to introduce an explicit link between polarisation and charge transfer.

$\Delta E_{\alpha\beta*}$ is the energy is the energy required to allow an electron transfer from an orbital $\alpha$ of molecule A towards a virtual $\beta*$ orbital on molecule B. It can be expressed as:

$$\mathbf{\Delta E_{\alpha\beta*} = (I_{L_\alpha} + \sum_C V_{C\rightarrow A}) - (A_{\beta*} + \sum_C V_{C\rightarrow B})} \qquad (6\text{-}38)$$

$I_{L_\alpha}$ is the ionization potential of the $L_\alpha$ lone pair as $A_{\beta*}$ is the electronic affinity of the electron acceptor. Here also, the introduction of the final iterated field of induce dipole allow to take into account the many-body properties of $E_{ct}$.

Here also, $E_{ct}$ has its ab initio counterpart within the CSOV framework. The sum of $E_{pol}$ and $E_{ct}$ matchs the $E_{OI}$ contribution at the HF and DFT level and the SAPT induction when $\delta E$ remains small (see Section 6.1).

- **Dispersion contribution**

$\mathbf{E_{disp}}$ [49, 50] is coupled to an exchange-dispersion term and is computed as an expansion of $C_n$ terms: $\mathbf{C_6/Z^6}$, $\mathbf{C_8/Z^8}$ et $\mathbf{C_{10}/Z^{10}}$, Z being expressed as:

$$\mathbf{Z = r_{ij}/\sqrt{vdw_A.vdw_B}}$$

$r_{ij}$ is the distance between atoms i and j; $vdw_A$ and $vdw_B$ are the effective radii of the involved atoms. The $C_6$, $C_8$ and $C_{10}$ coefficients are empirical parameters adjusted on $H_2$ dimers SAPT computations. Each one of the Cn terms are damped at short-range by the following function:

$$\mathbf{E_{damp}(n) = (1/R^n)L_{ij}} \exp(-\mathbf{a_{damp}(n)D(i,j))} \qquad (6\text{-}39)$$

where:

$$D(i,j) = ((vdw_A + vdw_B)b_{damp}/R_{ij}) - 1$$

$L_{ij}$, $a_{damp}$ and $b_{damp}$ are parameters; n can be 6, 8 ou10.

This dispersion energy is coupled to an exchange-dispersion component:

$$\mathbf{E_{exch\text{-}disp} = L_{ij}(1 - Q_i/N_{val}(i))(1 - Q_j/N_{val}(j))C_{exch}} \exp(-\mathbf{\beta_{ecxc}Z}) \qquad (6\text{-}40)$$

$Q_i$ and $Q_j$ are the net charges of atoms i and j; $N_{val}(i)$ and $N_{val}(j)$ their number of valence electrons. $C_{exch}$ and $\beta_{exch}$ are empirical parameters. Some additional refinements exist within SIBFA as explicit addition of lone pairs for the exchange term [50].

• **Ligand field contribution**

To correct the energy function from ligand field effects (presently in the case of open-shell cations), SIBFA uses the formalism of the Angular Overlap Model (AOM) [51–53]. The AOM [52] is based on the fact that the relative changes of **d** orbitals energies caused by ligand field effects can be associated to the overlap of these orbitals with the ligands orbitals. As intermolecular overlap integral can be factorized into radial and angular parts, it is then possible to consider the radial part as a constant for a given intermolecular distance. That way, it can be introduced as a parameter, specific of a given metal/ligand couple, whereas the angular part can be exactly computed as it depends only on the relative orientation of the metal **d** orbitals and of those of the ligands. More precisely, the AOM treatment can be seen as an effective Hamiltonian built on the basis of **d** orbitals. Its evaluation uses spherical coordinates ($\theta$, $\phi$) and requires the diagonalization of an energy matrix. For each computation, each ligand is considered separately as the total matrix reflects the sum of the local perturbation of the **d** orbitals due to ligands as each matrix element is the sum of the contribution of each ligands. The construction of the energy matrix uses angular coefficient denoted $D_i$ (see Table 6-6) which give the values of overlap of a ligand involved in a $\sigma$ interaction with the metal.

In order to be able to evaluate the radial part in all point of space and to adapt the AOM to the SIBFA intermolecular potential, we have introduced an exponential dependence of the radial overlap following a procedure introduced by Woodley et al. [54]:

$$e_\lambda = a + b.\exp(-\alpha.r) \tag{6-41}$$

r is the metal-ligand distance; a, b and $\alpha$ are parameters, specific of a given metal-ligand couple.

It is then possible to construct the energy matrix [51, 53]:

*Table 6-6.* angular coefficients for the different **d** orbitals involved d in a $\sigma$ interaction with ligands

| i | $D_i(\theta_i, \phi_i)$ |
| --- | --- |
| $z^2$ | $1/2\ (3\cos^2\theta - 1)$ |
| $yz$ | $1/2\ (\sqrt{3}\sin(2\theta)\sin\phi)$ |
| $xz$ | $1/2\ (\sqrt{3}\sin(2\theta)\sin\phi)$ |
| $xy$ | $1/4\ (\sqrt{3}(1-\cos 2\theta)\sin 2\phi)$ |
| $x^2\text{-}y^2$ | $1/4\ (\sqrt{3}(1-\cos 2\theta)\cos 2\phi)$ |

$$H_{dd'} = \sum_l e_\lambda^l <d|l><l|d'> \tag{6-42}$$

Its eigenvalues ($\varepsilon_i$) correspond to the relative energies of the **d** orbitals of the metal. These energies are used to computed the ligand field energy contribution denoted $E_{LF}$.

For a $d^n$ system:

$$E_{LF} = -2 \sum_{i=1}^{5} \varepsilon_i + \sum_{i=1}^{n} \rho_i \varepsilon_i \tag{6-43}$$

Where $\rho_i$ are the d orbital occupation numbers (0, 1 or 2).

### 6.3.2. The Gaussian Electrostatic Model (GEM)

As the SIBFA approach relies on the use of distributed multipoles and on approximation derived form localized MOs, it is possible to generalize the philosophy to a direct use of electron density. That way, the Gaussian electrostatic model (GEM) [2, 14–16] relies on ab initio-derived fragment electron densities to compute the components of the total interaction energy. It offers the possibility of a continuous electrostatic model going from distributed multipoles to densities and allows a direct inclusion of short-range quantum effects such as overlap and penetration effects in the molecular mechanics energies.

#### 6.3.2.1. *From Density Matrices to GEM*

This method relies on the use of an auxiliary gaussian basis set to fit the molecular electron density obtained from an ab initio one-electron density matrix:

$$\tilde{\rho} = \sum_{k=1}^{N} x_k k(r) \approx \rho = \sum_{\mu\nu} P_{\mu\nu} \phi\mu(r)\phi_\nu^*(r) \tag{6-44}$$

Do so, we use the formalism of the variational density fitting method [55, 56] where the Coulomb self-interaction energy of the error is minimized:

$$E_2 = \frac{1}{2} \iint \frac{[\rho(r_1) - \tilde{\rho}(r_1)][\rho(r_2) - \tilde{\rho}(r_2)]}{|r_1 - r_2|} dr_1 dr_2 = \langle \rho - \tilde{\rho} \| \rho - \tilde{\rho} \rangle \tag{6-45}$$

inserting the right hand of Eq. (6-44) into Eq. (6-45), we obtained:

$$E2 = \frac{1}{2} \sum_{\mu,\nu} \sum_{\sigma,\tau} P_{\mu\nu} P_{\sigma\tau} \langle \mu\nu \| \sigma\tau \rangle - \sum_l x_l \sum_{\mu,\nu} P_{\mu\nu} \langle \mu\nu \| l \rangle + \frac{1}{2} \sum_k \sum_l x_k x_l \langle k \| l \rangle \tag{6-46}$$

E$_2$ from Eq. (6-3) can be minimized with respect to the expansion coefficients x$_l$ and a linear system of equations can be obtained:

$$\frac{\partial E_2}{\partial x_l} = -\sum_{\mu\nu} P_{\mu\nu} \langle \mu\nu \| l \rangle + \sum_{k} xk \langle k \| l \rangle \tag{6-47}$$

Equation (6-5) is used to determine the coefficients:

$$\mathbf{x} = \mathbf{A}^{-1}\mathbf{b} \tag{6-48}$$

Where

$$b_l = \sum_{\mu\nu} P_{\mu\nu} \langle \mu\nu \| l \rangle \text{ and } A_{kl} = \langle k \| l \rangle$$

In a standard density fitting, the determination of the coefficients requires the use of a modified singular value decomposition (SVD) procedure in which the inverse of an eigenvalue is set to zero if it is below a certain cutoff. A cutoff value of $10^{-8}$ has been previously determined [14] to be acceptable for the molecules which will be under study. In addition to the SVD approach, we have also implemented noise reduction techniques for the fitting procedure as this method can produces numerical instabilities (noise) when the number of basis functions starts to grow and when higher angular momentum are used (these instabilities are also present when using only *s*-type spherical [14] functions albeit to a lower extent). Several strategies have been implemented [15]. Among them, we used the the Tikhonov regularization formalism and a damped Coulomb operator $\hat{O} = erfc(\beta r)/r$ procedure in order to localize the integrals to increase the calculation speed. Alternatively to the DF procedure, it is possible to perform such a fit using density and electrostatic grids [16]. That way, the ab initio calculated properties (density, electrostatic potential, and/or electric field) are fitted via a linear or nonlinear-least-squares procedure to the auxiliary basis sets (ABS). Neglecting the core contributions allows to perform more robust fit of the coefficients compared to the numerical grids and allows to reduce the number of functions and so the computational cost.

Using fitted densities expressed in a linear combination of Gaussian functions has the advantage that that the choice of Gaussian functions auxiliary basis set is not be restricted to Cartesian Gaussians. To use higher order angular momenta, normalized *Hermite Gaussian functions* can be preferred [2, 15, 16]. Indeed, the use of Hermite Gaussians in integral evaluation improves efficiency by the use of the McMurchie-Davidson (McD) recursion [57] since the expensive Cartesian-Hermite transformation is avoided. Obtaining the Hermite expansion coefficients from the fitted Cartesian coefficients is straightforward since Hermite polynomials form a basis for the linear space of polynomials. Moreover, Hermite Gaussians have a simple relation to elements of the Cartesian multipole tensor and can be used to multipoles distributed at the expansion sites. This smooth connection leads to a continuous

electrostatic model that can be used directly into second generation APMM such as SIBFA. It is important to note that unlike conventional multipole expansions, the spherical multipole expansion obtained from Hermite Gaussians has an intrinsic finite order, namely, the highest angular momentum in the ABS. This connection between multipoles and Hermite densities has its importance as, unlike s-type functions ($l = 0$), fitting coefficients with $l > 0$ (sp, spd . . .) are not invariant by rotation. These coefficients must be transformed for each molecular fragment orientation in order to compute interaction energies. Such a transformation can be achieved [15] by defining both a *global* orthogonal coordinate system frame and a *local* orthogonal coordinate frame for each fragment fitting site.

### 6.3.2.2.    *Computing Integrals for Molecular Mechanics*

The GEM force field follows exactly the SIBFA energy scheme. However, once computed, the auxiliary coefficients can be directly used to compute integrals. That way, the evaluation of the electrostatic interaction can virtually be exact for an perfect fit of the density as the three terms of the coulomb energy, namely the nucleus–nucleus repulsion, electron–nucleus attraction and electron–electron repulsion, through the use of $\tilde{\rho}$ [2, 14–16, 58].

$$\mathrm{E_{coulomb}} = \frac{Z_A Z_B}{r_{AB}} - \int \frac{Z_A \tilde{\rho}^B(\mathrm{r_B})}{r_{AB}} dr - \int \frac{Z_B \tilde{\rho}^A(\mathrm{r_A})}{r_{AB}} dr + \int \frac{\tilde{\rho}^A(\mathrm{r_A}) \tilde{\rho}^B(\mathrm{r_B})}{r_{AB}} dr \tag{6-49}$$

To complete the first order terms, the exchange–repulsion energy can be evaluated through an overlap model [14, 59] as:

$$\mathrm{E_{exch/rep}} \approx \mathrm{KS_\rho} \tag{6-50}$$

Where:

$$S_\rho = \int \rho a(r) \rho b(r) dr \approx \int \tilde{\rho} a(r) \tilde{\rho} b(r) dr$$

As electric fields and potential of molecules can be generated upon distributed $\tilde{\rho}$, the second order energies schemes of the SIBFA approach can be directly fueled by the density fitted coefficients. To conclude, an important asset of the GEM approach is the possibility of generating a general framework to perform Periodic Boundary Conditions (PBC) simulations. Indeed, such process can be used for second generation APMM such as SIBFA since PBC methodology has been shown to be a key issue in polarizable molecular dynamics with the efficient PBC implementation [60] of the multipole based AMOEBA force field [61].

### 6.3.2.3. *Using Periodic Boundary Conditions to Increase Computational Efficiency*

In this section we describe the methods to extend Ewald sum methodologies to accelerate the calculation of the intermolecular interactions using PBC. For simplicity, we begin with a generalization of Ewald sums to interacting spherical Hermite Gaussians (e.g. GEM-0 [14]). This is followed by the extension to arbitrary angular momentum. Finally, we describe the implementation of methods to speed up both the direct an reciprocal terms in the Ewald sum [62].

*N.1 Spherical charge densities*

To begin let $U$ denote a unit cell such that it contains the set of points $\mathbf{r}$ with associated fractional coordinates $s_1, s_2, s_3$ satisfying $-1/2 \leq s_i \leq +1/2, i = 1, 2, 3$. Then the idealized infinite crystal $C$ is generated by the union of all periodic translations $\mathbf{U_n}$ of $U$, using the set of general lattice vectors $\mathbf{n}$. For the Ewald sum to be convergent, extra conditions need to be imposed. To that end, consider a large but finite crystal, i.e., let $P$ denote a closed, bounded region in space, centered at the origin (e.g., sphere, cube, etc.). For a positive integer $K$, let $\Omega(P, K)$ denote the set of lattice vectors such that $|\mathbf{n}|/K$ is in $P$. The Coulomb interaction of a spherical Gaussian charge distribution $\rho_1$ in a unit cell $\mathbf{U}_0$ centered at point $\mathbf{R}_1 \in \mathbf{U}_0$ with an exponent $\alpha_1$, i.e. $\rho_1 = (\alpha_1/\pi)^{3/2} \exp(-\alpha_1(\mathbf{r} - \mathbf{R_1}))$, interacting with a second normalized Gaussian charge distribution $\rho_2$ centered at point $\mathbf{R}_2 \in \mathbf{U}_0$ with exponent $\alpha_2$ together with all images in $\mathbf{U_n}$ for $\mathbf{n} \in \Omega(P, K)$ centered at $\mathbf{R}+\mathbf{n}$ can be shown to be:

$$
\begin{aligned}
E_{12} = &\sum_{\mathbf{n} \in \Omega(P,K)} \frac{erfc(\mu_{12\alpha_0}|\mathbf{R}_{12} - \mathbf{n}|) - erfc(\mu_{12}|\mathbf{R}_{12} - \mathbf{n}|)}{|\mathbf{R}_{12} - \mathbf{n}|} \\
&+ \frac{1}{\pi V} \sum_{\mathbf{m} \neq 0} \frac{\exp(-\pi^2 \mathbf{m}^2/\mu_{12\alpha_0})}{\mathbf{m}^2} \exp(-2\pi i \cdot (\mathbf{R}_{12})) \\
&- \frac{\pi}{\mu_{12\alpha_0}} + \frac{1}{\pi} H_{P,K}(\mathbf{R}_{12}) + \varepsilon(K),
\end{aligned}
\tag{6-51}
$$

where the first term corresponds to the direct part of the Ewald sum, the second to the reciprocal part, $H_{P,K}(\mathbf{R}_{12})$ is the surface term which depends on the dipole of the unit cell ($\mathbf{D}$), $\varepsilon(K)$ denotes a quantity that converges to 0 as $K \to \infty$, $\mathbf{m}$ denotes the reciprocal lattice vectors, $1/\mu_{12\theta} = 1/\alpha_1 + 1/\alpha_2 + 1/\alpha_0$, $1/\mu_{12} = 1/\alpha_1 + 1/\alpha_2$ and $\alpha_0$ is the Ewald exponent [62].

Equation (6-51) can be generalized to calculate the energy $E_{P,K}$ of $U$ interacting with the entire crystal $P$. Let $\rho_1 \dots \rho_N$ be normalized spherical Gaussian charge distributions (e.g. GEM-0) centered at $\{\mathbf{R_1} \dots \mathbf{R}_N\} = \mathbf{R}^{\{N\}} \in U$, and let $q_1 + \dots + q_N = 0$ (neutral unit cell). Then the energy of the central unit cell $\mathbf{U}_0$ within a large spherical crystal, due to the interactions of the Gaussian charge distributions $q_i \rho_i$ with each other and all periodic images within the crystal is given by

$$E_{S,K}(\mathbf{R}^{\{N\}}) = \frac{1}{2} \sum_{\mathbf{n}}' \sum_{i,j=1}^{N} q_i q_j \left\{ \frac{erfc(\mu_{ij\alpha_0}^{1/2}|\mathbf{R}_{ij} - \mathbf{n}|) - erfc(\mu_{ij}^{1/2}|\mathbf{R}_{ij} - \mathbf{n}|)}{|\mathbf{R}_{ij} - \mathbf{n}|} \right\}$$

$$+ \frac{1}{2\pi V} \sum_{\mathbf{m} \neq 0} \sum_{i,j=1}^{N} q_i q_j \frac{\exp(-\pi^2 \mathbf{m}^2 / \mu_{ij\alpha_0})}{\mathbf{m}^2} \exp(-2\pi i \mathbf{m} \cdot (\mathbf{R}_{ij}))$$

$$- \frac{\pi}{2V} \sum_{i,j=1}^{N} \frac{q_i q_j}{\mu_{ij\alpha_0}} - \sum_{i=1}^{N} q_i^2 \left( \frac{\mu_{ii\alpha_0}}{\pi} \right)^{1/2} + \frac{2\pi \mathbf{D}^2}{3V} + \varepsilon(K),$$

(6-52)

where $\mathbf{R}_{ij} = \mathbf{R}_i - \mathbf{R}_j$ and $\mathbf{D} = q_1 \mathbf{R}_1 + \ldots + q_N \mathbf{R}_N$ is the unit cell dipole.

In order to be able to calculate the reciprocal contribution in Eq. (6-52) it is necessary to grid the Gaussian densities. However, for large exponents (compact Gaussians) this can become intractable. To overcome this problem, the first GEM implementation relied on a method inspired by Fusti-Molnar and Pulay [63]. In this method, the individual Gaussian charge densities are classified into compact and diffuse Hermite Gaussian functions for a given $\alpha_0$. Thus, all Hermites with an exponent $\alpha_i \geq \alpha_0$ are considered compact, and the rest are considered diffuse. In this way, the interaction energy expressions may be re-expressed in order to calculate the contributions involving diffuse Hermites completely in reciprocal space [64].

This method was subsequently improved by noting that the $\alpha_0$ in $\mu_{ij\alpha_0}$ can be different for each pair $ij$ [15]. In this way, the Hermite charge distributions $q_i \rho_i$ are separated into compact and diffuse sets based on their exponents $\alpha_i$. Subsequently, $\alpha_0$ is chosen to be infinite for $ij$ pairs where at least one of the two Gaussians is diffuse. This ensures that all pairs involving diffuse Hermites are evaluated in reciprocal space. For all compact $ij$ pairs, $\alpha_0$ is chosen so that $\mu_{ij\alpha_0}$ is constant, that is, given $\beta > 0$, a Gaussian distribution $q_i \rho_i$ is classified as compact if $\alpha_i \geq 2\beta$ ($C$ set) and diffuse otherwise ($D$ set). Then, for $i, j \in C$, choose $\alpha_0$ so that $1/\mu_{ij\alpha_0} = 1/\alpha_i + 1/\alpha_j + 1/\theta = 1/\beta$. Otherwise, $\alpha_0$ is set to infinity. From this, the Coulomb energy of the spherical unit cell can be re-expressed as:

$$E_{S,K}(\mathbf{R}^N) = \frac{1}{2} \sum_{\mathbf{n}}' \sum_{(i,j) \in C \times C} q_i q_j \left\{ \frac{erfc(\beta^{1/2}|\mathbf{R}_{ij} - \mathbf{n}|) - erfc(\mu_{ij}^{12}|\mathbf{R}_{ij} - \mathbf{n}|)}{|\mathbf{R}_{ij} - \mathbf{n}|} \right\}$$

$$+ \frac{1}{2\pi V} \sum_{\mathbf{m} \neq 0} \sum_{(i,j) \in C \times C}^{N} q_i q_j \frac{\exp(-\pi^2 \mathbf{m}^2 / \beta)}{\mathbf{m}^2} \exp(-2\pi i \mathbf{m} \cdot (\mathbf{R}_{ij}))$$

$$+ \frac{1}{2\pi V} \sum_{\mathbf{m} \neq 0} \sum_{(i,j) \notin C \times C}^{N} q_i q_j \frac{\exp(-\pi^2 \mathbf{m}^2 / \mu_{ij})}{\mathbf{m}^2} \exp(-2\pi i \mathbf{m} \cdot (\mathbf{R}_{ij}))$$

$$- \frac{\pi}{2V} \sum_{(i,j) \notin C \times C}^{N} q_i q_j \left( \frac{1}{\beta} + \frac{1}{\alpha_i} + \frac{1}{\alpha_j} \right)$$

$$
-\sum_{i\in C}^{N} q_i^2 \left(\frac{\beta}{\pi}\right)^{1/2} - \sum_{i\notin C}^{N} q_i^2 \left(\frac{\alpha_i}{\pi}\right)^{1/2}
$$
$$
+\frac{2\pi \mathbf{D}^2}{3V} + \varepsilon(K),
$$
(6-53)

### N.2 Higher angular momentum charge densities

In the case of GEM, the auxiliary bases employed for the fitting of the molecular fragment include higher angular momentum Gaussians. In this case, Eq. (6-53) can be extended to account for the higher order Gaussians. As explained above, the fitted densities are expanded in a linear combination of Hermite Gaussians $\Lambda_{tuv}(\mathbf{r}, \alpha, \mathbf{R})$. Here, the Gaussian charge distribution is given by $\rho_i(\mathbf{r}, \mathbf{R}_i, \propto) = \sum_{l=1}^{L}\sum_{tuv} \mathbf{c}_{i,l,tuv}\Lambda_{tuv}(\mathbf{r}, \alpha_l, \mathbf{R}_i)$, where $\mathbf{c}_{i,l,tuv}$ are the Hermite coefficients, and $L$ denotes the different exponents in the ABS on center $i$. Based on this, the Coulomb energy of the total density within the spherical crystal is given by

$$
E_{S,K}(\rho^{\{N\}}) = \frac{1}{2}\sum_{\mathbf{n}}'\sum_{i=1}^{N}\sum_{l_i\in C}\sum_{t_iu_iv_i}\mathbf{c}_{i,l_i,t_iu_iv_i}\sum_{j=1}^{N}\sum_{l_j\in C}\sum_{t_ju_jv_j}(-1)^{(t_j+u_j+v_j)}\mathbf{c}_{j,l_j,t_ju_jv_j}
$$
$$
\times \left(\frac{\partial}{\partial \mathbf{R}_{ijx}}\right)^{t_i+t_j}\left(\frac{\partial}{\partial \mathbf{R}_{ijy}}\right)^{u_i+u_j}\left(\frac{\partial}{\partial \mathbf{R}_{ijz}}\right)^{v_i+v_j}
$$
$$
\times \left\{\frac{erfc(\beta^{1/2}|\mathbf{R}_{ij}-\mathbf{n}|)-erfc(\mu_{l_il_j}^{1/2}|\mathbf{R}_{ij}-\mathbf{n}|)}{|\mathbf{R}_{ij}-\mathbf{n}|}\right\}
$$
$$
+\frac{1}{2\pi V}\sum_{\mathbf{m}\neq 0}\frac{1}{\mathbf{m}^2}\exp(-\pi^2\mathbf{m}^2/2\beta)\sum_{l_1\in C}^{N}S_{l_1}(\mathbf{m})
$$
$$
\times \exp(-\pi^2\mathbf{m}^2/2\beta)\sum_{l_2\in C}^{N}S_{l_2}(-\mathbf{m})
$$
$$
+\frac{1}{2\pi V}\sum_{\mathbf{m}\neq 0}\frac{1}{\mathbf{m}^2}\sum_{(l_1,l_2)\notin C\times C}\exp(-\pi^2\mathbf{m}^2/\alpha_{l_1})
$$
$$
\times \exp(-\pi^2\mathbf{m}^2/\alpha_{l_2})S_{l_1}(\mathbf{m})S_{l_2}(-\mathbf{m})
$$
$$
-\frac{\pi}{2V}\sum_{l_1\in C}^{N}\sum_{l_2\in C}^{N}\sum_{i=1}^{N}\sum_{j=1}^{N}\mathbf{c}_{i,l_1,000}\mathbf{c}_{j,l_2,000}\left(\frac{1}{\beta}-\frac{1}{\alpha_{l_1}}-\frac{1}{\alpha_{l_2}}\right)
$$
$$
-\sum_{i=1}^{N}E_{self}(\rho_i)+\frac{2\pi \mathbf{D}^2}{3V}+\varepsilon(K),
$$
(6-54)

where the structure factors $S_l(\mathbf{m})$ are given by

$$S_l(\mathbf{m}) = \sum_{i=1}^{N}\sum_{tuv} \mathbf{c}_{i,l,tuv} \left(\frac{\partial}{\partial \mathbf{R}_{ijx}}\right)^{t_i+t_j} \left(\frac{\partial}{\partial \mathbf{R}_{ijy}}\right)^{u_i+u_j} \left(\frac{\partial}{\partial \mathbf{R}_{ijz}}\right)^{v_i+v_j} \exp(-2\pi i\mathbf{m}\cdot(\mathbf{R}_i))$$

(6-55)

$E_{self}(\rho_i)$ is given by

$$\begin{aligned}
&- \lim_{\mathbf{R}\to 0} \sum_{t_1 u_1 v_1}\sum_{t_2 u_2 v_2} (-1)^{(t_j+u_j+v_j)}\mathbf{c}_{j,l_j,t_j u_j v_j}\\
&\times \left(\frac{\partial}{\partial \mathbf{R}_x}\right)^{t_i+t_j}\left(\frac{\partial}{\partial \mathbf{R}_y}\right)^{u_i+u_j}\left(\frac{\partial}{\partial \mathbf{R}_z}\right)^{v_i+v_j}\\
&\times \left\{\sum_{(l_1,l_2)\in C\times C} \mathbf{c}_{i,l_1,t_1 u_1 v_1}\mathbf{c}_{i,l_2,t_2 u_2 v_2} \frac{erfc(\beta^{1/2}|\mathbf{R}|)}{|\mathbf{R}|}\right.\\
&\left. - \sum_{(l_1,l_2)\notin C\times C} \mathbf{c}_{i,l_1,t_1 u_1 v_1}\mathbf{c}_{i,l_2,t_2 u_2 v_2} \frac{erfc(\mu_{12}^{1/2}|\mathbf{R}|)}{|\mathbf{R}|}\right\}
\end{aligned}$$

(6-56)

where, similar to the previous section above, $1/\mu_{12} = 1/\alpha_{l_1} + 1/\alpha_{l_2}$ and $\mathbf{D}$ is the unit cell dipole [62].

### N.2 Computational speedup for the direct and reciprocal sums

Computational speedups can be obtained for both the direct and reciprocal contributions. In the direct space sum, the issue is the efficient evaluation of the erfc function. One method proposed by Sagui et al. [64] relies on the McMurchie-Davidson [57] recursion to calculate the required erfc and higher derivatives for the multipoles. This same approach has been used by the authors for GEM [15]. This approach has been shown to be applicable not only for the Coulomb operator but to other types of operators such as overlap [15, 62].

In the case of the reciprocal sum, two methods have been implemented, smooth particle mesh Ewald (SPME) [65] and fast Fourier Poisson (FFP) [66]. SPME is based on the realization that the complex exponential in the structure factors can be approximated by a well behaved function with continuous derivatives. For example, in the case of Hermite charge distributions, the structure factor can be approximated by

$$
\begin{aligned}
S_l(\mathbf{m}) \approx \lambda(z_1)\lambda(z_2)\lambda(z_3) & \sum_{k_1=-\frac{M_1}{2}+1}^{\frac{M_1}{2}} \sum_{k_2=-\frac{M_2}{2}+1}^{\frac{M_2}{2}} \sum_{k_3=-\frac{M_3}{2}+1}^{\frac{M_3}{2}} \sum_{i=1}^{N} \sum_{tuv} \mathbf{d}_{i,l,tuv} \\
& \times \sum_{l_1=-\infty}^{\infty} \left(\frac{\partial}{\partial w_{1j}}\right)^t B(w_{1i}-k_1-l_1 M_1) \\
& \times \sum_{l_2=-\infty}^{\infty} \left(\frac{\partial}{\partial w_{2i}}\right)^u B(w_{2i}-k_2-l_2 M_2) \\
& \times \sum_{l_3=-\infty}^{\infty} \left(\frac{\partial}{\partial w_{3i}}\right)^v B(w_{3i}-k_3-l_3 M_3) \\
& \times \exp\left(2\pi i \left(\frac{m_1 k_1}{M_1}+\frac{m_2 k_2}{M_2}+\frac{m_3 k_3}{M_3}\right)\right)
\end{aligned}
\tag{6-57}
$$

where $\mathbf{d}_{i,l,tuv}$ are the transformed Hermite coefficients obtained from $\mathbf{c}_{i,l,tuv}$ by change of variable and $B(w)$ are B-splines [15]. Equation (6-57) is the approximation to $S_l(\mathbf{m})$ as a three dimensional discrete Fourier transform (3DFFT) times $\lambda(z_1)\lambda(z_2)\lambda(z_3)$. If the function $B(w)$ has finite support, then the structure factors can be calculated in O(N(log(N))) time.

The FFP method relies on the fact that the structure factors can be approximated by rewriting the reciprocal sum such that the structure factor is re-expressed as the 3DFFT evaluated at $\mathbf{m}$ of a Gaussian density $\rho(\mathbf{r}) = \sum_{i=1}^{N} q_i \rho_{2\alpha_0}(\mathbf{r}-\mathbf{r}_i)$ where $\rho_{2\alpha_0}$ is a normalized Gaussian with exponent $2\alpha_0$. This is very similar to the FFT methods used to accelerate structure factor and density map calculations in macromolecular structure determinations.

The efficiency of the methods outlined above has been tested by calculating the intermolecular Coulomb energies and forces for a series of water boxes (64, 128, 256, 512 and 1024) under periodic boundary conditions [15, 62]. The electron density of each monomer is expanded on five sites (atomic positions and bond mid-points) using two standard ABSs, A2 and P1.These sets were used to fit QM density of a single water molecule obtained at the B3LYP/6-31G* level. We have previously shown that the A1 fitted density has an 8% RMS force error with respect to the corresponding ab initio results. In the case of P1, this error is reduced to around 2% [15, 16]. Table 6-1 shows the results for the 5 water boxes using both ABSs (Table 6-7).

## 6.4.    CONCLUSION

As we have seen, Anisotropic Polarizable Molecular Mechanics (APMM) procedures such as SIBFA or GEM are more complex than usual classical approaches.

*Table 6-7.* Timing (in seconds) of the calculated water boxes using full Ewald, PME and FFP for two different accuracies. All calculations were performed on a single 3.g GHz Xeon processor[a]

| | RMS force = $10^{-3}$ | | | | | |
|---|---|---|---|---|---|---|
| | A1 | | | P1 | | |
| | Ewald | PME | FFP | Ewald | PME | FFP |
| 64 | 0.365 | 0.106 | 0.142 | 2.321 | 0.310 | 0.399 |
| 128 | 1.336 | 0.218 | 0.271 | 7.706 | 0.591 | 0.832 |
| 256 | 5.239 | 0.387 | 0.528 | 35.178 | 1.186 | 1.920 |
| 512 | 17.881 | 0.837 | 1.100 | 119.863 | 2.549 | 3.825 |
| 1024 | 71.513 | 1.701 | 2.236 | 486.384 | 4.953 | 6.890 |
| | RMS force = $10^{-4}$ | | | | | |
| | A1 | | | P1 | | |
| | Ewald | PME | FFP | Ewald | PME | FFP |
| 64 | 0.520 | 0.144 | 0.274 | 3.858 | 0.478 | 0.688 |
| 128 | 1.869 | 0.287 | 0.380 | 11.056 | 0.923 | 1.576 |
| 256 | 7.256 | 0.517 | 0.846 | 49.107 | 1.805 | 2.736 |
| 512 | 25.511 | 1.104 | 1.534 | 183.487 | 3.794 | 5.684 |
| 1024 | 108.158 | 2.249 | 3.152 | 714.644 | 6.947 | 11.307 |

[a] Additional speedups can be gained by reducing the size of the mesh for the sampling of the Gaussians by using the Gaussian split Ewald approach [67].

However, their parametrization is performed upon first principle energy decomposition schemes therefore lies on solid ground as any physical ingredients of the interaction can be added. We have seen that despite their variety, EDA schemes present significant convergence and can be easy used to calibrate MM approaches, especially as liner scaling techniques will allow performing large reference computations. In addition, methods able to unravel local electrostatic properties such as the ELF based DEMEP approach should help force field developers to build more realistic models. To conclude, we have seen that despite their different formulation, SIBFA and GEM share a common philosophy. That way, the GEM continuous electrostatic model will be used to replace SIBFA's distributed multipoles to produce a multiscale SIBFA-GEM approach. It will use the newly developed density based Periodic boundary conditions techniques giving access to an N.log(N) evaluation of integrals, a key issue to perform fast and accurate polarizable molecular dynamics simulations. As perspectives, it is important to point out that such APMM approaches should clearly be an asset for hybrid Quantum Mechanics/Molecular Mechanics (QM/MM) computations [2, 68, 69] as they embody short-range electrostatics and full induction. They will also help improving classical force fields by performing higher level reference calculations on relevant system [69].

## ACKNOWLEDGMENTS

## REFERENCES

1. Yoda T, Sugita Y, Okamoto Y (2004) Secondary-structure preferences of force fields for proteins evaluated by generalized-ensemble simulations. Chem Phys 307:269
2. Gresh N, Cisneros GA, Darden TA, Piquemal J-P (2007) Anisotropic, polarizable molecular mechanics studies of inter-, intra-molecular interactions, and ligand-macromolecule complexes. A bottom-up strategy. J Chem Theory Comput 3:1960
3. Kitaura K, Morokuma K (1976) A new energy decomposition scheme for molecular interactions within the Hartree-Fock approximation. Int J Quantum Chem 10:325
4. Bagus PS, Hermann K, Bauschlicher CW, Jr (1984) A new analysis of charge transfer and polarization for ligand–metal bonding: Model studies of $Al_4CO$ and $Al_4NH_3$. J Chem Phys 80:4378
5. Bagus PS, Illas F (1992) Decomposition of the chemisorption bond by constrained variations: order of the variations and construction of the variational spaces. J Chem Phys 96:8962
6. Piquemal J-P, Marquez A, Parisel O, Giessner-Prettre C (2005) A CSOV Study of the difference between HF and DFT Intermolecular Interaction Energy Values: the importance of the charge transfer contribution. J Comput Chem 26:1052
7. Stevens WJ, Fink WH (1987) Frozen fragment reduced variational space analysis of hydrogen bonding interactions. Application to water dimer. Chem Phys Lett 139:15
8. Ziegler T, Rauk A (1979) CO, CS, N2, PF3 and CNCH3 as donors and acceptors. A theoretical study by the Hartree-Fock-Slater transition state method. Inorg Chem 18:1755;
9. Bickelhaupt FM, Baerends EJ (2000) In: Lipkowitz KB, Boyd DB (ed) Rev Comp Chem 15(1), Wiley-VCH, New York
10. Jeziorski B, Moszynski R, Szalewicz K (1994) Perturbation theory approach to intermolecular potential energy surfaces of van der Waals complexes. Chem Rev 94:1887
11. Reinhardt P, Piquemal J-P, Savin A (2008) Fragment-localized Kohn-Sham orbitals via a Singles-CI procedure and application to local properties and intermolecular energy decomposition analysis. J Chem Theory Comput 4:2020
12. Pilmé J, Piquemal J-P (2008) Advancing beyond charge analysis using the electronic localization function: Chemically intuitive distribution of electrostatic moments. J Comput Chem 29:1440
13. Gresh N, Claverie P, Pullman A (1984) Theoretical studies of molecular conformation. Derivation of an additive procedure for the computation of intramolecular interaction energies. Comparison with ab initio SCF computations. Theor Chim 66(1):1–20
14. Piquemal J-P, Cisneros GA, Reinhardt P, Gresh N, Darden TA (2006) Towards a force field based on density fitting. J Chem Phys 24:104101
15. Cisneros GA, Piquemal J-P, Darden TA (2006) Generalization of the Gaussian electrostatic model: extension to arbitrary angular momentum, distributed multipoles and speedup with reciprocal space methods. J Chem Phys 125:184101
16. Cisneros GA, Elking D, Piquemal J-P, Darden TA (2007) Numerical fitting of molecular properties to Hermite Gaussians. J Phys Chem A 111:12049

17. Gordon MS, Jensen JH (1998) Wavefunctions and chemical bonding. In: Schleyer PvR, Allinger NL, Clark T, Gasteiger J, Kollman PA, Schaefer III, HF, Schreiner PR (eds) The encyclopedia of computational chemistry, 5:3198 John Wiley and Sons, Chichester

18. Piquemal J-P, Chevreau H, Gresh N (2007) Towards a separate reproduction of the contributions to the Hartree-Fock and DFT intermolecular interaction energies by polarizable molecular mechanics with the SIBFA potential.J Chem Theory Comput 3:824

19. Piquemal J-P, Chelli R, Procacci P, Gresh N (2007) Key role of the polarization anisotropy of water in modeling classical polarizable force fields. J Phys Chem A 111:8170

20. Claverie P, PhD Thesis, CNRS library number A.O. 8214, Paris (1973)

21. Claverie P (1976) Localization and delocalization in quantum chemistry. In: Chalvet O, Daudel R, Diner S, Malrieu JP (eds) 21:127, Reidel, Dordrecht

22. Hess O, Caffarel M, Huiszoon C, Claverie P (1990) Second-order exchange effects in intermolecular interactions. The water dimer. J Chem Phys 92:6049

23. Khaliullin RZ, Cobar EA, Lochan RC, Bell AT, Head-Gordon M (2007) Unravelling the origin of intermolecular interactions using absolutely localized molecular orbitals. J Phys Chem A 111:8753

24. Daudey J-P (1974) Direct determination of localized SCF orbitals. Chem Phys Lett 24:574

25. Ochsenfeld C, Kussmann J, Lambrecht DS (2007) Linear-scaling methods in quantum chemistry. In: Lipkowitz KB, Cundari TR (ed) Rev Comp Chem 23(1), VCH, New York

26. Rubio J, Povill A, Malrieu J-P, Reinhardt PJ (1997) Direct determination of localized Hartree-Fock orbitals as a step toward Nscaling procedures. J Chem Phys 107:10044

27. Pople JA, Gill PMW, Johnson BG (1992) Kohn-Sham density-functional theory within a finite basis set. Chem Phys Lett 199:557

28. Becke AD, Edgecombe KE (1990) A simple measure of electron localization in atomic and molecular systems. J Chem Phys 92:5397

29. Savin A, Flad OJ, Andersen J, Preuss H,Von Schering HG, Angew (1992) Electron localization in solid-state. Structures for the elements: the diamond. Chem Int 31:187

30. Silvi B, Savin A (1994) Classification of chemical bonds based on topological analysis of electron localization functions. Nature 371:683

31. Piquemal J-P, Pilmé J, Parisel O, Gérard H, Fourré I, Bergès J, Gourlaouen C, de la Lande A, van Severen MC, Silvi B (2008) What can be learnt on biological or biomimetic systems with the topological analysis of the electron localization function? Int J Quant Chem 108:1951

32. Popelier PLA (2000) Atoms in molecules: An introduction, Prentice-Hall, Harlow, UK

33. Bader RFW (1990) Atoms in molecules: A quantum theory, Oxford University Press, Oxford

34. Vigné-Maeder F, Claverie P (1988) The exact multicenter multipolar part of a molecular charge distribution and its simplified representations. J Chem Phys 88:4934

35. Garmer DR, Stevens WJ (1989) Transferability of molecular distributed polarizabilities from a simple localized orbital based method. J Phys Chem 93:8263

36. Piquemal J-P int. J. Quant Chem. (2010) submitted (Theobio 09, special issue)

37. Piquemal J-P, Gresh N, Giessner-Prettre C (2003) Improved formulas for the calculation of the electrostatic contribution to intermolecular interaction energy from multipolar expansion of the electronic distribution. J Phys Chem A 107:10353

38. Murrel JN, Teixeira D (1970) J Mol Phys 19:521

39. Williams DR, Schaad LJ, Murrel JN (1967) J Chem Phys 47:4916

40. Piquemal J-P (2004) Evaluation of molecular interactions in bioinorganic systems: from ab initio computations to polarizable force fields. PhD Thesis, UPMC number 2004PA066267, Université Pierre et Marie Curie, Paris, France

41. Gresh N, Piquemal J-P, Krauss M (2005) Representation of Zn(II) complexes in polarizable molecular mechanics. Further refinements of the electrostatic and short-range contribution of the intermolecular interaction energy. Comparisons with parallel ab initio computations. J Comput Chem 26:1113

42. Mulliken R, Rieke C, Orloff D, Orloff H (1948) Formulas and numerical tables for overlap integrals. J Chem Phys 17:1248

43. Roothan CJ (1951) A study of two-center integrals useful in calculations on molecular structure. J Chem Phys 19:1445

44. Foster JM, Boys SF (1960) Canonical configurational interaction procedure. Rev Mod Phys 32:300

45. Gresh N (1995) Energetics of $Zn^{2+}$ binding to a series of biologically relecant ligands: A molecular mechanics investigation grounded on ab initio self-consistent field supermolecular computations. J Comput Chem 16:856

46. Murrel JN, Randic M, Williams DR (1966) Proc Roy Soc A 284:566 (London)

47. Gresh N, Claverie P, Pullman A (1982) Computations of intermolecular interactions: Expansion of a charge-transfer energy contribution in the framework of an additive procedure. Applications to hydrogen-bonded systems. Int J Quant Chem 22:199

48. Gresh N, Claverie P, Pullman A (1986) Intermolecular interactions: Elaboration on an additive procedure including an explicit charge-transfer contribution. Int J Quant Chem 29:101

49. Creuzet S, Langlet J, Gresh N (1991) J Chim Phys 88:2399

50. Piquemal J-P, Gresh N (to be published)

51. Piquemal J-P, Williams-Hubbard B, Fey N, Deeth RJ, Gresh N, Giessner-Prettre C (2003) Inclusion of the ligand field contribution in a polarizable molecular mechanics: SIBFA LF. J Comput Chem 24:1963

52. Schäffer CE, Jørgensen CK (1965) The angular overlap model, an attempt to revive the ligand field approaches. Mol Phys 9:401

53. Deeth RJ (2001) The ligand field molecular mechanics model and the stereoelectronic effects of d and s electrons. Coord Chem Rev 212:11

54. Woodley M, Battle PD, Catlow CRA, Gale JD (2001) Development of a new interatomic potential for the modeling of ligand field effects. J Phys Chem B 105:6824

55. Boys SF, Shavit I (1959) A fundamental calculation of the energy surface for the system of three hydrogens atoms, NTIS, Springfield, VA, AD212985

56. Dunlap BI, Connolly JWD, Sabin JR (1979) On first row diatomic molecules and local density models. J Chem Phys 71:4993

57. McMurchie LE, Davidson ER (1978) One- and two-electron integrals over cartesian gaussian functions. J Comp Phys 26:218

58. Cisneros GA, Piquemal J-P, Darden TA (2005) Intermolecular electrostatic energies using density fitting. J Chem Phys 123:044109

59. Wheatley RJ, Price S (1990) An overlap model for estimating the anisotropy of repulsion. Mol Phys 69:507

60. Piquemal J-P, Perera L, Cisneros GA, Ren P, Pedersen LG, Darden TA (2006) Towards accurate solvation dynamics of divalent cations in water using the polarizable Amoeba force field: from energetics to structure. J Chem Phys 125:054511

61. Ren P, Ponder JW (2003) Temperature and pressure dependence of AMOEBA water model. J Phys Chem B 107:5933

62. Darden TA (2008) Dual bases in crystallographic computing. In International Tables of Crystallography, vol. B, 3rd edn. Kluwer Academic Publishers, Dordrecht, The Netherlands, (in press)

63. Fusti-Molnar L, Pulay P (2002) The Fourier transform Coulomb method: Efficient and accurate calculation of the Coulomb operator in a Gaussian basis. J Chem Phys 117:7827

64. Sagui C, Pedersen LG, Darden TA (2004) Towards an accurate representation of electrostatics in classical force fields: Efficient implementation of multipolar interactions in biomolecular simulations. J Chem Phys 120:73–87
65. Essmann U, Perera L, Berkowitz ML, Darden T, Lee H, Pedersen LG (1995) A smooth particle mesh Ewald method. J Chem Phys 103:8577
66. York D, Yang W (1994) The fast Fourier Poisson method for calculating Ewal sums. J Chem Phys 101:3298
67. Shan YB, Klepeis JL, Eastwood MP, Dror RO, Shaw DE (2005) Gaussian split Ewald: A fast Ewald mesh method for molecular simulation. J Chem Phys 122:054101
68. Cisneros GA, Piquemal J-P, Darden TA (2006) QM/MM electrostatic embedding with continuous and discrete functions. J Phys Chem B 110:13682
69. Cisneros GA, Na-Im Tholander S, Elking D, Darden TA, Parisel O, Piquemal J-P (2008) Simple formulas for improved point-charge electrostatics in classical force fields and hybrid Quantum Mechanical/Molecular Mechanical embedding. Int J Quant Chem 108:1905

# CHAPTER 7

# "MULTI-SCALE" QM/MM METHODS
# WITH SELF-CONSISTENT-CHARGE
# DENSITY-FUNCTIONAL-TIGHT-BINDING (SCC-DFTB)

QIANG CUI[1] AND MARCUS ELSTNER[2]

[1]*Department of Chemistry and Theoretical Chemistry Institute, University of Wisconsin, Madison, 1101 University Ave, Madison, WI 53706, USA, e-mail: cui@chem.wisc.edu*
[2]*Department of Physical and Theoretical Chemistry TU Braunschweig, Hans-Sommer-Straße 10 D-38106 Braunschweig, Germany*

**Abstract:**   Recent QM/MM developments based on the SCC-DFTB method as the QM level are discussed and illustrated using several biological applications. An important theme in these developments concerns the treatment of long-range electrostatics in a "multi-scale" fashion in which the system is partitioned into QM, MM and continuum electrostatic regions. The value of carefully benchmarking the computational model using microscopic *pKa* calculations and the importance of conducting sufficient conformational sampling are emphasized with applications to long-range proton transfer problems

## 7.1.   INTRODUCTION

For chemical reactions occurring in biological systems, entropic contributions to the reactive free energy can be as important as total energy contributions. In many cases, fluctuations in the environment of the active site lead to large variations of reaction energetics, emphasizing the importance of conformational sampling [1–3]. Therefore, in productive QM/MM studies, a focus on the accuracy of the QM description needs to be balanced with computational efficiency to allow for sufficient sampling. In most QM/MM applications to biological systems [4–7], the QM level is typically either Density Functional Theory (DFT) or a semi-empirical (SE) method, although highly correlated methods have also been used for improving the energetics based on single point energy calculations [8]. With DFT methods as the QM level, QM/MM

simulations are often limited to the pico-second regime, which is likely too short for most biological problems except for ultra-fast processes involving electronically excited states [9]. In recent years, however, a promising framework for calculating free energies with DFT based QM/MM methods have been proposed [10–12], where the key idea is to decouple the fluctuation of QM and MM atoms and to sample fluctuation of the environment at the MM level with the frozen QM part treated as effective charges.

Semi-empirical methods are roughly three orders of magnitude faster than DFT due to the approximation and even neglect of specific interaction integrals [13]. As a result, SE methods routinely allow for at least nano seconds of simulations and therefore a meaningful probe of free energies while treating fluctuations of both the QM and MM atoms at equal footing [14, 15]. This is an important advantage of SE methods, which underlines the value of exploring and pushing forward the applicability of SE based QM/MM methods. This is the main topic of discussion in this chapter, although some of the technical insights (e.g., the treatment of van der Waals interactions and long-range electrostatics) are, in principle, applicable to DFT and ab initio based QM/MM methods as well.

Most SE methods are derived from the Hartree-Fock theory by applying various approximations resulting in, for example, the so-called Neglect of Differential Diatomic Overlap (NDDO) type of methods, the most well known being the MNDO, AM1 and PM3 models. Some of the more recent models include the inclusion of orthogonalization corrections in the OMx model [16], the PDDG/PM3 model [17], and the NO-MNDO model [18]. The retained integrals in these SE methods are treated as parameters and derived by either fitting to experimental data or pre-calculation from first principles. Due to these approximations, SE methods usually have an overall lower accuracy than DFT and other ab initio methods, although this may be reversed for specific applications. For example, an interesting option is to develop "specific reaction parametrizations" [19] for a SE method (e.g., PM3), which may provide a very accurate description for the reaction of interest at a level even unmatched by popular DFT methods. However, parameterization of a SRP that works well for condensed phase simulations is not as straightforward as for gas-phase applications and a large number of carefully constructed benchmark calculations are needed [20, 21]. Therefore, it remains an interesting challenge to develop generally robust SE methods that properly balance computational efficiency and accuracy.

We focus on the self-consistent-charge density-functional-tight-binding (SCC-DFTB) method [22], which is an approximation to density functional theory (DFT) and derived from a second order expansion of the DFT total energy expression. In recent years, SCC-DFTB has been successfully applied to a wide range of problems involving structures and dynamics of biomolecules and biocatalysis in several enzymes; for comprehensive reviews see, for example, Refs. [23–27]. With respect to its computational efficiency, SCC-DFTB is comparable to the NDDO type semi-empirical methods, i.e., being 2–3 orders of magnitude faster than DFT(HF) methods (with small to medium-sized basis sets). The speed-up with respect to DFT is achieved without much loss of accuracy in the description of molecular geometries,

while reaction energies and vibrational frequencies are usually less reliable [22, 28]. This is confirmed by two recent thorough studies that evaluated SCC-DFTB for heats of formation, molecular structures etc. on large sets of molecules [29, 30]. While the most sophisticated NDDO methods [13] are slightly superior to SCC-DFTB for heats of formation, the strength of SCC-DFTB is the overall excellent prediction of molecular structures, in particular for larger (bio-) molecular systems [31–34].

Due to their fast computational speed, the SE methods may not be the computational bottleneck for SE based QM/MM simulations when the system is very large in size, such as membrane proteins or large RNA (e.g., the ribosome!) molecules. In these context, there is great interest in pushing forward QM/MM techniques in a "multi-scale" framework [35]. This is important not only for computational efficiency but also for better benefiting from enhanced sampling techniques for regions of most relevance. The most straightforward concept is to treat the reactive fragments with QM, the immediate environment (e.g., within 20–25 Å) with MM, and the rest with continuum electrostatics. Although this scheme has been envisioned many years ago by a number of researchers, a flexible implementation applicable to biomolecules has only been reported in recent years. For example, we have implemented the generalized solvent boundary (GSBP) condition approach of Roux and co-workers [36] in a QM/MM framework [37]. A related but different formulation based on the boundary element approach has been reported by York and co-workers [38]. With simpler QM methods, such a framework has been explored by Warshel and co-workers in their pioneering studies [39].

In the following, we first briefly review the SCC-DFTB method and comment on a few issues related to the implementation of SCC-DFTB/MM, such as the "multi-scale" SCC-DFTB/MM-GSBP protocol. Next, a few specific examples of SCC-DFTB/MM simulations are given. The basic motivation is to highlight a number of issues that might impact either the quantitative or even qualitative nature of the result. We hope that the chapter is particularly instructive to researchers who are relatively new to the field and able to help them carry out meaningful QM/MM simulations.

## 7.2.    QM/MM METHODOLOGY

### 7.2.1.    SCC-DFTB

In DFT the total energy is expressed as a functional of the electron density $\rho$ of the molecular system of interest. The derivation of the SCC-DFTB method starts by choosing a reference density $\rho_0$ as a superposition of densities $\rho_0^\alpha$ of the neutral atoms $\alpha$ constituting the molecular system,

$$\rho_0 = \sum_\alpha \rho_0^\alpha \qquad (7\text{-}1)$$

and a density fluctuation $\delta\rho$, also built up from atomic contributions

$$\delta\rho = \sum_\alpha \delta\rho^\alpha, \tag{7-2}$$

in order to represent the ground state density

$$\rho = \rho_0 + \delta\rho. \tag{7-3}$$

Then, the DFT total energy functional is expanded up to second order with respect to the charge density fluctuations $\delta\rho$ around the reference $\rho_0$ [22] ($\rho_0' = \rho_0(\vec{r}')$, $\int' = \int d\vec{r}'$):

$$E = \sum_i^{occ}\langle\Phi_i|\hat{H}^0|\Phi_i\rangle + \frac{1}{2}\iint' \left(\frac{1}{|\vec{r}-\vec{r}'|} + \frac{\delta^2 E_{xc}}{\delta\rho\,\delta\rho'}\bigg|_{\rho_0}\right)\delta\rho\,\delta\rho'$$

$$-\frac{1}{2}\iint' \frac{\rho_0'\rho_0}{|\vec{r}-\vec{r}'|} + E_{xc}[\rho_0] - \int V_{xc}[\rho_0]\rho_0 + E_{cc} \tag{7-4}$$

$\hat{H}^0 = \hat{H}[\rho_0]$ is the effective Kohn-Sham Hamiltonian evaluated at the reference density $\rho_0$ and the $\Phi_i$ are Kohn-Sham orbitals. $E_{xc}$ and $V_{xc}$ are the exchange-correlation energy and potential, respectively, and $E_{cc}$ is the core-core repulsion energy (an extension up to third order has been presented recently [40, 41]).

The derivation of the DFTB model involves three major approximations of Eq. (7-4):

- The Kohn-Sham orbitals in the first term are expanded in a minimal LCAO basis set $\Phi_i = \sum_\mu c_\mu^i \eta_\mu$. The resulting matrix-elements in the AO basis are subjected to a two-center approximation, which reduces the number of integrals to be evaluated significantly and allows to calculate and tabulate these parameters [42, 43].
- The density fluctuations in the second term are approximated by atomic charge monopoles $\Delta q_\alpha$, which allows to approximate the functional derivatives by an analytical function $\gamma_{\alpha\beta}$ (for a more detailed discussion, see Refs. [22, 40, 41]).
- The remaining four energy contributions depend on the *reference density only*. This is an important observation, which allows to combine these contribution into the so called 'repulsive energy' term $E_{rep}[\rho_0]$, which is treated in a simplified way by approximating it by a sum of two-body potentials [44], $E_{rep}[\rho_0] = \sum_{\alpha\beta} U_{\alpha\beta}(R_{\alpha\beta})$

These approximations result in the final energy of the SCC-DFTB model:

$$E = \sum_{i\mu\nu}^{occ} c_\mu^i c_\nu^i < \eta_\mu|\hat{H}^0|\eta_\nu > + \sum_{\alpha<\beta} U_{\alpha\beta}(R_{\alpha\beta}) + \frac{1}{2}\sum_{\alpha\beta}\Delta q_\alpha\Delta q_\beta\gamma_{\alpha\beta} \tag{7-5}$$

Having these severe approximations in mind, SCC-DFTB performs surprisingly well for many systems of interest, as discussed above. However, it has a lower overall accuracy than DFT or post HF methods. Therefore, applying it to new classes of systems should be only done after careful examination of its performance. This can be done e.g. by conducting reference calculations on smaller model systems with DFT or ab initio methods. A second source of errors is related to some intrinsic problems with the GGA functionals also used in popular DFT methods (SCC-DFTB uses the PBE functional), which are inherited in SCC-DFTB. This concerns the well known GGA problems in describing van der Waals interactions [32], extended conjugate $\pi$ systems [45, 46] or charge transfer excited states [47, 48].

### 7.2.2.    QM/MM Implementations with DFTB

DFTB has been combined with various force field programs, including CHARMM [49], Amber [50, 51], SIGMA [52], TINKER [34] and GROMACS (to be published). All these combined methods have used the same strategy for the interface with respect to the treatment of bonding interactions at the interface and non-bonded Coulomb and van der Waals (vdW) terms. The concepts will not be reviewed in detail here, for this purpose we would like to refer to an excellent and extensive recent review [4].

In principle, one can distinguish between "additive" and "subtractive" QM/MM methods. In both approaches, the system is divided into two (or more) parts, where one part is described by the QM method (e.g. the active site in protein) while the other is treated by the MM method. In the "subtractive" scheme, the MM method is applied to the entire system and the MM description for the QM region is corrected by subtracting the results of a QM and a MM calculation for that region. In this method, the MM method has to treat the QM region as well, although its contribution is "subtracted out". In the "additive" schemes, the QM method is applied to the QM region and the MM method treats the remainder of the system. The two methods in the subsystems are then coupled by bonding and non-bonding (electrostatic and vdW) terms. All QM/MM implementations of DFTB mentioned above have used the additive scheme. A recent implementation into the GAUSSIAN [53] software packages uses the subtractive scheme (to be published). The DFTB QM/MM methodology has been described in several reviews as well [24–27].

### 7.2.2.1.    QM/MM Frontier

When combining QM with MM methods, the partitioning of the system will often intersect a chemical bond. This bond is usually chosen to be a carbon-carbon single bond (whenever possible) and three major coupling methods have been developed, which are referred to as the "link-atom [54]", "pseudo-atom/bond [55]" and "hybrid-orbital [56]" approach, respectively. In the "link atom" approach the open valency at the border is capped by a hydrogen atom, and most DFTB QM/MM implementations are based on this simple scheme [49, 50] or related variations [57]. Recently,

a "hybrid-orbital" approach has been implemented into the DFTB/CHARMM interface [58]. Several systematic comparisons of different approaches indicate that all linking schemes work more or less equally well if the MM charges at the boundary are treated carefully [57, 59, 60].

### 7.2.2.2.    *QM/MM Coulomb Interaction*

A second approximation in the SCC-DFTB/MM implementation is related to the treatment of the Coulomb interaction between the QM and MM regions. First principle methods have to evaluate additional core-electron integrals [61], while SE methods facilitate the coupling using the integral approximations of the respective SE level [54]. The DFTB QM/MM interface has been implemented by making use of the Mulliken charges for the QM region, calculating the Coulomb sum between Mulliken charges from the QM atoms and MM charges. In principle, this interaction should not use the point charge approximation for charges near the QM region, which may lead to a severe overpolarization of the QM region, as frequently discussed for QM/MM implementations using first principles methods [62, 63]. As a consequence, a short range damping should be applied, which, however is neglected in practical SE QM/MM implementations since QM Mulliken charges are usually small in magnitude (for a more detailed discussion, see Ref. [24]). In principle, a more elaborate charge schemes like the CM3 charges [64] may also be superior; however, adopting such charge schemes is more tedious since the charge calculation would have to be considered in both the SCF and parameter fitting procedures.

### 7.2.2.3.    *QM/MM van der Waals (vdW) Interactions*

The vdW interaction between QM and MM regions represents both the short-range repulsion and long-range dispersion interaction between the QM and MM atoms. Since for most force fields, the non-bonding parameters (vdW parameters and charges) are fitted together as a set, directly "borrowing" vdW parameters from existing force fields for the QM atoms may lead to a unbalanced description. Therefore, vdW parameters may have to be refitted when used for the QM/MM interface [65, 66], especially when solvation free energy and solvent structure around the solute are of interest [67]. For *relative* energetic properties (e.g., differential solvation), however, error cancellation often leads to only mild dependence of the result on the choice of vdW parameters.

### 7.2.2.4.    *The Choice of Boundary Condition and Treatment of Electrostatics*

As in classical simulations of biomolecules, there are two general frameworks for setting up QM/MM simulations for a biological system: periodic boundary condition (PBC) and finite-size boundary condition (FBC). When the system of interest is small (∼200–300 amino acids), PBC is well suited because the entire system can be completely solvated and therefore structural fluctuations ranging from the residue level to domain scale can potentially be treated at equal footing, within the limit

of the sampling time-scale. To avoid periodic artifacts [68], however, large simulation cells are needed, which makes the PBC set-up computationally intensive even if the QM region is small. A key issue in PBC simulations is the treatment of long-range electrostatics, which has a rich literature for classical simulations [69, 70]. For QM/MM simulations, a number of groups have implemented Ewald and Particle-Mesh-Ewald algorithms for either semi-empirical QM [71–73] or ab initio QM methods [74]. Therefore, especially with semi-empirical level QM, it has become very feasible to carry out PBC simulations with QM/MM potentials.

If the process of interest is highly localized, FBC is preferred because only a modest number of atoms (e.g., a spherical region of 25 Å radius) need to be simulated explicitly. The reliability of FBC based simulations, however, depends on the details of the set-up procedure, such as the treatment of interface and interactions between the explicitly simulated and simplified regions. In the context of QM/MM simulations, a popular scheme is the stochastic boundary condition (SBC) developed by Brooks and co-workers [75] originally for classical simulations. In the most recently refined version [76, 77], the atoms in a spherical region are treated explicitly and surrounded by a buffer region of typically 2–4 Å thick; to consider the screening effect due to the bulk solvent, partial charges of charged residues are scaled based on Poisson-Boltzmann (PB) calculations in the vacuum and solution. The subsequent QM/MM simulations are done using the scaled partial charges with the system in vacuum, and the effects of the full charges can be considered by PB "post-correction" calculations using a collection of snapshots. A number of applications have shown that for modest-size systems, SBC works rather well provided that the PB corrections are done carefully [78, 79].

A potential problem of SBC is that the charge-scaling factors are calculated for a single structure, which may not be most appropriate for all structures sampled during the simulation. In addition, for reactions that involve significant variations in charge distributions, the PB "post-corrections" tend to involve large numbers of opposite signs and therefore subject to significant numerical noises. Motivated by these limitations, we have pursued an alternative "multi-scale" protocol based on the Generalized Solvent Boundary Potential (GSBP) approach [36] developed by Roux and co-workers for classical simulations.

Similar to the standard SBC, GSBP partitions the system into inner and outer regions and the effects of the outer region on the inner, reaction region are represented implicitly within the total effective potential (potential of mean force) [36],

$$W_{\text{GSBP}} = U^{(\text{ii})} + U^{(\text{io})}_{\text{int}} + U^{(\text{io})}_{\text{LJ}} + \Delta W_{\text{np}} + \Delta W^{(\text{io})}_{\text{elec}} + \Delta W^{(\text{ii})}_{\text{elec}}, \qquad (7\text{-}6)$$

where $U^{(\text{ii})}$ is the complete inner–inner potential energy, $U^{(\text{io})}_{\text{int}}$ and $U^{(\text{io})}_{\text{LJ}}$ are the inner–outer internal (bonds, angles, and dihedrals) and Lennard–Jones potential energies, respectively, and $\Delta W_{\text{np}}$ is the non-polar confining potential. The last two terms in Eq. (7-6) are the core of GSBP, representing the long-range electrostatic interaction between the outer and inner regions. The contribution from distant protein

charges (screened by the bulk solvent) in the outer region, $\Delta W_{\text{elec}}^{(\text{io})}$, is represented in terms of the corresponding electrostatic potential in the inner region, $\phi_s^{(o)}(\mathbf{r}_\alpha)$,

$$\Delta W_{\text{elec}}^{(\text{io})} = \sum_{\alpha \in \text{inner}} q_\alpha \phi_s^{(o)}(\mathbf{r}_\alpha) \tag{7-7}$$

The dielectric effect on the interactions among inner region atoms is represented through a reaction field term,

$$\Delta W_{\text{elec}}^{(\text{ii})} = \frac{1}{2} \sum_{mn} Q_m M_{mn} Q_n \tag{7-8}$$

where $\mathbf{M}$ and $Q$ are the generalized reaction field matrix and generalized multipole moments, respectively, in a basis set expansion [36].

The advantage of the GSBP method lies in its ability to include these contributions explicitly while sampling configurational space of the reaction region during a simulation at minimal additional cost. The static field potential, $\phi_s^{(o)}(\mathbf{r})$, and the generalized reaction field matrix $\mathbf{M}$ are computed only once based on PB calculations and stored for subsequent simulations. The only quantities that need to be updated during the simulation are the generalized multipole moments, $Q_n$,

$$Q_n = \sum_{\alpha \in \text{inner}} q_\alpha b_n(\mathbf{r}_\alpha) \tag{7-9}$$

where $b_n(\mathbf{r}_\alpha)$ is the $n$th basis function at nuclear position $\mathbf{r}_\alpha$.

As described in Ref. [37], the implementation of GSBP in a combined QM/MM framework requires three additional terms,

$$\frac{1}{2} \sum_{mn} Q_m^{\text{QM}} M_{mn} Q_n^{\text{QM}} + \sum_{mn} Q_m^{\text{QM}} M_{mn}(Q_n^{\text{MM}} - Q_n^{\text{EX}}) + \int d\mathbf{r}\, \rho^{\text{QM}}(\mathbf{r}) \phi_s^{(o)}(\mathbf{r}), \tag{7-10}$$

corresponding to the QM–QM and QM–MM (corrected for exclusions due to link host schemes [37, 57]) reaction field, and the QM-static field terms, respectively. For the GSBP combined with SCC-DFTB, these terms take on a simple form because $\rho^{\text{QM}}(\mathbf{r})$ is expressed in terms of Mulliken charges, [22] $\rho^{\text{QM}}(\mathbf{r}) = \sum_{A \in \text{QM}} \Delta q^A \delta(\mathbf{r} - \mathbf{R_A})$, resulting in expressions,

$$Q_m^{\text{QM}} = \sum_{A \in \text{QM}} \Delta q^A b_m(\mathbf{R_A}), \tag{7-11}$$

for the generalized QM multipoles and,

$$\Delta W_{\text{elec}}^{\text{QM(io)}} = \sum_{A \in \text{QM}} \Delta q^A \phi_s^{(o)}(\mathbf{R}_A), \tag{7-12}$$

for the interaction of the QM region with the static field due to the outer region. These terms are included in the SCC-DFTB matrix elements during the SCF iteration [37].

Although the formulation of GSBP is self-consistent, the validity of the approach depends on many factors especially the size of the inner region and the choice of the dielectric "constant" for the outer region. Therefore, for any specific application, the simulation protocol has to be carefully tested using relevant benchmarks such as $pK_a$ of key residues (see examples below in Sections 7.3.1 and 7.3.2).

### 7.2.2.5. *Free Energy Simulations Using QM/MM Potentials*

Due to its low computational cost, SCC-DFTB makes it possible to perform QM/MM simulations with extensive sampling, which is crucial for condensed phase systems. In general, two types of free energy simulations are of interest: potential of mean force (PMF) and free energy perturbation (FEP) simulations. PMF is relevant for studying the free energy profile for a chemical reaction and the result can be related to the rate constant through rate theories [80]. In addition to adequate sampling, the value of the computed PMF depends critically on the choice of the reaction coordinate(s). For chemical reactions that involve a large number of groups, such as long-range proton transfers, for example, the design of a proper reaction coordinate is important for meaningful simulation of the mechanism (see example below in Section. 7.3.3).

Although FEP is mostly useful for binding type of simulations rather than chemical reactions, it can be valuable for reduction potential and $pK_a$ calculations, which are of interest from many perspectives. For example, prediction of reliable $pK_a$ values of key groups can be used as a criterion for establishing a reliable microscopic model for complex systems. Technically, FEP calculation with QM/MM potentials is complicated by the fact that QM potentials are non-seperable [78]. When the species subject to "perturbation" $(A \rightarrow B)$ differ mainly in electronic structure but similar in nuclear connectivity (e.g., an oxidation-reduction pair), we find it is beneficial to use the same set of nuclear geometry for the two states [78], i.e., the coupling potential function has the form,

$$U^\lambda(\lambda) = (1 - \lambda)U_A(\mathbf{X}_A, \mathbf{X}_C) + \lambda U_B(\mathbf{X}_A, \mathbf{X}_C) + U_{CC}(\mathbf{X}_C) \tag{7-13}$$

We note that using the same set of geometry for the two states, per se, does not introduce any approximation due to the state character of free energy (i.e., as $\lambda$ evolves from 0 to 1, the geometry evolves from that for $A$ to $B$). In practice, error may occur when SHAKE [81] is used to constrain bond distances involving hydrogen to be the same for all $\lambda$ values, although the magnitude is expected to be very small for all practical cases.

This "dual-topology-single-coordinate" (DTSC) protocol apparently works only when $A$ and $B$ have very similar (but not identical) nuclear geometries. In the special case of a small difference between the two "solute" states, e.g., $AH$ and $A^-$ for $pK_a$ problems, special thermodynamic cycles can be constructed to take advantage of the DTSC protocol [73, 82]. When $A$ and $B$ are very different (e.g., for relative binding free energy simulations), one could either use the conventional dual-topology-dual-coordinate potential or use the elegant "chaperone" approach proposed by Yang and co-workers [83].

## 7.3.     EXAMPLES AND DISCUSSIONS

In this section, we use examples to illustrate several key issues that may significantly impact the reliability of QM/MM simulations of biological systems. In addition, we also discuss calculations that are useful for validating the simulation protocols in realistic applications.

### 7.3.1.     The Importance of Consistent Electrostatic Treatment

As discussed in many previous studies of biomolecules, the treatment of electrostatic interactions is an important issue [69, 70, 84]. What is less widely appreciated in the QM/MM community, however, is that a balanced treatment of QM–MM electrostatics and MM–MM electrostatics is also an important issue. In many implementations, QM–MM electrostatic interactions are treated without any cut-off, in part because the computational cost is often negligible compared to the QM calculation itself. For MM–MM interactions, however, a cut-off scheme is often used, especially for finite-sphere type of boundary conditions. This imbalanced electrostatic treatment may cause over-polarization of the MM region, as was first discussed in the context of classical simulations with different cut-off values applied to solute-solvent and solvent–solvent interactions [85]. For QM/MM simulations with only energy minimizations, the effect of over-polarization may not be large, which is perhaps why the issue has not been emphasized in the past. As MD simulations with QM/MM potential becomes more prevalent, this issue should be emphasized.

As an illustration, we briefly discuss the SCC-DFTB/MM simulations of carbonic anhydrase II (CAII), which is a zinc-enzyme that catalyzes the interconversion of $CO_2$ and $HCO_3^-$ [86]. The rate-limiting step of the catalytic cycle is a proton transfer between a zinc-bound water/hydroxide and the neutral/protonated His64 residue close to the protein/solvent interface. Since this proton transfer spans at least 8–10 Å depending on the orientation of the His 64 sidechain ("in" vs. "out", both observed in the X-ray study [87]), the transfer is believed to be mediated by the water molecules in the active site (see Figure 7-1). To carry out meaningful simulations for the proton transfer in CAII, therefore, it is crucial to be able to describe the water structure in the active site and the sidechain flexibility of His 64 in a satisfactory manner.

*Figure 7-1.* The active site of CAII rendered from the crystal structure (PDB ID: 2CBA [87]). All dotted lines correspond to hydrogen-bonding interactions with distances ≤3.5 Å. E117 and E106 are in close proximity to H119, and E106 also interacts with T199 through the presumed hydroxyl proton of T199 (not shown for clarity). H64 is resolved to partially occupy both the "in" and "out" rotameric states

As shown in Figure 7-2, rather different orientations of the His 64 sidechain and water distributions are accessible in SCC-DFTB/MM simulations with different electrostatic treatments; in all simulations, no cutoff is used for QM–MM interactions. When the "default" force-shift cutoff scheme is used for interaction among MM atoms, the His 64 side chain consistently flipped to the "out" position (Figure 7-2(a)) in the very early stage of the trajectory in all COHH [88] simulations (protonated His 64, zinc-bound hydroxide) [37]. This behavior is likely due to the unbalanced QM–MM and MM–MM interactions, which overestimate the repulsion between the protonated His 64 and zinc-hydroxide (both groups bear a charge of +1). For example, the favorable interaction between Glu 106 in the active site (Figure 7-1) and the protonated His 64, which counteracts the repulsion between His 64 and the zinc site, is truncated in the MM interactions. When the MM interactions are treated with the extended electrostatics scheme in which interactions beyond 12 Å are computed using multipolar expansion, both "in" and "out" configurations were sampled more evenly during the simulations [37]. However, if the bulk solvation effect is not taken into account, as done in many conventional stochastic boundary condition based simulations, the His 64 side chain was observed [37] to reach a region very close (∼5.0 Å) compared to the value ∼8.0 Å in the X-ray data) to the zinc atom, as shown in Figure 7-2(a). Indeed, without the proper treatment of the bulk solvation, the attraction between the protonated His 64 and active site residues is overestimated (i.e., opposite to the force-shift cutoff simulations). In the GSBP based SCC-DFTB/MM simulations, since QM–MM and MM–MM electrostatics

*Figure 7-2*. Properties of CAII active site in the COHH state (zinc-bound hydroxide and protonated His 64). (**a**) Superposition of a few key residues from two stochastic boundary SCC-DFTB/MM simulations with the X-ray structure [87] (colored based on atom-types); the two sets of simulations did not have any cut-off for the electrostatic interactions between SCC-DFTB and MM atoms but used different treatments for the electrostatic interactions among MM atoms: group-based extended electrostatics (in *yellow*) and atom-based force-shift cut-off (in *green*). Extended electrostatics simulations sampled configurations with the protonated His 64 too close to the zinc moiety while force-shift simulations consistently sampled the "out" configuration of His 64 in multiple trajectories. (**b**) Statistics for productive water-bridges (only from two and four shown here) between the zinc bound water and His 64 with different electrostatics protocols

were treated in a balanced manner and the effect of bulk solvation (as well as protein atoms outside of the explicitly sampled region) is considered, the system was well behaved and no artifactual behavior was observed.

Another property relevant to the current discussion is the distribution of water in the active site. Specifically, we characterize the population of various "water wires" connecting the zinc-bound water/hydroxide and His 64 found in the SCC-DFTB/MM simulations. These wires were identified following a definition of hydrogen-bond in terms of both distance (O—O $< 3.5$ Å) and angle (O—H—O $\geq 140°$) and care

was taken to classify the wire as productive if a proton can be successfully passed from donor to acceptor through the wire and unproductive otherwise (e.g., one water molecule serves as a double-donor or double-acceptor); when more than one productive water bridge was present in a frame, only the shortest one was counted. As shown in Figure 7-2(b), the agreement between Ewald and GSBP based SCC-DFTB/MM results was overall very good, while the agreement between GSBP and stochastic boundary simulations depended on the protonation state and orientation of the His 64 sidechain. For example, for the CHOH state (zinc-bound water and neutral His 64) with His 64 adopting the "in" conformation, extended electrostatics and force-shift cutoff simulations distinctively favored four- and two-water bridges, respectively; GSBP simulations, on the other hand, produced rather even distributions of two- to six-water bridges. The Ewald simulations produced results similar to the GSBP calculations, although a slight preference over four-water bridge was seen (Figure 7-2(b)).

The problem with the cutoff simulation was also apparent when analyzing the number of water molecules in the active site; within 10 Å from the zinc ion, the average number of water molecules in the CHOH simulations (His 64 "in") is 14.2, 14.7, 15.1 and 23.1 for the GSBP, Ewald, extended electrostatics and force-shift cutoff protocols, respectively [27]. This issue of water "flooding" the active site when imbalanced QM-MM and MM-MM interactions are used has also been found in the studies of several systems in our group [89, 90]. For example, in the simulation of $\beta$-lactamase, the substrate dissociated from the active site in a SCC-DFTB/MM simulation using force-shift cutoff for MM–MM interactions as a result of active-site flooding; with a balanced electrostatics treatment in the QM/MM-GSBP framework, the active-site remained stable during nanoseconds of simulations [90].

It is worth noting that although the GSBP protocol provided results in favorable agreement with Ewald simulations in the above discussions, the validity of GSBP based simulations depends on several key factors. First, the fact that the "outer region" protein atoms are fixed may suppress the fluctuation of atoms in the "inner region" because collective motions, which make significant contributions to thermal fluctuations, are removed. As shown in Figure 7-3, with a 20 Å-inner-region simulations for CAII [91], even the largest RMSF is smaller than 0.5 Å, which is considerably lower compared to results from Ewald simulations or experimental B-factors. Upon closer inspection, the atomic fluctuations close to the center of the sphere (zinc ion in CAII simulations) are reproduced in GSBP calculations but the damping effect starts to increase steeply when the buffer region is approached. For residues within 13.5 Å from the zinc, for example, the RMSD between the RMSFs from the 20 Å-inner-region GSBP simulations and those from the Ewald simulations is 0.11 Å; the RMSD then quickly shoots up beyond that region. The RMSDs of atomic RMSF in the 25 Å-inner-region GSBP simulations are generally smaller than those from the 20 Å set-up; to reach the same RMSD of 0.11 Å, for example, the region extends to 17 Å from the zinc ion. Similarly, the diffusion constant of water molecules in the active site is described well in GSBP simulations for those very close to the center of the sphere; for water molecules close to the inner/outer boundary, substantial decrease

*Figure 7-3.* Active site properties of CAII from SCC-DFTB/MM-GSBP simulations [91]. (**a**) The root mean square differences between the RMSFs calculated from GSBP simulations (WT-20 and WT-25 have an inner radius of 20 and 25 Å respectively) and those from Ewald simulation, for atoms within a certain distance from the zinc, plotted as functions of distance from the zinc ion; that the center of the sphere in GSBP simulations is the position of the zinc ion in the starting (crystal) structure. (**b**) The diffusion constant for TIP3P water molecules as a function of the distance from the zinc ion in different simulations

in the diffusion constant as compared to the Ewald simulations was observed (Figure 7-3). Therefore, it is important to keep in mind that the GSBP protocol is best suited for describing localized processes and the proper size of the "inner region" needs to be established with careful benchmark calculations for the specific system of interest.

### 7.3.2.    $pK_a$ Calculations as Crucial Benchmark

Although water structure and sidechain flexibilities are useful gauges for the simulation protocol, more quantitative measures are needed for reliable QM/MM simulations of enzyme systems. In this regards, we have found that reduction potential [78] and $pK_a$ [73, 91] calculations are particularly useful benchmark calculations because the results are likely very sensitive to the simulation details.

With the QM/MM potential and the free energy perturbation (thermodynamic integration) approach discussed above, it is in principle possible to compute absolute $pK_a$ values. In practice, however, this is very challenging because of uncertainties in the proton solvation free energy and also the need of predicting highly accurate gas-phase proton affinity. With a series of amino acid sidechain analogues, we showed [73] that reliable absolute $pK_a$ values are obtained with SCC-DFTB/MM simulations only when highly-correlated ab initio methods (e.g., CCSD) are used to correct for gas-phase proton affinity and ab initio/*DFT* QM methods are used to correct for QM–MM interactions. For the purpose of validating the simulation protocol, such as

the GSBP set-up, however, it is sufficient to compute $pK_a$ shift, i.e., the change in a group's $pK_a$ going from bulk solution to the protein, as commonly done in MM or continuum electrostatics calculations.

It is worth emphasizing that when the titritable group is in the interior of the protein, even if the area is rather polar (such as the active-site of the CAII, which is full of water molecules) the structural response of the protein/solvent during "in silico" titration can be substantial. This means that long simulations or enhanced sampling techniques are crucial for obtaining reliable $pK_a$ values. As an example, the results of the free energy derivatives and statistical analyses of the data from two 20 Å-inner-region simulations for the E106Q mutant of CAII (denoted as "E106Q-20", the group under titration is the zinc-bound water) are shown in Table 7-1. It is clear that while the free energy derivatives converge with *statistical* errors ~1 kcal/mol, there are differences between the two independent runs on the order of 10(!) kcal/mol for some λ windows (compare Columns 4 and 7 of Table 7-1). Therefore, the free energy derivatives for these runs appear to have equilibrated to sample different regions of the configuration space. Inspection of the trajectories suggests that the difference is likely due to the different orientations of the Thr 199 sidechain sampled in separate simulations, which leads to substantial variation in the interaction between Thr 199 and the zinc-bound water and therefore change in the free energy derivatives.

The $pK_a$ calculations also illustrate the importance of carefully considering the treatment of the QM/MM frontier. As discussed in several previous studies [57, 59, 60, 62], the treatment of the QM/MM frontier is particularly important for reactions where the net charge of the QM region changes significantly, as in the titration process. As shown in Table 7-2, re-evaluating the free energy derivatives with the more robust "link-host-group" exclusion for the three zinc-bound

*Table 7-1.* Representative results from statistical analyses used to determine the values and statistical errors of $\partial \Delta G_{\mathrm{CH}(D)\mathrm{OH}}/\partial \lambda$ for the $pK_a$ calculation of the zinc-bound water in CAII[a]

| | SET 1 | | | SET 2 | | |
|---|---|---|---|---|---|---|
| λ | Length[b] | Block size[c] | $\partial \Delta G_{\mathrm{CH}(D)\mathrm{OH}}/\partial \lambda$[d] | Length[b] | Block size[c] | $\partial \Delta G_{\mathrm{CH}(D)\mathrm{OH}}/\partial \lambda$[d] |
| 0.00 | 1.4 (0.3) | 13 (86) | 212.6 (1.0) | 1.3 (0.6) | 15 (51) | 211.1 (1.1) |
| 0.25 | 1.1 (0.6) | 5 (89) | 184.0 (0.7) | 1.1 (0.8) | 6 (44) | 176.1 (1.1) |
| 0.50 | 1.4 (0.6) | 8 (103) | 156.8 (0.8) | 1.2 (0.2) | 9 (112) | 144.9 (0.6) |
| 0.75 | 1.1 (0.7) | 6 (65) | 123.0 (1.4) | 1.1 (0.3) | 8 (94) | 113.9 (1.1) |
| 1.00 | 1.4 (0.6) | 14 (58) | 69.0 (1.6) | 1.2 (0.6) | 18 (33) | 67.5 (1.8) |

[a] As examples, results from two independent sets of 20 Å-inner-region simulations for the E106Q mutant of CAII are shown.
[b] Total simulation length (in *ns*) for each λ window, the number in parentheses is the length of equilibration identified by trend analysis.
[c] Number without any parentheses is the size of the block (in *ps*), number with parentheses is the number of blocks; these are determined after the equilibration sections of the trajectories are removed.
[d] In kcal/mol; the number in parentheses is the statistical error evaluated based on block average.

*Table 7-2.* Calculated $pK_a$s of the zinc-bound water in the wild type and E106Q mutant of CAII using thermodynamic integration[a]

| Calculation | $pK_a$[b] | $\Delta G_{CH(D)OH}$[c] | $\Delta\Delta G^{EXGR}$[d] |
|---|---|---|---|
| WT-20 | 7.1 | 151.3 (1.8) | −8.0 |
| WT-25 | 5.4 | 149.0 (1.6) | −8.8 |
| E106Q-20 | −3.0 | 137.5 (4.5) | −8.5 |
| E106Q-25 | −3.4 | 136.9 (1.5) | −8.9 |

[a] The free energies and contributions are in kcal/mol, the numbers in parentheses are statistical uncertainties.

[b] The $pK_a$ is calculated using 4-methyl-imidazole in solution (4-MI) as the reference compound; i.e., $pK_a^{CAII} = 7.0$ (exp.[115]) $+[\Delta G_{CH(D)OH}^{CAII} - \Delta G_{CH(D)OH}^{4-MI}]/1.370$; the computed value of $\Delta G_{CH(D)OH}^{4-MI}$ is 152.5 kcal/mol.

[c] The values here include the contribution from the EXGR correction ($\Delta\Delta G^{EXGR}$).

[d] The effect due to switching the QM/MM frontier treatment from "link-host-atom" exclusion to "link-host-group" exclusion (EXGR).

His residues at configurations sampled using the popular "link-host-atom" exclusion scheme changes the free energy derivatives by 8–9 kcal/mol despite that the QM/MM frontiers are far from the zinc-bound water. With this effect taken into account, the calculated $pK_a$ value for the zinc-bound water in the WT CAII is in encouraging agreement with experiment: the value is 7.1 (5.4) for the 20 (25) Å-inner-region simulations, as compared to the experimental value of around 7 [86].

An interesting finding in the $pK_a$ calculations for CAII is that the E106Q mutant is computed to reduce the $pK_a$ for the zinc-bound water by around ∼9 $pK$ units. Although the shift is very large, this result is reasonable considering that the mutation neutralizes a negative charge near the zinc-bound water, which is supported by a perturbative analysis of the free energy derivatives [91]. Strikingly, the $pK_a$ determined experimentally from the pH profile of $k_{cat}/K_M$ yielded no shift [92]. Considering that there is little structural change between the WT and E106Q mutant [93], the calculation result suggests that there may be a change in the mechanism for the step manifested by $k_{cat}/K_M$, perhaps similar to the behaviors of the cobalt substituted CAII [94] where, unlike the Zn(II) containing CAII, $k_{cat}/K_M$ depends on the concentration of bicarbonate. Therefore, $pK_a$ calculations can not only provide a stringent benchmark for the simulation protocol but also provide additional mechanistic insight that may stimulate new experimental investigations. In the investigation of proton pumping in complex biomolecules, for example, where electrostatics are crucial [14, 95] and major ambiguities exist concerning the titration states of various groups, we argue that $pK_a$ analysis of key residues is an indispensable step before the proton transfer pathways are explored.

   In this context, we note that $pK_a$ in proteins has also been calculated by minimization type of studies using QM/MM protocols with continuum electrostatics for the bulk solvent [96, 97]. Although rather impressive results have been obtained, the cases considered rarely involve residues deep in the active site or core region of proteins. In these latter cases, the proper sampling of conformational space is expected to be extremely important since structural response to titration is significant [73, 98]. By adopting a proper reference molecule in solution, errors in the proton affinity and interaction with the environment are largely cancelled out, which is an important reason that encouraging $pK_a$ shifts have been obtained using approximate QM methods such as SCC-DFTB in a QM/MM-FEP framework.

### 7.3.3.     The Importance of Sampling for Analyzing Chemistry

For the analysis of chemical reaction mechanisms, it is a common practice to use the minimum energy path (MEP) based calculations. Although MEP, or intrinsic reaction coordinate, is an extremely useful concept for small molecule reactions, it is less suitable for condensed phase reactions. This is largely because the existence of astronomically number of *local* stationary points and minimum energy paths in condensed phase systems, thus collecting a few paths is unlikely to capture the behavior of the system at finite temperature. In favorable cases, the MEPs can provide a qualitative understanding of the effect of the enzyme environment on the reaction of interest, which is indeed valuable [3, 89, 99, 100]; we will not review these cases here. In some cases, however, where *collective* rearrangements in the protein/solvent during the reaction is important, *local* MEP calculations may provide even misleading results as we discuss below.

   A useful example is the long-range proton transfer in CAII. As discussed above, there are multiple water wires of different length in the active site that connect the donor/acceptor groups (zinc-bound water, His 64). A question of interest is whether specific length of water wire dominates the proton transfer or all wires have comparable contributions. As the first approach, a large number of MEPs have been collected starting from different snapshots collected from equilibrium MD simulations at the SCC-DFTB/MM level. Since essentially a positive charge is transferred over a long-distance in the proton transfer, it was not surprising that the MEP energetics were found to depend on the origin of the starting structure, which reflects the fact that the active-site residues/solvent respond significantly to the proton transfer. For example, when the starting structure came from a CHOH equilibrium simulation, the proton transfer from the zinc-water to His 64 is largely *endothermic* according to MEPs (on average by as much as ∼13 kcal/mol, see Table 7-3). By contrast, when the starting structure came from a COHH simulation, the same proton transfer reaction was found largely *exothermic* based on MEP results (Table 7-3). Analysis of the interaction energies in the structures along the MEPs found that water molecules in the active site (within 7.5 Å from the zinc), which could reorient more easily than protein residues (dipoles), are largely responsible for this striking variation [27].

*Table 7-3.* Statistics for the barrier and reaction energy (in kcal/mol) associated with minimum energy paths based on different samplings of carbonic anhydrase[a]

| Sampling | | Two-water | Three-water | Four-water |
|---|---|---|---|---|
| CHOH[b] | $\Delta E_{act}$ | 16.2 (4.1) | 20.8 (4.9) | 23.4 (5.8) |
| | $\Delta E_{rxn}$ | 13.2 (5.4) | 14.7 (5.1) | 13.7 (6.0) |
| COHH[b] | $\Delta E_{act}$ | 3.7 (1.8) | 5.4 (3.6) | 8.7 (3.7) |
| | $\Delta E_{rxn}$ | −8.2 (5.0) | −17.3 (5.5) | −13.6 (5.4) |
| "TS-reorganized"[b] | $\Delta E_{act}$ | 6.8 (2.2) | 12.6 (1.9) | 17.4 (2.0) |
| | $\Delta E_{rxn}$ | −0.2 (3.5) | 3.3 (2.8) | 3.0 (2.7) |
| PMF[c] | $\Delta G_{act}$ | 13.0 "in" | | 12.8 "out" |
| | $\Delta G_{rxn}$ | 0.1 "in" | | 0.1 "out" |

[a] The simulations were carried out using the SCC-DFTB/MM-GSBP approach as discussed in details in Ref. [27] Independent of the chemical state simulated in the MD, all the energetics reported were determined relative to the chemical state involving the zinc-bound water. Numbers without parentheses are average values and those with parentheses are the standard deviations. The typical sample size for different simulations includes fifty minimum energy paths starting from independent snapshots.
[b] Indicate the chemical state for the QM region used in the MD simulations. CHOH: zinc-bound water and neutral His 64; COHH: zinc-bound hydroxide and doubly protonated His 64; "TS-reorganized": transferring protons along the water wire were constrained to be half-way between neighboring oxygen atoms, which is the protocol that approximately samples the protein/solvent configurations so that the proton-localized states are nearly degenerate.
[c] The details for the PMF calculations are given in Ref. [14], which involve using a collective reaction coordinate and on the order of 7–10 *ns* of simulations for each His 64 sidechain configuration ("in" vs. "out").

As an attempt to capture the "intrinsic barrier" for the proton transfer reaction, which is known to be close to be thermoneutral experimentally [101], we generated configurations from equilibrium MD simulations in which protons along a specific type of water wire were restrained to be equal distance from nearby heavy atoms (e.g., oxygen in water or $N\epsilon$ in His 64). In this way, the charge distribution associated with the reactive components is midway between the CHOH and COHH states, thus the active-site configuration was expected to facilitate a thermoneutral proton transfer process as confirmed by MEP calculations using such generated configurations as the starting structure (see Table 7-3). Interestingly, the barriers in such "TS-reorganized" MEPs showed a steep dependence on the length of the water wire; it was small ($\sim$6.8$\pm$2.2 kcal/mol) with short wires but substantially higher than the experimental value ($\sim$10 kcal/mol) with longer water wires (e.g., 17.4$\pm$2.0 kcal/mol for four-water wires).

This steep wire-length dependence is in striking contrast with the more rigorous PMF calculations of the same proton transfer process [14, 102]. In the PMF calculations, a collective coordinate [103] is used to monitor the progress of the proton transfer without enforcing specific sequence of events involving individual protons along the wire; the use of a collective coordinate is important because this allows

averaging over different water wire configurations, which is proper since the life-time of various water wires is on the pico-second time scale [27, 37], much faster than the time scale of the proton transfer ($\mu s$) [101]. In the PMF calculations, the wire-length dependence is examined by comparing results with different His 64 orientations ("in" and "out", which is about 8 and 11 Å from the zinc, respectively); both configurations are associated with multiple lengths of water wires but different relative populations (clearly, longer water wires have higher population for the "out" configuration). As shown in Table 7-3, the two sets of PMF calculations produced barriers of very similar values, which suggests that the length of the water wire (or donor–acceptor distance) is unlikely an important factor in determining the proton transfer rate. Physically speaking, this makes sense for the following reason. The $pK_a$ of both the donor and acceptor groups in CAII are shown experimentally to be around 7.0 [86] (which was also reproduced by the simulations discussed above), thus the dominant energetic component is related to the proton exchange between water in the wire and the donor/acceptor group; either with the conventional Grotthius mechanism [104, 105] (zinc-bound water first transfers a proton to the next water, generating a hydronium) or the "proton-hole" mechanism we proposed recently [14] (see below, which involves first transfer a water proton to His 64, generating a hydroxide), the energetics change is approximately 7 $pK_a$ unit, which is close to the experimentally observed barrier of $\sim 10$ kcal/mol. In other words, once either a hydronium or a hydroxide is generated, the species can move without a major barrier (consistent with the high mobility of hydronium and hydroxide in solution [106]) over a variable distance and therefore the dependence on the length of the water wire is not expected to be large (for more complete discussions, see Refs. [14, 102]).

What is then the origin for the striking difference between the MEP and PMF results? To illustrate the configurations accessed in the MEP and PMF simulations, the "excess coordination plots [14]" are compared in Figure 7-4. In such plots, the key information is the protonation state of the heavy atoms (donor/acceptor atoms plus mediating water oxygen) as a function of the reaction progress; positive "excess coordination" indicates extra proton (e.g., hydronium) while negative value indicates "proton hole" (e.g., hydroxide). The PMF simulations (Figure 7-4(b)) exhibit negative peaks along the anti-diagonal of the plot with both the donor/acceptor protonated during most of the reaction; as explained in Ref. [14], this is indicative of a "proton-hole" mechanism that involves the generation and *sequential* propagation of hydroxide rather than the classical Grotthuss mechanism *assumed* in all previous proton transfer studies of CAII and related systems. As shown in Figure 7-4(a), the configurations accessed in the MEP calculations show discrete positive peaks throughout the reaction, which indicates an almost fully *concerted* proton transfer. While sequential transfers of hydroxide (or hydronium), as discussed above, are not expected to exhibit a strong dependence on the distance of transfer, concerted transfers are expected to be more costly for long-distance transfers because many O—H bonds are broken simultaneously. We emphasize that the concerted nature of the MEPs is unlikely an artifact of the path search algorithm because multiple guesses are used and they all converged to the same path. The concerted nature is likely dictated by the

*Figure 7-4.* Excess coordination number plots for (**a**) MEP (three-water-bridge) [27] and (**b**) PMF (H64 "in") [14] simulations for the proton transfer in CAII. Note that the MEP calculation follows a very concerted mechanism while PMF simulation follows a step-wise "proton hole" mechanism

protein/solvent configurations used in the MEP calculations. As discussed above, to generate the "TS-reorganized" configurations, all transferring protons along the wire are constrained to be half-way between the neighboring heavy atoms; therefore, such sampled protein/solvent configurations would favor a concerted over step-wise proton transfers. Although *all* atoms in the inner region are allowed to move in the MEP searches, the local nature of MEPs does not allow collective reorganization of the active site residues/solvent molecules thus the "memory" of the sampling procedure is not erased.

In short, this example clearly illustrates that care must be exercised when using MEP to probe the mechanism of chemical reactions in biomolecules, especially when collective rearrangements in the environment are expected (e.g., reactions involving charge transport). In addition, when doing PMF type of simulations, the choice of a proper reaction coordinate is equally important; in this context, understanding the relative time-scale of various motions (e.g., water-wire rearrangements vs. proton transfer) is crucial in deciding what degrees of freedom should be properly averaged over [107], which is not always straightforward. Therefore, the validity of the reaction coordinate should always be tested with activated dynamics [80, 108] or transition path sampling [109].

## 7.4.    CONCLUSIONS AND PERSPECTIVES: WHAT'S THE NEXT STEP?

As many chapters in this book and recent reviews [4–6, 27, 110] indicate, substantial progresses are being made in the QM/MM community for both DFT/ab initio and SE based QM/MM methods. For relatively localized chemical reactions, very reliable

results can now be obtained as compared to experimentally measured barrier heights [8] and relative kinetic parameters such as kinetic isotope effects [7]. Regarding future challenges, in addition to the pressing need to further improve the QM and MM methods themselves (e.g., fast SE methods [17, 21, 41], more accurate DFT [111] and polarizable MM force fields [112]), there are several topics of particular interest, which we briefly mention below.

- First, although nanoseconds of straightforward MD simulations appear to be sufficient for dealing with rather localized chemical processes, such samplings are still too limited for reactions that are coupled to significant structural rearrangements in the biomolecular/solvent environment. Good examples include long-range proton/electron transfers and nucleotide hydrolysis in molecular motors and titration of residues in the interior of the protein. To treat those processes adequately, using enhanced sampling techniques based on either replica exchange or other sophisticated Monte Carlo schemes is likely required.
- Second, most QM/MM free energy simulations of chemical reactions have been based on potential of mean force simulations along a (set of) reaction coordinate(s) chosen based on chemical intuition. Selecting a proper reaction coordinate is increasingly difficult as the reactive process becomes more "delocalized" and tightly coupled to significant reorganizations in the environment. Adopting more sophisticated path sampling techniques [109] is likely another important necessity in QM/MM studies of increasingly complex biological problems. It's worth emphasizing that such path sampling techniques are computationally demanding, which once again highlights the importance of developing effective SE based QM/MM methods.
- Third, as the size and complexity of the biomolecular systems at hand further expand, there are more uncertainties in the molecular model itself. For example, the resolution of the X-ray structure may not be sufficiently high for identifying the locations of critical water molecules, ions and other components in the system; the oxidation states and/or titration states of key reactive groups might be unclear. In those cases, it is important to couple QM/MM to other molecular simulation techniques to establish and to validate the microscopic models before elaborate calculations on the reactive mechanisms are investigated. In this context, $pK_a$ and various spectroscopic calculations [113, 114] can be very relevant.
- Finally, in the spirit of this book's subject, it remains a major challenge to properly describe the structural and dynamic features of the environment at a coarse-grained but robust level. For example, in the SCC-DFTB/MM-GSBP protocol, the outer region is described as fixed with an apparent dielectric constant; how to include the slow motions for the outer region and adequately describe the relevant dielectric contribution are important issues that need to be solved. In this context, it is perhaps worth emphasizing again that one of the major advantages of adopting a multi-scale(resolution) type of model is that enhanced sampling techniques can be more effectively applied when the number of active degrees of freedom is relatively small.

## ACKNOWLEDGMENTS

## REFERENCES

1. Klaehn M, Braun-Sand S, Rosta E, Warshel A (2005) J Phys Chem B 109:15645
2. Zhang Y, Kua J, McCammon JA (2003) J Phys Chem B 107:4459
3. Cui Q, Karplus M (2002) J Am Chem Soc 124:3093
4. Senn HM, Thiel W (2007) Topics Curr Chem 268:173
5. Shurki A, Warshel A (2003) Adv Prot Chem 66:249
6. Gao J, Truhlar DG (2002) Annu Rev Phys Chem 53:467
7. Friesner RA, Guallar V (2005) Annu Rev Phys Chem 56:389
8. Claeyssens F, Harvey JN, Manby FR, Mata RA, Mulholland AJ, Ranaghan KE, Schutz M, Thiel S, Thiel W, Werner HJ (2006) Angew Chem Int Ed 45:6856
9. Hayashi S, Tajkhorshid E, Schulten K (2003) Biophy J 85:1440
10. Zhang YK, Liu HY, Yang WT (2000) J Chem Phys 112:3483
11. Hu H, Lu ZY, Yang WT (2007) J Chem Theor Comp 3:390
12. Kastner J, Senn HM, Thiel S, Otte N, Thiel W (2006) J Chem Theor Comp 2:452
13. Thiel W (1996) Adv Chem Phys 93:703
14. Lipkowitz KB, Boyd DB (eds) (1995) J Gao, In reviews in computational chemistry VII. VCH, New York
15. Riccardi D, König P, Prat-Resina X, Yu H, Elstner M, Frauenheim T, Cui Q (2006) J Am Chem Soc 128:16302
16. Weber W, Thiel W (2000) Theor Chem Acc 103:495
17. Tubert-Brohman I, Guimaraes CRW, Jorgensen WL (2005) J Chem Theor Comput 1:817
18. Sattelmeyer KW, Tubert-Brohman I, Jorgensen WL (2006) J Chem Theor Comput 2:413
19. Gonzalez-Lafont A, Truong T, Truhlar DG (1991) J Phys Chem 95:4618
20. Cui Q, Karplus M (2002) J Phys Chem B 106:1768
21. Nam K, Cui Q, Gao J, York DM (2007) J Theor Comp Chem 3:486
22. Elstner M, Porezag D, Jungnickel G, Elstner J, Haugk M, Frauenheim T, Suhai S, Seifert G (1998) Phys Rev B 58:7260
23. Elstner M, Frauenheim T, Kaxiras E, Seifert G, Suhai S (2000) Phys Stat Solid B 217:357
24. Elstner M, Frauenheim T, Suhai S (2003) Theochem 632:29
25. Elstner M (2006) Theor Chem Acc 116:316
26. Cui Q (2006) Theor Chem Acc 116:51
27. Riccardi D, Schaefer P, Yang Y, Yu H, Ghosh N, Prat-Resina X, Konig P, Li G, Xu D, Guo H et al (2006) J Phys Chem B 110:6458
28. Kruger T, Elstner M, Schiffels P, Frauenheim T (2005) J Chem Phys 122:114110

29. Sattelmeyer KW, Tirado-Rives J, Jorgensen W (2006) J Phys Chem A 110:13551
30. Otte N, Scholten M, Thiel W (2007) J Phys Chem A 111:5751
31. Elstner M, Jalkanen KJ, Knapp-Mohammady M, Frauenheim T, Suhai S (2001) Chem Phys 263:203
32. Elstner M, Hobza P, Frauenheim T, Suhai S, Kaxiras E (2001) J Chem Phys 114:5149
33. Mohle K, Hofmann HJ, Thiel W (2001) J Comp Chem 22:509
34. Liu H, Elstner M, Kaxiras E, Frauenheim T, Hermans J, Yang W (2001) Proteins 44:484
35. Cui Q (2006) Theor Chem Acc 116:51
36. Im W, Bernéche S, Roux B (2001) J Chem Phys 114:2924
37. Schaefer P, Riccardi D, Cui Q (2005) J Chem Phys 123:014905
38. Gregersen BA, York DM (2005) J Phys Chem B 109:536
39. Warshel A (1991) Computer modeling of chemical reactions in enzymes and solution. Wiley, New York
40. Elstner M (2007) J Phys Chem A 111:5614
41. Yang Y, Yu H, York D, Cui Q, Elstner M (2007) J Phys Chem A 111:10861–10873
42. Porezag D, Frauenheim T, Kohler T, Seifert G, Kaschner R (1995) Phys Rev B 51:12947
43. Seifert G (2007) J Phys Chem A 111:5609
44. Foulkes W, Haydock R (1989) Phys Rev B 40:12520
45. Wanko M, Hoffman M, Frauenheim T, Elstner M (2006) J Comput Aided Mol Des 20:511
46. Bondar N, Suhai S, Fischer S, Smith J, Elstner M (2007) J Struct Biol 157:454
47. Wanko M, Garavelli M, Bernardi F, Niehaus TA, Frauenheim T, Elstner M (2004) J Chem Phys 120:1674, cited By (since 1996): 42, URL www.scopus.com
48. Wanko M, Hoffman M, Strodel P, Koslowski A, Thiel W, Neese F, Frauenheim T, Elstner M (2005) J Phys Chem B 109:3606
49. Cui Q, Elstner M, Kaxiras E, Frauenheim T, Karplus M (2001) J Phys Chem B 105:569
50. Han W, Elstner M, Jalkanen KJ, Frauenheim T, Suhai S (2000) Int J Quant Chem 78:459
51. Saebra G, Walker R, Elstner M, Case D, Roitberg A (2007) J Phys Chem A 111:5655
52. Hu H, Elstner M, Hermans J (2003) Proteins Struct Funct Genet 50:451
53. Frisch MJ et al (2004) Gaussian 03, Revision B.05 (2003). Gaussian, Inc., Wallingford
54. Field MJ, Bash PA, Karplus M (1990) J Comput Chem 11:700
55. Zhang Y, Lee T, Yang W (1999) J Chem Phys 110:46
56. Gao J, Amara P, Alhambra C, Field MJ (1998) J Phys Chem A 102:4714
57. Konig PH, Hoffmann M, Frauenheim T, Cui Q (2005) J Phys Chem B 109:9082
58. Pu JZ, Gao JL, Truhlar DG (2004) J Phys Chem A 108:5454
59. Antes I, Thiel W (1999) J Phys Chem A 103:9290
60. Reuter N, Dejaegere A, Maigret B, Karplus M (2000) J Phys Chem A 104:1720
61. Lyne PD, Hodoscek M, Karplus M (1999) J Phys Chem A 103:3462
62. Das D, Eurenius KP, Billings EM, Sherwood P, Chattfield DC, Hodošček M, Brooks BR (2002) J Chem Phys 117:10534
63. Laio A, VandeVondele J, Rothlisberger U (2002) J Chem Phys 116:6941
64. Kalinowski JA, Lesyng B, Thompson JD, Cramer CJ, Truhlar DG (2004) J Phys Chem A 108:2545
65. Freindorg M, Gao J (1996) J Comput Chem 17:386
66. Bash PA, Ho LL, MacKerell AD Jr, Levine D, Hallstrom P (1996) Proc Natl Acad Sci 93:3698
67. Riccardi D, Li G, Cui Q (2004) J Phys Chem B 108:6467
68. Weber W, Hunenberger PH, McCammon JA (2000) J Phys Chem B 104:3668
69. Davis ME, Mccammon JA (1990) Chem Rev 90:509
70. Sagui C, Darden TA (1999) Annu Rev Biophys Biomol Struct 28:155
71. Gao J, Alhambra C (1997) J Chem Phys 107:1212

72. Nam K, Gao J, York DM (2005) J Chem Theor Comp 1:2
73. Riccardi D, Schaefer P, Cui Q (2005) J Phys Chem B 109:17715
74. Laino T, Mohamed F, Laio A, Parrinello M (2006) J Chem Theor Comp 2:1370
75. Brooks CL III, Karplus M (1989) J Mol Biol 208:159
76. Simonson T, Archontis G, Karplus M (1997) J Phys Chem B 101:8349
77. Dinner AR, Lopez X, Karplus M (2003) Theor Chem Acc 109:118
78. Li G, Zhang X, Cui Q (2003) J Phys Chem B 107:8643
79. Dinner AR, Blackburn GM, Karplus M (2001) Nature 413:752
80. Chandler D (1987) Introduction to modern statistical mechanics. Oxford University Press, New York
81. Rychaert JP, Ciccotti G, Berendsen HJ (1977) J Comput Phys 23:327
82. Li G, Cui Q (2003) J Phys Chem B 107:14521
83. Yang W, Bitetti-Putzer R, Karplus M (2004) J Chem Phys 120:9450
84. Warshel A, Sharma PK, Kato M, Parson WW (2006) Biochim Biophys Acta Proteins Proteom 1764:1647
85. Ashbaugh HS, Wood RH (1997) J Chem Phys 106:8135
86. Silverman DN, Lindskog S (1988) Acc Chem Res 21:30
87. Håkansson K, Carlsson M, Svensson LA, Liljas A (1992) J Mol Biol 227:1192
88. Toba S, Colombo G, Merz KMJ (1999) J Am Chem Soc 121:2290
89. Zhang X, Harrison D, Cui Q (2002) J Am Chem Soc 124:14871
90. Xu D, Guo H, Cui Q (2007) J Phys Chem A 111:5630
91. Riccardi D, Cui Q (2007) J Phys Chem A 111:5703
92. Liang Z, Xue Y, Behravan G, Jonsson B, Lindskog S (1993) Eur J Biochem 211:821
93. Xue Y, Liljas A, Jonsson B, Lindskog S (1993) Proteins Struct Funct Genet 17:93
94. Tu C, Tripp BC, Ferry JG, Silverman DN (2001) J Am Chem Soc 123:5861
95. Kato M, Pisliakov AV, Warshel A (2006) Proteins Struct Funct Bioinfor 64:829
96. Li H, Hains AW, Everts JE, Robertson AD, Jensen JH (2002) J Phys Chem B 106:3486
97. Jensen JH, Li H, Robertson AD, Molina PA (2005) J Phys Chem A 109:6634
98. Kato M, Warshel A (2006) J Phys Chem B 110:11566
99. Shaik S, Kumar D, de Visser SP, Altun A, Thiel W (2005) Chem Rev 105:2279
100. Bondar AN, Fischer S, Smith JC, Elstner M, Suhai S (2004) J Am Chem Soc 126:14668
101. Silverman DN (1995) Methods Enzymol 249:479
102. Riccardi D, König P, Guo H, Cui Q (2008) Biochem 47:2369–2378
103. Koenig P, Ghosh N, Hoffman M, Elstner M, Tajkhorshid E, Frauenheim T, Cui Q (2005) J Phys Chem B 110:548
104. de Grotthuss CJT (1806) Ann Chim 58:54
105. Mohammed OF, Pines D, Dreyer J, Pines E, Nibbering ETJ (2005) Science 310:83
106. Voth GA (2006) Acc Chem Res 39:143
107. Berezhkovskii AM, Szabo A, Weiss GH, Zhou HX (1999) J Chem Phys 111:9952
108. Karplus M (2000) J Phys Chem B 104:11
109. Bolhuis PG, Chandler D, Dellago C, Geissler PL (2002) Annu Rev Phys Chem 53:291
110. Zhang YK (2006) Theor Chem Acc 116:43
111. Mori-Sanchez P, Cohen AJ, Yang WT (2006) J Chem Phys 124:091102
112. Ponder JW, Case DA (2003) Adv Prot Chem 66:27
113. Cui Q, Karplus M (2000) J Phys Chem B 104:3721
114. Sproviero EM, Gascon JA, McEvoy JP, Brudvig GW, Batista VS (2007) Curr Opin Struct Biol 17:173
115. Lide DR (ed) (2005) CRC handbook chemistry, physics, 85th edn. CRC Press, Boca Raton

CHAPTER 8

# COARSE-GRAINED INTERMOLECULAR POTENTIALS DERIVED FROM THE EFFECTIVE FRAGMENT POTENTIAL: APPLICATION TO WATER, BENZENE, AND CARBON TETRACHLORIDE

GAURAV PRANAMI[1], LYUDMILA SLIPCHENKO[2], MONICA H. LAMM[1], AND MARK S. GORDON[2]

[1]*Department of Chemical and Biological Engineering, Iowa State University, Ames, IA 50011, USA*
[2]*Department of Chemistry, Iowa State University, Ames, IA 50011, USA,*
*e-mail: mark@si.msg.chem.iastate.edu (M.S. Gordon)*

**Abstract:**    A force matching technique based on previous work by Voth and co-workers is developed and employed to coarse grain intermolecular potentials for three common solvents: carbon tetrachloride, benzene, and water. The accuracy of the force-matching approach is tested by comparing radial distribution functions (RDF) obtained from simulations using the atomistic and coarse-grained potentials. Atomistic molecular dynamics simulations were performed using the effective fragment potential method (EFP). The RDFs obtained from molecular dynamics simulations of EFPs for carbon tetrachloride, benzene and water are in a good agreement with the corresponding experimental data. The coarse-grained potentials reproduce the EFP molecular dynamics center-of-mass RDFs with reasonable accuracy. The biggest discrepancies are observed for benzene, while the coarse-graining of water and spherically symmetric carbon tetrachloride is of better quality

**Keywords:**    Coarse-graining, Force-matching, Effective fragment potential method, Molecular dynamics, Radial distribution functions, Multiscale modeling

## 8.1.    INTRODUCTION

In the molecular dynamics (MD) [1, 2] technique, a system of particles evolves in time according to the equation of motion, $F_i = m_i \ddot{x}_i$, where $F_i$ is the net force acting on particle $i$, and $m_i$ and $\ddot{x}_i$ are the mass and acceleration of particle $i$, respectively. In a molecular system, typical bond lengths are of the order of angstroms while bond vibrations take place at the time scale of $10^{-13}$ s. Therefore, the equations of motion for atoms have to be integrated with time steps on the order of $10^{-15}$ s. However, many important chemical and biological phenomena in macromolecules take place at much larger time scales, as shown in Table 8-1.

*Table 8-1.* Characteristic time scales of different events in macromolecular systems

| System/Phenomena | Characteristic time scales |
|---|---|
| Fusion of micelles [3] | $10^{-2}$ s |
| Self-assembly of diblock copolymers [4] | $10^{-6}$ s |
| Entanglement of a polymer chain [5] | $10^{-5}$ s |
| Protein folding [6] | $10^{-6}$ s |
| DNA replication [6] | $10^{-3}$ s |
| Membrane fusion [6] | $10^{-1}$ s |

Large time and length scales of characteristic events in macromolecular systems, such as polymers, lipids and nanoparticles, prohibit a molecular simulation study at the atomistic level due to enormous CPU time and memory requirements. Moreover, when studying phenomena occurring at timescales on the order of $10^{-6}$ to $10^{-1}$ s, the behavior of the fast degrees of freedom, like bond fluctuations, are not always of interest. Therefore, a systematic approach to coarse-graining is needed, in which the unimportant degrees of freedom are eliminated, but the underlying physics governing the phenomena at larger length and time scales is retained.

Thus, the aim of coarse-graining techniques is to determine the effective interaction potentials between the coarse-grained sites such that the simulation of the coarse-grained (CG) system using the CG potentials yields properties that compare favorably with those of the corresponding atomistic system. An effective coarse-graining method should be systematic, automatic, and fast. Moreover, it should be flexible enough to handle different kinds of potentials and capable of generating CG potentials that would reproduce properties matching experimental data or the properties obtained from an atomistic simulation.

In general, coarse-graining an atomistic system requires a two step process, i.e., first, grouping the atoms into CG sites, as shown in Figure 8-1, and second, determining the effective bonded and non-bonded potentials between the CG sites. Most of the coarse-graining procedures reported in the literature can be classified into three categories: (i) optimization of potential parameters by fitting them to a desired property [7–11], (ii) structure matching [5, 12–22], and (iii) force matching [23, 24]. Systematic structure matching and force matching methods are preferable to the method of ad hoc parameter guessing.

In structure matching methods, potentials between the CG sites are determined by fitting structural properties, typically radial distribution functions (RDF), obtained from MD employing the CG potential (CG-MD), to those of the original atomistic system. This is often achieved by either of two closely related methods, Inverse Monte Carlo [12–15] and Boltzmann Inversion [5, 16–22]. Both of these methods refine the CG potentials iteratively such that the RDF obtained from the CG-MD approaches the corresponding RDF from an atomistic MD simulation.

Both the inverse Monte Carlo and iterative Boltzmann inversion methods are semi-automatic since the radial distribution function needs to be re-evaluated at

*Figure 8-1.* Coarse-graining procedure. (**a**) A snapshot of an atomistic system. (**b**) Groups of atoms of the atomistic system are combined into CG sites in order to reduce the number of degrees of freedom. (**c**) The atomistic system of (**a**) as represented by CG sites

each iteration. Moreover, convergence can be a problem if the potential of mean force (PMF) does not serve as a good initial guess. The inverse Monte Carlo method evaluates all the potentials at the same time and requires adequate sampling of the four-particle correlation functions during each iteration. The latter becomes time consuming if there are several different kinds of pair potentials. On the other hand, in the Boltzmann inversion method, potentials can be refined one at a time by keeping the rest of them constant. However, as the potentials depend upon each other, it is important to ensure that each potential does not change during the optimization of others. The advantage of the structure matching methods is that they result in potentials that will reproduce the correct structural properties; however, their main disadvantages are the necessity for frequent re-evaluation of radial distribution functions and the increasing complexity for a system with more than five coarse-grained sites.

In the force matching method, the effective pair-forces between coarse-grained sites are derived from the net force acting on chosen CG sites along an MD trajectory obtained from a short atomistic MD. The force-matching method has the advantage of being systematic and automatic because the CG pair forces are evaluated from the data gathered along the atomistic trajectory and there is no need to run multiple simulations. The force matching method has been successfully applied to study condensed phase liquids [23–25], ionic liquids [26, 27], $C_{60}$ nanoparticles [28, 29] and dimyristoylphosphatidylcholine (DMPC) lipid bilayers [30, 31].

In order to ensure accurate CG potentials, one needs to conduct MD simulations with a reliable atomistic potential model. The most desirable theoretical approach for the atomistic-scale simulations would be to use a level of quantum mechanics (QM) that can treat both intermolecular and intramolecular interactions with acceptable accuracy. Realistically, the minimal QM levels of theory that can adequately treat all different types of chemical forces are second order perturbation theory [32] (MP2)

and, preferably, coupled cluster (CC) theory with some accounting for triples; i.e., CCSD(T) [33]. Unfortunately, a sufficiently high level of QM comes at a significant computational cost; for example, CCSD(T) scales $\sim N^7$ with a problem size, where N is the number of atomic basis functions. This places serious limitations on the sizes of accessible molecular systems. Moreover, in order to obtain the CG potentials including the three-body terms, one needs to perform an extensive sampling of the atomistic-level system, which at present is impractical even for very short time-scales. Successful examples of QM-based coarse-graining have been presented based on Car-Parrinello simulations of atomistic systems [23].

An alternative approach is to replace an accurate but expensive first-principle-based technique by a reliable model potential. Such potentials, broadly referred to as molecular mechanics (MM), generally cannot account for bond-breaking, but can, in principle, account for the range of intermolecular interactions. However, using a fitted pair-wise potential may result in losing quantitative accuracy, predictability, and the underlying physics.

This contribution pursues a different approach for preserving the accuracy of the atomistic level, by using a model potential that is exclusively derived from first principles, the effective fragment potential method (EFP). The original EFP1 method [34, 35] was developed specifically to describe aqueous solvent effects on biomolecular systems and chemical reaction mechanisms, and contains fitted parameters for the repulsive term. A general (EFP2) method [36] is applicable to any solvent; it includes all of the essential physics and has no empirically fitted parameters. A force matching technique is applied to derive a coarse-grained potential from the molecular trajectories generated with EFP MD simulations. The quality of the EFP force matching is tested on carbon tetrachloride, benzene, and water systems. This contribution is the first application of a coarse-graining procedure to the EFP method.

## 8.2.    THEORY

### 8.2.1.    Effective Fragment Potential Method

The effective fragment potential method is a first-principles based model potential for describing intermolecular forces. The interaction energy in EFP1, specifically designed for modeling water, consists of electrostatic, induction, and fitted exchange-repulsion terms. Presently, three different EFP1 models are available, with fitting done to represent Hartree-Fock (HF), DFT/B3LYP, and MP2 levels of theory. These models are called EFP1/HF [34], EFP1/DFT [37], and EFP1/MP2 [38], respectively. In EFP1/MP2, fitted dispersion terms are also included.

The general EFP2 model can be applied to any solvent and includes electrostatic, induction, exchange-repulsion, dispersion, and charge-transfer components, all of which are derived from first-principles using long- and short-range perturbation theory. Charge-transfer interactions are not included in this work, since they are primarily important for charged species. All of the EFP2 parameters are generated during a MAKEFP run, performed for each unique molecule; e.g., benzene and $CCl_4$. Once

EFP parameters for a particular fragment are generated for a given atomic basis set, they can be used in a variety of applications. The various components of the non-bonded interactions between molecules are evaluated using the EFP2 generated parameters. The procedure has been described in elsewhere [36]; only the main points are summarized below.

The electrostatic energy is calculated using the distributed multipolar expansion introduced by Stone [39, 40], with the expansion carried out through octopoles. The expansion centers are taken to be the atom centers and the bond midpoints. So, for water, there are five expansion points (three at the atom centers and two at the O-H bond midpoints), while in benzene there are 24 expansion points. The induction or polarization term is represented by the interaction of the induced dipole on one fragment with the static multipolar field on another fragment, expressed in terms of the distributed localized molecular orbital (LMO) dipole polarizabilities. That is, the number of polarizability points is equal to the number of bonds and lone pairs in the molecule. One can opt to include inner shells as well, but this is usually not useful. The induced dipoles are iterated to self-consistency, so some many body effects are included.

The Coulomb point multipole model breaks down when fragments approach too closely, since then the actual electron density on the two fragments is not well approximated by point multipoles. Thus, electrostatic interactions become too repulsive whereas the induction energy is too attractive if fragments approach each other too closely. In order to avoid this unphysical behavior, electrostatic and induction energy terms are modulated by exponential damping functions with parameters being obtained from fitting the damped multipole potential to the Hartree-Fock one [41, 42]. In EFP2, the induction energy terms are damped in a similar way [43].

The exchange repulsion energy in EFP2 is derived as an expansion in the intermolecular overlap. When this overlap expansion is expressed in terms of frozen LMOs on each fragment, the expansion can reliably be truncated at the quadratic term [44]. This term does require that each EFP carries a basis set, and the smallest recommended basis set is 6-31++G(d,p) [45] for acceptable results. Since the basis set is used only to calculate overlap integrals, the computation is very fast and quite large basis sets are realistic.

The dispersion interaction can be expressed as the familiar inverse R expansion,

$$E_{disp} = \sum_n C_n R^{-n} \tag{8-1}$$

The coefficients $C_n$ may be derived from the (imaginary) frequency dependent polarizabilities summed over the entire frequency range [46]. If one employs only dipole polarizabilities the dispersion expansion is truncated at the leading term, with $n = 6$. In the current EFP2 code, an estimate is used for the $n = 8$ term, in addition to the explicitly derived $n = 6$ term. Rather than express a molecular $C_6$ as a sum over atomic interaction terms, the EFP2 dispersion is expressed in terms of LMO-LMO

interactions. In order to ensure that the dispersion interaction goes to zero at short distances, the damping term proposed by Tang and Toennies [47] is employed.

The effective fragment potential method is several orders of magnitude less computationally expensive than ab-initio methods because it evaluates intermolecular interactions by simplified formulas derived from perturbation series in terms of intermolecular distances and orbital overlap integrals. The most time-consuming terms in EFP2 are the charge-transfer (omitted in the current work) and the exchange-repulsion, which is evaluated using calculated on the fly orbital overlap integrals between different fragments. EFP2 can be used in MD simulations of moderately sized systems. For example, calculation of the energy and gradient for a system of 64 waters with periodic boundary conditions (PBC) requires about 2 s on one Opteron 2600 MHz processor. Despite its low computational cost, the accuracy of EFP in predicting structures and binding energies in weakly-bonded complexes and liquids is very high and comparable with that of MP2 [42].

## 8.2.2.     Force Matching Procedure

The aim of the force matching procedure is to obtain the effective pair-force between CG sites using the force data obtained from a detailed atomistic molecular dynamics (MD) trajectory. The current implementation of the force-matching method closely follows the formulation from Refs. [23, 24].

The first step of the systematic force matching procedure is to define the CG sites, which are generally the centers of mass or geometric centers of groups of atoms, as illustrated in Figure 8-1(b), thus eliminating the group's internal degrees of freedom. Following the coarse graining scheme depicted in Figure 8-1(b), the snapshot of the MD trajectory of Figure 8-1(a) will look like the one shown in Figure 8-1(c). In the next step, the forces and positions of atoms from the detailed atomistic MD are converted to forces and positions of CG sites as depicted in Figure 8-2.

Assume that there are a total of $N$ coarse-grained sites in the system for any one MD snapshot (p $=$ 1), with coordinates ($r_i = x_i, y_i, z_i$) and net forces, $F_i$, (where $i = 1 - N$) acting on them, and that these are known from the atomistic MD trajectory data. If $f_{ij}(r_i, r_j)$ represents the force acting on the $i$th CG site due to



*Figure 8-2.* Conversion of forces from atomistic MD to forces on coarse-grained sites

the $j$th CG site, then each snapshot from the MD trajectory results in the following $n(= N$ *for one snapshot*) equations:

$$\sum_{j=1}^{N} f_{ij} = F_i \quad i = 1, 2, 3 \ldots N \tag{8-2}$$

Here the pair-force $f_{ij}(r_i, r_j)$ is unknown, so a model pair-force $f_{ij}(r_i, r_j, p_1, p_2 \ldots p_m)$ is chosen, which depends linearly upon $m$ unknown parameters $p_1, p_2 \ldots p_m$. Consequently, the set of Eq. (8-2) is a system of linear equations with $m$ unknowns $p_1, p_2 \ldots p_m$. The system (8-2) can be solved using the singular value decomposition (SVD) method if $n > m$ (over-determined system), and the resulting solution will be unique in a least squares sense. If $m > n$, more equations from later snapshots along the MD trajectory should be added to the current set so that the number of equations is greater than the number of unknowns. Mathematically, $n = qN > m$ where $q$ is the number of MD snapshots used to generate the system of equations.

It is important to note that model pair-forces for the interactions *A-A, A-B, A-C*, etc. are different from each other although they may have the same functional form. If required, the interaction between two non-bonded $A$ CG sites ($A$-$A_{non-bonded}$) can be treated differently from the interaction between two bonded $A$ CG sites ($A$-$A_{bonded}$). In a system with $A, B \ldots E$ as chosen coarse-grained sites,

$$f_{ij} = \begin{cases} f_{ij}^{AA}(p_1^{AA}, p_2^{AA} \ldots p_a^{AA}) & if \quad ij = AA \\ f_{ij}^{AB}(p_1^{AB}, p_2^{AB} \ldots p_b^{AB}) & if \quad ij = AB \text{ or } BA \\ \quad \vdots & \quad \vdots \\ f_{ij}^{EE}(p_1^{EE}, p_2^{EE} \ldots p_z^{EE}) & if \quad ij = EE \end{cases} \tag{8-3}$$

Clearly, the total number of unknowns that need to be determined is $m = a + b + \ldots + z$ and a solution set for parameters $p_1, p_2 \ldots p_m$ is determined using the singular value decomposition or any other suitable method. The mean pair-force corresponding to the "potential of mean force" can be obtained in a systematic manner by averaging a number of sets of solutions for parameters $p_1, p_2 \ldots p_m$ obtained along the atomistic MD trajectory in which the phase space is sampled extensively.

A convenient and systematic way to represent $f_{ij}(r_{ij})$ ($r_{ij}$ is the distance between particles $i$ and $j$) as a linear function of unknowns is to employ cubic splines [48], as shown in Figure 8-3. The advantage of using cubic splines is that the function is continuous not only across the mesh points, but also in the first and second derivatives. This ensures a smooth curvature across the mesh points. The distance $r_{ij}$ is divided into 1-dimensional mesh points, thus, $f_{ij}(r_{ij})$ in the $k$th mesh ($r_k \leq r_{ij} \leq r_{k+1}$) is described by Eqs. (8-4), (8-5) and (8-6) [48].

*Figure 8-3.* $f_{ij}(r_{ij})$ as cubic splines. The distance, $r_{ij}$, between atoms $i$ and $j$ is divided into the mesh as shown. In each mesh, the pair-force $f_{ij}$ is modeled as a cubic polynomial

$$f_{ij}\left(r_k \leq r_{ij} \leq r_{k+1}\right) = A\left(r_k, r_{ij}, r_{k+1}\right) f\Big|_k + B\left(r_k, r_{ij}, r_{k+1}\right) f\Big|_{k+1}$$
$$+ C\left(r_k, r_{ij}, r_{k+1}\right) f''\Big|_k + D\left(r_k, r_{ij}, r_{k+1}\right) f''\Big|_{k+1} \tag{8-4}$$

$$A = \frac{r_{k+1} - r_{ij}}{r_{k+1} - r_k} \qquad B = 1 - A \qquad C = \frac{1}{6}\left(A^3 - A\right)\left(r_{k+1} - r_k\right)^2$$
$$D = \frac{1}{6}\left(B^3 - B\right)\left(r_{k+1} - r_k\right)^2 \tag{8-5}$$

Here, $f\Big|_k$, $f'\Big|_k$ and $f''\Big|_k$ are the values of the pair-force $f_{ij}(r_{ij})$ and its first and second derivatives, respectively, at mesh point $r_k$. Equation (8-4) ensures the continuity of the function and its second derivative at the mesh points. In order to make the first derivatives continuous across the mesh points $r_k$, an additional set of Eq. (8-6) is needed:

$$\frac{r_k - r_{k-1}}{6} f''\Big|_{k-1} + \frac{r_{k+1} - r_{k-1}}{3} f''\Big|_k + \frac{r_{k+1} - r_k}{6} f''\Big|_{k+1} = \frac{f\Big|_{k+1} - f\Big|_k}{r_{k+1} - r_k} - \frac{f\Big|_k - f\Big|_{k-1}}{r_k - r_{k-1}} \tag{8-6}$$

At the end points of the mesh, Eq. (8-6) cannot apply. Instead, one needs to introduce boundary conditions, for instance, at large $r_{ij}$ the pair-force $f_{ij}$ is usually zero. It is important to note that the mesh sizes should not necessarily be uniform. For example, at those separations for which the pair-force varies rapidly with distance the mesh size can be chosen to be small enough to capture all of the variations.

If some meshes do not get sampled in the set of Eq. (8-2), i.e., if the coefficients of the correspondingare $f\big|_k$ and $f''\big|_k$ are zero, then these unknowns are removed from the set of equations and set equal to zero. Equation (8-6) and boundary conditions need to be satisfied exactly, however, in this work we have solved all the equations in least squared sense. By solving the combined set of Eqs. (8-2– 8-6). Solutions obtained for a large number of such sets of equations are averaged to reduce statistical noise in the CG pair-force. Then, a suitable analytic function should be fitted to the tabulated $f_{ij}(r_{ij})$. If no distinction is made between $A$-$A_{bonded}$ and $A$-$A_{non\text{-}bonded}$ interactions in the force matching procedure the resulting coarse-grained $A$-$A$ pair-force will have the combined effect of bonded and non-bonded interactions. In a typical atomistic MD simulation, bonded and non-bonded interactions both occur at short $A$-$A$ distances. Therefore, the coarse-grained $A$-$A$ interaction force may not be physically correct. In general, separate treatments of bonded and non-bonded interactions are preferred even though it increases the total number of unknowns and the size of the linear least squared problem.

Once all the coarse-grained interactions are determined using the force-matching procedure, they need to be validated by running a MD simulation of the coarse-grained system. Comparing properties such as pair correlation function(s) obtained from the coarse-grained and original atomistic MD is a direct test of the quality of coarse-graining.

The reduction in the number of degrees of freedom can lead to an incorrect pressure in the simulation of the coarse-grained systems in NVT ensembles or to an incorrect density in NPT ensembles [24]. The pressure depends linearly on the pair-forces in the system, hence the effect of the reduced number of degrees of freedom can be accounted for during the force matching procedure [24]. If $T$ is the temperature, $V$ the volume, $N$ the number of degrees of freedom of the system, and $k_b$ the Boltzmann constant then the pressure $P$ of a system is given by

$$P = \frac{Nk_bT}{3V} + \frac{1}{3V}\sum_{i<j} \vec{f}_{ij} \bullet \vec{r}_{ij} \tag{8-7}$$

In order to compensate for the reduced number of degrees of freedom in the coarse-grained system, the following constraint should be added to the set of Eq. (8-2):

$$\sum_{i<j} \vec{f}_{ij} \bullet \vec{r}_{ij} = 3P^{At-MD}V^{CG} - N^{CG}k_bT \tag{8-8}$$

Here, $P^{At\text{-}MD}$ is the pressure in the system in detailed atomistic MD and $V^{CG}$ and $N^{CG}$ are the volume and number of degrees of freedom of the coarse-grained system. The left side of this equation is evaluated for the coarse-grained system for each snapshot during force-matching.

## 8.3.    COMPUTATIONAL DETAILS

### 8.3.1.    EFP MD Simulations

Molecular dynamics simulations of liquid carbon tetrachloride, benzene, and water were performed using the effective fragment potential method, as implemented in the GAMESS (General Atomic and Molecular Electronic Structure System) electronic structure package [49, 50]. EFP2 parameters for benzene were obtained using the 6-311++G(3df,2p) basis set at the MP2/aug-cc-pVTZ [51] geometry of the benzene monomer, with C—C and C—H bond lengths of 1.3942 Å and 1.0823 Å, respectively. EFP2 parameters for $CCl_4$ were generated by using the 6-311++G(d,p) basis [52–54], with the monomer geometry optimized at the MP2/6-311G(d,p) level (C—Cl bond length of 1.772 Å). The EFP1/MP2 [38] potential was used for water.

The MD simulations were carried out in an NVT ensemble at ambient conditions; each simulation contained 64 molecules in a cubic box with periodic boundary conditions. The temperature is set to 300K in all the simulations. Table 8-2 summarizes details of the EFP-MD simulations used for force-matching. In particular, the type of the potential, box size, time step of integration, frequency of data sampling, total number of sampled configurations, and the total time of the equilibrated MD simulation are listed for each system. Initial equilibration of the systems was performed before recording the data for force matching. In order to ensure good energy conservation in the MD simulations, switching functions were employed for all EFP interaction terms at long distances [55]. Additionally, in simulations of water, Ewald summations were used to treat long-range electrostatic interactions (charge-charge, charge-dipole, dipole-dipole, and charge-quadrupole).

Since EFP employs frozen internal geometries of fragments, during the MD simulations, $CCl_4$, benzene and water molecules are treated as rigid bodies with a net force and torque acting on each center of mass (COM). Thus, the net forces acting on COMs required for force matching are directly available from the EFP MD simulations. The information about torques is not used in force-matching because each molecule is represented as a point at its COM.

*Table 8-2.*  EFP-MD simulation parameters[a]

| System | Potential | Box Length (Å) | Timestep (fs) | Frequency (fs) | Samples | Total simulation time (ps) |
|---|---|---|---|---|---|---|
| $CCl_4$ | EFP2 | 21.77 | 0.3 | 30 | 1200 | 36 |
| Benzene | EFP2 | 21.20 | 0.5 | 50 | 500 | 25 |
| Water | EFP1/MP2 | 12.40 | 0.3 | 30 | 1000 | 30 |

[a]For each system, columns specify the type of the potential, simulation box size, time step of MD integration, frequency at which data is sampled for force-matching, total number of configurations sampled in MD simulations and the total time of equilibrated MD simulations.

### 8.3.2. Force Matching

Carbon tetrachloride, benzene and water molecules have been coarse-grained to their COM using the force matching technique described in Section 8.2.1. The effective COM pair-force was modeled using cubic splines over a range of distances, described by an inner cutoff and an outer cutoff, with the mesh sizes summarized in Table 8-3. The outer cutoff of the model pair-force was set such that it never exceeded half of the simulation box length and large enough to ensure that the effective pair-force obtained from force-matching naturally approaches zero at the chosen outer cutoff. The inner cut-off can be safely chosen as zero or it can be approximated as a distance which is smaller than the smallest separation between a pair of CG sites sampled in the atomistic MD. The mesh-size should be small enough to capture all the features of the effective pair-force but, as mentioned earlier, smaller mesh-sizes result in more unknowns. Consequently, a smaller mesh is used in the regions where the CG pair-force is sharply repulsive and varies rapidly with distance. A total of $k$ meshes are used to model an interaction that is expressed in $2k + 2$ unknowns.

*Table 8-3.* Force matching details[a]

| System | Distance (Å) | Mesh-Size (Å) | Unknowns | Configurations per Set | Number of sets averaged |
|---|---|---|---|---|---|
| $CCl_4$ | 4–6 | 0.05 | 178 | 3 | 400 |
| | 6–10.8 | 0.1 | | | |
| Benzene | 3–6 | 0.05 | 202 | 3 | 150 |
| | 6–10 | 0.1 | | | |
| Water | 2–3.5 | 0.025 | 222 | 4 | 250 |
| | 3.5–6 | 0.05 | | | |

[a] For each studied system, the range of distances at which pair-forces are modeled as cubic splines is given, as well as mesh sizes and the number of resulting unknowns, the number of configurations included in a set to generate an over-determined system of equations, and the number of sets for which the least squared solution is averaged.

The net forces acting on the COMs of all molecules in a given MD snapshot were equated to the corresponding net force obtained from the model pair-force; consequently, each configuration yields 64 equations since each EFP-MD simulation contains 64 molecules. Three or four (see Table 8-3) MD configurations were used to generate a set of equations such that the number of equations was greater than the number of unknowns. Solutions for a number of such sets were averaged to obtain the effective mean COM pair-force. In the results reported here, the pressure is not constrained.

At short distances, approximately equal to the excluded volume diameter, effective pair forces obtained from force matching exhibit unphysically large fluctuations. This is largely due to inadequate sampling of configurations at short distances in

the EFP-MD simulation. These short-range pair-force data were ignored in further analysis. Ignoring the force data may lead to some inconsistency in the agreement of properties of atomistic and CG systems; however, as very few CG sites exist at such small separations in the EFP-MD simulation, this should not lead to significant error if averaging is done over a large number of sets during force matching. The remaining pair-force data, $f(r)$, obtained from force-matching, are fitted to the following function:

$$F(r) = \int \sum_{n=2}^{16} \frac{A_n}{r^n} \tag{8-9}$$

The fitting coefficients $A_n$ for CCl4, benzene, and water are listed in Table 8-4. The corresponding effective COM pair-potential, $U(r)$, was obtained by integrating $F(r)$ with the condition that the potential is zero at the outer cutoff:

$$U(r) = -\int F(r)\, dr \tag{8-10}$$

The coarse-grained pair-force and pair-potential were used for carrying out the molecular dynamics simulations (coarse-grained MD, CG-MD) of 64 points, each point representing a COM of a $CCl_4$ or benzene or water molecule, at the same conditions as used for the corresponding atomistic MD (Table 8-2). All CG-MD simulations were run using the LAMMPS [56] (Large-scale Atomic/Molecular Massively Parallel Simulator) molecular simulation code available at http://lammps.sandia.gov. LAMMPS is capable of running MD simulations using tabulated pair-forces and

*Table 8-4.* Fitting coefficients $A_n$ corresponding to Eq. (8-9). The units of $r$ and $F(r)$ are Å and kcal/mol-Å, respectively, for curve fitting

|          | Carbon Tetrachloride        | Benzene                     | Water                        |
|----------|-----------------------------|-----------------------------|------------------------------|
| $A_2$    | $-1.969703152101180E+21$    | $1.244993048322720E+21$     | $2.175775867955310E+17$      |
| $A_3$    | $4.306989690407250E+21$     | $-2.620983313167290E+21$    | $-8.551913708210820E+17$     |
| $A_4$    | $-4.333455273390900E+21$    | $2.545558326687080E+21$     | $1.550778775883470E+18$      |
| $A_5$    | $2.659644152606080E+21$     | $-1.511552672176440E+21$    | $-1.719444254190780E+18$     |
| $A_6$    | $-1.112575488728160E+21$    | $6.129876548482810E+20$     | $1.302147676537040E+18$      |
| $A_7$    | $3.356022412317150E+20$     | $-1.795710744899950E+20$    | $-7.124434785792410E+17$     |
| $A_8$    | $-7.528548090161490E+19$    | $3.918212052328320E+19$     | $2.903928257920110E+17$      |
| $A_9$    | $1.276021907548210E+19$     | $-6.468394131217870E+18$    | $-8.957222760163820E+16$     |
| $A_{10}$ | $-1.642124116314870E+18$    | $8.117691446975380E+17$     | $2.100873205120480E+16$      |
| $A_{11}$ | $1.596839968332750E+17$     | $-7.706208900534920E+16$    | $-3.728428139785960E+15$     |
| $A_{12}$ | $-1.155063457121670E+16$    | $5.446871839354590E+15$     | $4.928243298286980E+14$      |
| $A_{13}$ | $6.026936982240990E+14$     | $-2.779449283288870E+14$    | $-4.704569177419660E+13$     |
| $A_{14}$ | $-2.144530055247100E+13$    | $9.678942266368310E+12$     | $3.066010958528650E+12$      |
| $A_{15}$ | $4.658131620368480E+11$     | $-2.058787501037060E+11$    | $-1.221016058567130E+11$     |
| $A_{16}$ | $-4.660182908774720E+09$    | $2.018075773802740E+09$     | $2.241814292894480E+09$      |

pair-potentials. Therefore, the CG pair-forces obtained from force matching can be directly used to run the CG MD simulations. The integration time step was 1 fs. Each equilibrated CG-MD simulation was 3 ns long and the position data was collected every 1 ps. In order to test the ability of coarse-grained potentials to reproduce properties of atomistic systems, RDFs obtained from CG-MD are compared below to the corresponding COM-COM RDFs from atomistic EFP-MD.

## 8.4. RESULTS AND DISCUSSION

### 8.4.1. EFP-MD Simulations

EFP radial distribution functions for liquid $CCl_4$, benzene, and water are presented in Figures 8-4, 8-5 and 8-6. EFP2 and experimental [57] C—Cl and Cl—Cl RDFs for liquid carbon tetrachloride are shown in Figure 8-4. The Cl—Cl EFP2 RDFs are in good agreement with the experimental data. The discrepancies in the C—Cl RDF curves are more significant, although the qualitative features of the experimental RDF are reproduced. It is possible that the strong structural enhancement observed in the experimental $CCl_4$ RDFs is an artifact that arises due to numerical instabilities when specific atom-atom RDFs are obtained from X-ray and neutron analysis data [58]. To confirm this the so-called $G_d(r)$ functions were calculated. The $G_d^X(r)$ and $G_d^n(r)$ functions are Fourier transforms of the X-ray and neutron diffraction distinct structure functions, respectively; the latter are unambiguously determined experimentally [57]. For $CCl_4$, $G_d^X(r)$ and $G_d^n(r)$ functions are connected to specific atom-atom RDFs in the following way:

$$G_d^X(r) \approx 0.00 g_{CC}(r) - 0.12 g_{CCl}(r) - 0.88 g_{ClCl}(r)$$
$$G_d^n(r) \approx 0.02 g_{CC}(r) + 0.25 g_{CCl}(r) + 0.75 g_{ClCl}(r),$$

where $g_{CC}(r)$, $g_{CCl}(r)$, and $g_{ClCl}(r)$ are C—C, C—Cl, and Cl—Cl RDFs, respectively.

Experimental and EFP-MD $G_d^n(r)$ and $G_d^X(r)$ functions are shown in Figure 8-4c and 8-4d, respectively. The agreement between the EFP2 and experimental $G_d$-functions is better than that between specific RDFs, although some discrepancies remain. EFP2 overestimates the heights of the peaks at 4.0 Å in both $G_d$ graphs, and the peaks at 6.2 Å are slightly shifted to longer distances.

Figure 8-5 shows the EFP2 and experimental [59] C—C RDFs for liquid benzene. The EFP2 RDF is in reasonable agreement with the experimental curve, with three distinct peaks in the 4–7 Å region. These peaks might correspond to different orientations of neighboring benzene molecules in solution, e.g., T-shaped-like and parallel-displaced configurations are possible. Compared to experiment, the EFP2 RDF features are slightly more pronounced, suggesting that EFP2 over-structures

*Figure 8-4.* Comparison of EFP2 and experimental RDFs and $G_d$-functions for liquid carbon tetrachloride: (**a**) C−Cl, and (**b**) Cl−Cl RDFs, (**c**) $G_d^n$, (**d**) $G_d^X$



*Figure 8-5.* Comparison of EFP2 and experimental C−C RDFs for liquid benzene

*Figure 8-6.* Comparison of EFP1/MP2 and experimental O−O RDFs for liquid water

liquid benzene. This may be due to the fact that EFP2 slightly overestimates the interactions between benzene molecules [42].

The EFP1/MP2 oxygen-oxygen RDF for water is shown in Figure 8-6. The positions of the peaks in the EFP1/MP2 and experimental RDFs are in excellent agreement [60], but the intensities of the EFP peaks are overestimated, i.e., EFP1/MP2 over-structures the water RDF. Some degree of over-structuring has been attributed to omitting quantum affects [61], although such affects are likely to be very small when no H atoms are involved. Over-structuring could arise due to intrinsic inaccuracies in the EFP1/MP2 potential, for example, water-water interactions that are too strong. A detailed analysis of the performance of different EFP models for liquid water can be found elsewhere [62].

## 8.4.2. Coarse-Graining

Because carbon tetrachloride is a spherically symmetric molecule, it is logical to coarse-grain it to its COM and represent it as a single point. The effective pair-force and pair-potential for the $CCl_4$ COM obtained from force matching are shown in Figure 8-7. The CG pair-potential was obtained by integrating the pair-force according to Eq. (8-10). The potential, $U(r)$, becomes sharply repulsive below $r \sim 4$Å indicating that the $CCl_4$ excluded volume corresponds to the diameter $\sim 4$Å. This is reasonable given that the C—Cl bond length in $CCl_4$ is 1.767Å and the excluded volume diameter should be slightly larger than twice the C—Cl bond length ($\approx 3.534$ Å). Therefore, the force matching method has taken excluded volume into account. The potential $U(r)$ is smooth and slowly varying with a wide minimum at $r \sim 7$Å.

*Figure 8-7.* Coarse-graining of CCl$_4$. (**a**) *Black*: COM-COM RDF from EFP-MD, (**b**) *red*: RDF from CG-MD, (**c**) *orange circles*: the effective COM pair-force, (**d**) *green*: polynomial fit of the force matching data, (**e**) *blue*: the effective COM pair potential

The comparison of the CG and EFP RDFs clearly indicates that the coarse-grained potential is able to reproduce the liquid structure of CCl$_4$ reasonably well. The locations of the CG RDF peaks are in good agreement with those in the EFP RDF, although the CG peaks are a bit higher. This may be attributed to an overly steep repulsive CG pair-force (at r $\sim$ 4.5 Å) used in the CG MD. There is an uncertainty about the nature of the repulsive CG pair-force at short distances, where the pair-force obtained from force-matching exhibits large unphysical fluctuations due to insufficient sampling of pairs at short separations (r $\sim$ 4.5 Å) in the atomistic MD. Ignoring the data with large unphysical fluctuations in CG pair-force and replacing it with a fit through the remaining pair-force data may make the pair-force strongly repulsive at r $\sim$ 4.5 Å as shown in Figure 8-7. This repulsion may be stronger than the repulsion in the corresponding EFP-MD. Strong repulsion at close separations (r $\sim$ 4.5 Å) in the CG MD at the same density as the EFP-MD probably results in a more structured liquid, so sharper peaks are observed.

Favorable coarse-graining results for CCl$_4$ are not surprising because this molecule is spherically symmetric. Planar benzene presents a more stringent test for the force-matching approach. Figure 8-8 shows the effective pair-force and pair-potential for the benzene COM.

The locations of the peaks in the benzene CG RDF are in good agreement with the EFP RDF, but the heights of the peaks consistently exceed that of the EFP RDF. Moreover, the first peak in the CG RDF is not as broad as the first peak in the EFP RDF. This indicates that the coarse-grained pair-potential produces a more structured

*Figure 8-8.* Coarse-graining of benzene. (**a**) *Black*: COM-COM RDF from EFP-MD, (**b**) *red*: RDF from CG-MD, (**c**) *orange circles*: the effective COM pair-force, (**d**) *green*: polynomial fit of the force matching data, (**e**) *blue*: the effective COM pair potential

liquid compared to that of the EFP. It is possible that these discrepancies arise due to the use of a system (64 molecules) that is too small and a 25 ps EFP MD run that is too short. This might result in inadequate configuration sampling. For example, as noted above for $CCl_4$, inadequate configuration sampling at short distances, r $\sim$ 4–4.5 Å, the CG leads to unphysically large fluctuations in the CG pair force. Therefore the data in this range has to be neglected. This leads to the loss of information relating to the minimum energy parallel displaced benzene dimer configuration, see Figure 8-9(b). Fitting a curve using the remaining *f(r)* data makes the interaction at short distances more repulsive, resulting in a larger excluded volume and a narrower first peak. The 5.0 Å shoulder in the first peak of the EFP RDF can be associated with the T-shaped benzene dimer structure, Figure 8-9(a). Due to inadequate sampling and repulsion at short distances, the peak in the CG RDF lacks this shoulder and is narrower and higher than the corresponding EFP RDF peak. Additionally, due to the short EFP run, the coarse-grained potential has been averaged over only 150 sets, compared to 250 and 400 sets for water and $CCl_4$, respectively (Table 8-3). This is because the EFP MD simulations of liquid benzene are more computationally demanding than the simulations for water or $CCl_4$ (see Table 8-5).

Water is a very important and widely used solvent. Many of the unique properties of water are the result of the complex interactions that occur among water molecules. A water molecule is planar and highly polar, so it is an important system to test with the force matching approach. The coarse-graining results for water are summarized

*Figure 8-9.* The minimum energy configurations of benzene dimer: (**a**) T-shaped, (**b**) parallel-diplaced, and (**c**) edge-to-edge structures

in Figure 8-10. The CG pair-force and pair potentials match qualitatively with the results reported in the literature [24].

The CG RDF is in reasonable agreement with the EFP MD COM-COM RDF. In particular, the first peak is in excellent agreement, while the second peak is slightly off (4.2 vs. 4.5 Å and slightly lower compared to the EFP RDF). After about 5 Å there is almost no structure in the CG RDF. Water is a complicated molecule to coarse-grain to a single site due to the presence of Van der Waals and coulombic interactions. Moreover, as it is a highly polar molecule, coarse-graining it to a single point at its COM may not be the best choice. Despite these shortcomings, the one-site coarse-grained potential is able to reproduce the first and second peaks, indicating that the force matching technique works reasonably well even for polar non-symmetric molecules.

The real advantage of coarse-graining is the speed up due to the reduction in the number of degrees of freedom and due to substituting a complex EFP potential by

*Table 8-5.* CPU time per timestep for EFP-MD and CG-MD simulations and CPU speed-up due to coarse-graining[a]

| | EFP-MD | | CG-MD | | |
|---|---|---|---|---|---|
| | CPU time (s) | Timestep (fs) | CPU time (s) | Timestep (fs) | Speed-up[b] |
| $CCl_4$ | 24.80 | 0.3 | 0.000238 | 1.0 | $3.47 \cdot 10^5$ |
| Benzene | 117.30 | 0.5 | 0.000249 | 1.0 | $9.42 \cdot 10^5$ |
| Water | 1.03 | 0.3 | 0.000239 | 1.0 | $1.29 \cdot 10^4$ |

[a] All CPU times reported are for MD simulations carried out on one AMD 280 Opteron 2.4 GHz processor.
[b] Speed-up = EFP-MD (CPU time/time step)/CG-MD (CPU time/time step)

*Figure 8-10.* Coarse-graining of water. (**a**) *Black*: COM-COM RDF from EFP-MD, (**b**) *red*: RDF from CG-MD, (**c**) *orange circles*: the effective COM pair-force, (**d**) *green*: polynomial fit of the force matching data, (**e**) *blue*: the effective COM pair potential

a simpler polynomial one. The speed-ups achieved for the three studied systems are listed in Table 8-5. Additional speed-up can be achieved by using larger timestep for running CG-MD, which is reasonable to do because the CG potentials vary less rapidly with distance compared to underlying interaction potentials used in atomistic MD.

## 8.5.    CONCLUSIONS

In the work presented here, the force matching technique is used to coarse-grain three typical solvents, carbon tetrachloride, benzene, and water, to their centers of mass. The accuracy of the force-matching is tested by comparing structural properties, namely, the radial distribution functions, of the underlying atomistic and coarse-grained systems. The atomistic MD simulations were performed using the effective fragment potential method. The EFP is a first-principles-based method designed for describing intermolecular interactions.

The RDFs for all three systems obtained from coarse-grained MD compare favorably with RDFs from EFP MD. Owing to its spherical symmetry, $CCl_4$ was found to be the most amenable to coarse-graining using the force matching method. For benzene, the coarse-grained MD produced a more structured liquid than that obtained with the atomistic MD. This might be attributed to a limited sampling of configuration space in the atomistic MD that is required for force

matching. This issue needs to be explored further. On the other hand, the coarse-grained RDF of water is in reasonably good agreement with the corresponding atomistic RDF.

The quality of coarse-grained potentials critically depends on the accuracy of the underlying atomistic MD from which the data required for force matching are generated. EFP, used for "atomistic" MD simulations in this work, is a promising technique for capturing the chemistry of liquids and solvents. Most of the previous applications of the EFP method focused on analysis of reactions and properties in complexes and clusters. This work presents the results of EFP MD simulations on liquid $CCl_4$, benzene, and water. In all cases, the EFP RDFs are in reasonable agreement with the available experimental data. EFP does tend to produce sharper peaks in RDFs, suggesting that some overstructuring of liquids may occur. For coarse-graining, the quality of the sampling of conformational space in an EFP MD simulation can be an issue. For example, more extensive sampling of conformational space in benzene could potentially improve the quality of its coarse-graining. Because it is a first principles-based technique, EFP is significantly more expensive than other force fields. This makes long EFP-MD simulations computationally demanding. These issues will be addressed in future work. Future contributions will also extend the methodology presented here to coarse-graining polymers, in order to study the mechanisms of their aggregation.

## ACKNOWLEDGEMENT

## REFERENCES

1. Allen MP, Tildesley DJ (1996) Computer simulations of liquids. Oxford University Press.,Oxford
2. Frenkel D, Smit B (2001) Understanding molecular simulations. Academic Press, San Diego.
3. Smit B, Esselink K, Hilbers PAJ, Vanos NM, Rupert LAM, Szleifer I (1993) Langmuir, 9(1):9–11
4. Chushak Y, Travesset A (2005) J Chem Phys 123(23)
5. Baschnagel J, Binder K, Paul W, Laso M, Suter UW, Batoulis I, Jilge W, Burger T (1991) J Chem Phys 95(8):6014–6025
6. Nielsen SO, Lopez CF, Srinivas G, Klein ML (2004) J Phys Condens Matter 16(15):R481–R512
7. Shelley JC, Shelley MY, Reeder RC, Bandyopadhyay S, Klein ML (2001) J Phys Chem B 105(19):4464–4470
8. Shelley JC, Shelley MY, Reeder RC, Bandyopadhyay S, Moore PB, Klein ML (2001) J Phys Chem B 105(40):9785–9792
9. Marrink SJ, de Vries AH, Mark AE (2004) J Phys Chem B 108(2):750–760

10. Marrink SJ, Mark AE (2003) J Am Chem Soc 125(49):15233–15242
11. Harmandaris VA, Adhikari NP, van der Vegt NFA, Kremer K (2006) Macromolecules 39(19): 6708–6719
12. Lyubartsev AP, Laaksonen A (1995) Phys Rev E 52(4):3730–3737
13. Lyubartsev AP, Laaksonen A (1999) J Chem Phys 111(24):11207–11215
14. Murtola T, Falck E, Patra M, Karttunen M, Vattulainen I (2004) J Chem Phys 121(18):9156–9165
15. Murtola T, Falck E, Karttunen M, Vattulainen I (2007) J Chem Phys 126(7):075101-1–075101-14
16. Soper AK (1996) Chem Phys 202(2–3):295–306
17. Reith D, Putz M, Muller-Plathe F (2003) J Comput Chem 24(13):1624–1636
18. Tschop W, Kremer K, Batoulis J, Burger T, Hahn O (1998) Acta Polymerica 49(2–3):61–74
19. Li XJ, Kou DZ, Rao SL, Liang HJ (2006) J Chem Phys 124(20):204909-1–204909-7
20. Li XJ, Ma XJ, Huang L, Liang HJ (2005) Polymer 46(17):6507–6512
21. Depa PK, Maranas JK (2005) J Chem Phys 123(9):094901-1–094901-7
22. Depa PK, Maranas JK (2007) J Chem Phys 126(5):054903-1–054903-8
23. Izvekov S, Parrinello M, Burnham CJ, Voth GA (2004) J Chem Phys 120(23):10896–10913
24. Izvekov S, Voth GA (2005) J Chem Phys 123(13):134105-1–134105-13
25. Izvekov S, Voth GA (2005) J Phys Chem B 109(14):6573–6586
26. Wang YT, Izvekov S, Yan TY, Voth GA (2006) J Phys Chem B 110(8):3564–3575
27. Wang YT, Voth GA (2006) J Phys Chem B 110(37):18601–18608
28. Violi A, Voth GA (2005) In High Performance Computing and Communications, Proceedings, pp 938–947
29. Izvekov S, Violi A, Voth GA (2005) J Phys Chem B 109(36):17019–17024
30. Izvekov S, Voth GA (2005) J Phys Chem B 109(7):2469–2473
31. Izvekov S, Voth GA (2006) J Chem Theory Comput 2(3):637–648
32. Moller C, Plesset S (1934) Phys Rev 46:618
33. Raghavachari K, Trucks GW, Pople JA, Head-Gordon M (1989) Chem Phys Lett 157(6): 479–483
34. Gordon MS, Freitag MA, Bandyopadhyay P, Jensen JH, Kairys V, Stevens WJ (2001) J Phys Chem A 105(2):293–307
35. Jensen JH, Day PN, Gordon MS, Basch H, Cohen D, Garmer DR, Kraus M, Stevens WJ (1994) Modeling the Hydrogen Bond 569:139–151
36. Gordon MS, Slipchenko LV, Li H, Jensen JH (2007) Ann Rep Comp Chem 3:177–193
37. Adamovic I, Freitag MA, Gordon MS (2003) J Chem Phys 118(15):6725–6732
38. Song J, Gordon MS unpublished
39. Stone AJ (1981) Chem Phys Lett 83(2):233–239
40. Stone AJ (1996) The theory of intermolecular forces, Oxford University Press, Oxford
41. Freitag MA, Gordon MS, Jensen JH, Stevens WJ (2000) J Chem Phys 112(17):7300–7306
42. Slipchenko LV, Gordon MS (2007) J Comput Chem 28(1):276–291
43. Slipchenko LV, Gordon MS unpublished results
44. Jensen JH, Gordon MS (1996) Mol Phys 89(5):1313–1325
45. Jensen JH, Gordon MS (1998) J Chem Phys 108(12):4772–4782
46. Adamovic I, Gordon MS (2005) Mol Phys 103(2–3):379–387
47. Tang KT, Toennies JP (1984) J Chem Phys 80(8):3726–3741
48. William H Press SAT, William T Vetterling, Brian P Flannery (2002) Numerical recipes: the art of scientific computing, Cambridge University Press, Cambridge
49. Gordon MS, Schmidt MW (2005) In: Dykstra CE, Frenking G, Kim KS, Scuseria GE (eds) Theory and applications of computational chemistry, Ch 41, Elsevier, Amsterdam

50. Schmidt MW, Baldridge KK, Boatz JA, Elbert ST, Gordon MS, Jensen JH, Koseki S, Matsunaga N, Nguyen KA, Su SJ, Windus TL, Dupuis M, Montgomery JA (1993) J Comput Chem 14(11): 1347–1363
51. Woon DE, Dunning TH (1993) J Chem Phys 98(2):1358–1371
52. Hariharan PC, Pople JA (1973) Theoretica Chimica Acta 28(3):213–222
53. Krishnan R, Binkley JS, Seeger R, Pople JA (1980) J Chem Phys 72(1):650–654
54. Clark T, Chandrasekhar J, Spitznagel GW, Schleyer PV (1983) J Comput Chem 4(3):294–301
55. Li H, Netzloff HM, Gordon MS (2006) J Chem Phys 125(19):194103-1–194103-9
56. Plimpton S (1995) J Comput Phys 117(1):1–19
57. Narten AH (1976) J Chem Phys 65(2):573–579
58. Steinhauser O, Neumann M (1980) Mol Phys 40(1):115–128
59. Narten AH (1977) J Chem Phys 67(5):2102–2108
60. Sorenson JM, Hura G, Glaeser RM, Head-Gordon T (2000) J Chem Phys 113(20):9149–9161
61. Allesch M, Schwegler E, Gygi F, Galli G (2004) J Chem Phys 120(11):5192–5198
62. Netzloff HM, Gordon MS (2004) J Chem Phys 121(6):2711–2714

# CHAPTER 9

# FORMALISMS FOR THE EXPLICIT INCLUSION OF ELECTRONIC POLARIZABILITY IN MOLECULAR MODELING AND DYNAMICS STUDIES

PEDRO E. M. LOPES[1], EDWARD HARDER[2], BENOÎT ROUX[2], AND ALEXANDER D. MACKERELL, JR.[1]

[1] *Department of Pharmaceutical Sciences, School of Pharmacy, University of Maryland, 20 Penn Street, Baltimore, MD 21230, USA, e-mail: amackere@rx.umaryland.edu*
[2] *Department of Biochemistry and Molecular Biology, Center for Integrative Science, University of Chicago, Chicago, Illinois, 60637*

**Abstract:**    Current methodologies for modelling electronic polarization effects in empirical force fields are presented. Emphasis is placed on the mathematical details of the methods used to introduce polarizability, namely induced dipoles, Drude oscillators or fluctuating charge. Overviews are presented on approaches used to damp short range electrostatic interactions and on Extended Langrangian methods used to perform Molecular Dynamics simulations. The final section introduces the polarizable methods under development in the context of the program CHARMM

## 9.1.    INTRODUCTION

Molecular mechanical (MM) force fields (FF) are widely used in molecular modeling studies of systems with thousands on up to millions of atoms. To date they have proved to be surprisingly accurate in many applications despite their simplified functional forms. This accuracy is, to some extent surprising, due to the number of approximations in FFs, the biggest of which is typically the method by which the charge distribution of the molecules is treated. In additive force fields, which represents the bulk of current FFs [1, 2] this is done by assigning partial fixed charges to the atoms, thus creating a force field whose electrostatic properties are not capable of reacting to changes in the environment. Additive force fields common for biomolecular simulations [3–6] share the same general functional form:

$$U = U_{bond} + U_{LJ} + U_{elect} \tag{9-1a}$$

$$U_{elect} = \sum_{i=1}^{N} \sum_{j \neq i} \frac{q_i q_j}{r_{ij}} \tag{9-1b}$$

where $U_{bond}$ (**r**) represents the bonded terms (bonds, angles, dihedrals, etc), $U_{LJ}$ (**r**) represents the van der Waals (vdW) interactions and is typically treated as a 6–12 Lennard-Jones (LJ) term and $U_{elect}$ (**r**) is the electrostatic term of the Coulomb form (Eq. 9-1b). In recent years it has became that the additive treatment of electrostatic interactions is limiting and that electronic polarization will play a central role in the next generation of force fields for molecular simulations, including both molecular dynamics (MD) and Monte Carlo (MC) methods [7–9]. Accordingly, much effort is being devoted to the development of methods and parameters required for the implementation of polarizable force fields.

In nature, the charge distribution of a molecule can be significantly influenced or 'polarized' by its surroundings, a phenomenon that may be included in FFs via the explicit inclusion of electronic polarizability. A large motivating factor towards polarizable force fields is the fact that it is increasingly important to simulate heterogeneous environments, which requires that a given model is able to provide an environment-dependent response. For example, it is commonly accepted that modeling a water molecule with fixed point charges is not adequate to simultaneously describe bulk water molecules as well as water in highly hydrophobic environments [10, 11]. This is more critical if it is considered that a given molecule may visit these environments within the course of a single simulation. Therefore, the inclusion of molecular polarizability seems a basic requirement in order to develop force fields suitable for a wide range of hetergeneous environments.

Electronic polarizability is often included in force fields via the use of induced dipoles. Assuming that hyperpolarization effects are absent, the induced dipoles respond linearly relative to the electric field. In this case, the induced dipole $\boldsymbol{\mu}$ on an atom is the product of the total electric field **E** and the atomic polarizability tensor $\boldsymbol{\alpha}$.

$$\boldsymbol{\mu} = \boldsymbol{\alpha} \cdot \mathbf{E} \tag{9-2}$$

The total electric field, **E**, is composed of the external electric field from the permanent charges $\mathbf{E}^0$ and the contribution from other induced dipoles. This is the basis of most polarizable force fields currently being developed for biomolecular simulations. In the present chapter an overview of the formalisms most commonly used for MM force fields will be presented. It should be emphasized that this chapter is not meant to provide a broad overview of the field but rather focuses on the formalisms of the induced dipole, classical Drude oscillator and fluctuating charge models and their development in the context of providing a practical polarization model for molecular simulations of biological macromolecules [12–21]. While references to works in which the different methods have been developed and applied are included throughout the text, the major discussion of the implementation of these models focuses

on efforts in the context of the program CHARMM. For additional information the reader is referred to issue 6 of volume 3 of the Journal of Chemical Theory and Computation which was a "Special Issue on Polarization."

## 9.2. METHODS TO INCLUDE POLARIZATION IN CLASSICAL FORCE FIELDS

### 9.2.1. Induced Dipoles

One method for treating polarizability is the assignment of both partial atomic charges and induced dipoles on the atoms in a molecule. In its most common implementation in biomolecular simulations, inducible point dipoles are added to some or all atomic sites in the molecule [22–25]. An alternative methodology proposed by Allinger and co-workers is the use of bond dipoles [26].

The dipole moment, $\boldsymbol{\mu}_i$, induced on a site $i$ is proportional to the electric field at that site, $\mathbf{E}_i$. The proportionality constant is the polarizability tensor, $\boldsymbol{\alpha}_i$. The dipole feels an electric field both from the permanent charges of the system and from the other induced dipoles. The expression for $\boldsymbol{\mu}_i$ is

$$\boldsymbol{\mu}_i = \boldsymbol{\alpha}_i \cdot \mathbf{E}_i = \boldsymbol{\alpha}_i \cdot \left[ \mathbf{E}_i^0 - \sum_{j \neq i} \mathbf{T}_{ij} \boldsymbol{\mu}_j \right] \tag{9-3}$$

where $\mathbf{E}_i^0$ is the field from the permanent charges. (Note that permanent dipoles or higher multipoles, when present, contribute to $\mathbf{E}^0$.) The induced dipoles interact through the dipole field tensor, $\mathbf{T}_{ij}$

$$\mathbf{T}_{ij} = \frac{1}{r_{ij}^3} \left[ \mathbf{I} - 3 \frac{\mathbf{r}_{ij} \mathbf{r}_{ij}}{r_{ij}^2} \right] \tag{9-4}$$

where $\mathbf{I}$ is the identity matrix, $\mathbf{r}$ is the distance between $i$ and $j$, and x, y, and z are the Cartesian components of the vector between $i$ and $j$.

The inducing field responsible for the energy of the induced dipoles, $U_{ind}$, has contributions from three terms: the permanent or static field, $U_{stat}$, the induced dipole–induced dipole interaction, $U_{dip}$, and the polarization energy, $U_{pol}$,

$$U_{ind} = U_{stat} - U_{dip} - U_{pol} \tag{9-5}$$

The energy $U_{stat}$ is the interaction energy of the $N$ induced dipoles with the static electric field,

$$U_{stat} = - \sum_{i=1}^{N} \boldsymbol{\mu}_i \cdot \mathbf{E}_i^0, \tag{9-6}$$

the energy $U_{dip}$ represents the induced dipole–induced dipole interaction via the dipole field tensor, $\mathbf{T}_{ij}$

$$U_{elect} = \frac{1}{2} \sum_{i=1}^{N} \sum_{j \neq i} \boldsymbol{\mu}_i \cdot \mathbf{T}_{ij} \cdot \boldsymbol{\mu}_j \qquad (9\text{-}7)$$

and the polarization energy, $U_{pol}$,

$$U_{pol} = \frac{1}{2} \sum_{i=1}^{N} \boldsymbol{\mu}_i \cdot \mathbf{E}_i \qquad (9\text{-}8)$$

is defined as the energy cost needed to induce the dipoles [22, 27]. By using Eq. (9-3) the electric field can be replaced by $\boldsymbol{\alpha}_i^{-1} \cdot \boldsymbol{\mu}_i$, and $U_{pol}$ becomes

$$U_{pol} = \frac{1}{2} \sum_{i=1}^{N} \boldsymbol{\mu}_i \cdot \boldsymbol{\alpha}_i^{-1} \cdot \boldsymbol{\mu}_i \qquad (9\text{-}9)$$

where $\boldsymbol{\alpha}_i^{-1}$ is the inverse of the polarization tensor. All polarizable models in which dipole moments, charges, or other multipoles can be modified by their environment will have a polarization term corresponding to $U_{pol}$. Inserting the expressions for $U_{stat}$, $U_{elect}$ and $U_{pol}$ into Eq. (9-5) yields

$$U_{ind} = \sum_{i=1}^{N} \boldsymbol{\mu}_i \cdot \left[ -\mathbf{E}_i^0 + \frac{1}{2} \sum_{j=i} \mathbf{T}_{ij} \cdot \boldsymbol{\mu}_j + \frac{1}{2} \mathbf{E}_i \right] \qquad (9\text{-}10)$$

and using the relationship of Eq. (9-3) ($\mathbf{E}_i = \mathbf{E}_i^0 - \sum \mathbf{T}_{ij} \cdot \boldsymbol{\mu}_j$), reduces the dependence of $U_{ind}$ to a function of the static field, $\mathbf{E}^0$,

$$U_{ind} = -\frac{1}{2} \sum_{i=1}^{N} \boldsymbol{\mu}_i \cdot \mathbf{E}_i^0 \qquad (9\text{-}11)$$

It is noteworthy that the induced energy is the dot product of the induced dipole and the static field and not the total field. The interpretation of Eq. (9-11) is that a static field is required to originate induced dipoles.

Interesting properties of the induced dipole polarizable model can be derived by simple mathematical manipulation. A particularly important one relates the minimum of the energy with converged values of the induced dipole. By combining Eqs. (9-9) and (9-10), the induction energy can be rewritten as

$$U_{ind} = -\sum_{i=1}^{N} \boldsymbol{\mu}_i \cdot \mathbf{E}_i^0 + \frac{1}{2}\sum_{i=1}^{N}\sum_{j\neq i} \boldsymbol{\mu}_j \cdot \mathbf{T}_{ij} \cdot \boldsymbol{\mu}_j + \frac{1}{2}\sum_{i=1}^{N} \boldsymbol{\mu}_j \cdot \boldsymbol{\alpha}_i^{-1} \cdot \boldsymbol{\mu}_i \qquad (9\text{-}12)$$

and the derivative of $U_{ind}$ with respect to the induced dipoles is

$$\nabla_{\boldsymbol{\mu}_i} U_{ind} = -\mathbf{E}_i^0 + \sum_{j\neq i} \mathbf{T}_{ij} \cdot \boldsymbol{\mu}_j + \boldsymbol{\alpha}_i^{-1} \cdot \boldsymbol{\mu}_i = 0 \qquad (9\text{-}13)$$

Taking into consideration Eq. (9-3), makes the derivative in Eq. (9-13) zero because $\boldsymbol{\alpha}_i^{-1} \cdot \boldsymbol{\mu}_i = \mathbf{E}_i^0 - \sum \mathbf{T}_{ij} \cdot \boldsymbol{\mu}_j$. The converged values of the induced dipoles are those that minimize the energy. Consecutive (iterative) adjustments of $\boldsymbol{\mu}_i$ are referred to as self-consistent field (SCF) calculations and represent a systematic optimization of the polarization degrees of freedom to attain values of those degrees of freedom that minimize the energy. More details are presented in Section 9.4. Other polarizable models also have auxiliary variables, analogous to $\boldsymbol{\mu}$, which are adjusted to minimize the energy in similar ways.

An approximation to the induced point dipole model that uses induced charges was proposed by Ferenczy and Reynolds [28–36]. This induced charge method involves point charges only and those depend on the environment. In this sense the method is related to the fluctuating charge model (see below). It is based on the idea of representing a series of multipole moments by several lower rank multipole moments on neighboring sites (e.g. representing a dipole by two individual monopoles) [29, 30, 37, 38]. Such a model was shown to be efficient in accounting for electrostatic interactions, and preliminary extensions to polarization using point charges have been described [32, 34]. The method shares similarities with the approaches of Zhu et al. [39] and Sprik [15] but without the requirements for a regular geometry or a molecular dynamics implementation, respectively. The method was systematically extended so that both the polarization energy and its derivatives can be determined.

More recently, Ponder and co-workers [40–48] developed the AMOEBA force field based on a modification of the formulation of Applequist [49] and Thole [50]. It is based on a modification of Eq. (9-3) with the static electric field due to permanent charges replaced by permanent multipoles:

$$\boldsymbol{\mu}_i = \boldsymbol{\alpha}_i \cdot \mathbf{E}_i = \boldsymbol{\alpha}_i \cdot \left[ \sum_{j\neq i} \mathbf{T}_{ij}^{\alpha} \mathbf{M}_j + \sum_{k\neq i} \mathbf{T}_{ik}^{\alpha\beta} \boldsymbol{\mu}_k \right] \qquad (9\text{-}14)$$

where $\mathbf{M}(M_i = \left(q_i, \mu_{i,x}, \mu_{i,y}, \mu_{i,z}, Q_{i,xx}, Q_{i,xy}, Q_{,xz} \cdots Q_{i,zz}\right)^T)$ is the vector of permanent atomic multipole components, up to quadrupole, and $\mathbf{T}$ is the interaction matrix, defined previously.

In the AMOEBA force field the "permanent" atomic multipoles are determined from QM calculations [40]. The prescription considers that the resulting multipoles on the atoms result from two components, the "permanent" and the "induced" moments:

$$\mathbf{M}_i = \mathbf{M}_i^{perm} + \mathbf{M}_i^{ind} \tag{9-15}$$

$\mathbf{M}_i^{ind}$ results from induction from all sites in the absence of an external field and is defined by an expression similar to Eq. (9-3)

$$\mathbf{M}_{i,\alpha}^{ind} = \alpha_i \cdot \left[ \sum_{j \neq i} \mathbf{T}_{ij}^{\alpha} \mathbf{M}_j^{perm} + \sum_{k \neq i} \mathbf{T}_{ik}^{\alpha} \mathbf{M}_k^{ind} \right] \tag{9-16}$$

Substituting Eq. (9-15) into Eq. (9-16) results in

$$\mathbf{M}_{i,\alpha}^{ind} = \alpha_i \cdot \left[ \sum_{j \neq i} \mathbf{T}_{ij}^{\alpha} \left( \mathbf{M}_j - \mathbf{M}_j^{ind} \right) + \sum_{k \neq i} \mathbf{T}_{ik}^{\alpha} \mathbf{M}_k^{ind} \right]. \tag{9-17}$$

In the case $j = k$, Eq. (9-17) simplifies greatly and becomes

$$\mathbf{M}_{i,\alpha}^{ind} = \alpha_i \sum_{j \neq i} \mathbf{T}_{ij}^{\alpha} \mathbf{M}_j \tag{9-18}$$

Otherwise, if $j \neq k$ and after rearrangement and introduction of different fractional scaling factors the following equation results:

$$\begin{aligned}
\mathbf{M}_{i,\alpha}^{ind} &= \alpha_i \cdot \left[ \sum_{Allsites} s_j^p \mathbf{T}_{ij}^{\alpha} \left( \mathbf{M}_j - \mathbf{M}_j^{ind} \right) + \sum_{Allsites} s_j^m \mathbf{T}_{ik}^{\alpha} \mathbf{M}_k^{ind} \right] \\
&= \alpha_i \cdot \left[ \sum_{Allsites} s_j^p \mathbf{T}_{ij}^{\alpha} \mathbf{M}_j + \sum_{Allsites} \left( s_j^m - s_j^p \right) \mathbf{T}_{ij}^{\alpha} \mathbf{M}_j^{ind} \right]
\end{aligned} \tag{9-19}$$

The scaling factor $s_j$ can take any value between 0 and 1 and is applied to site $j$. The superscripts $p$ and $m$ indicate permanent and mutual induction, respectively. Equation (9-19) can be solved iteratively using similar procedures to those used to solve Eq. (9-3). The formal "permanent" moments can be calculated by subtracting induced moments from moments from ab initio calculations. For any conformation of a given compound the atomic multipoles can be determined from Distributed Multipole Analysis (DMA) [51].

Calculation of the energy and forces acting on a molecular system requires knowledge of the magnitude of the inducible dipoles. The forces associated with the dipoles (spatial derivatives of the potential) [13], can be computed from Eq. (9-12), and on atomic site $k$ are

$$F_k = -\nabla U_{ind} = \sum_{i=1}^{N} \mu_i \nabla_k E_i^0 + \frac{1}{2} \sum_{i=1}^{N} \mu_i \cdot \sum_{j=1}^{N} \left( \nabla_k T_{ij} \right) \cdot \mu_j \tag{9-20}$$

All contributions to the forces from terms involving derivatives with respect to the dipoles are zero because of the condition in Eq. (9-13) [13, 52]. Attaining this condition is a necessary prerequisite for using induced dipoles, or any electronic polarization, in force field calculations.

The electric field that each dipole feels depends on all other induced dipoles requiring the use of a self-consistent method as mentioned above. To achieve this, Eq. (9-13) can be rearranged and written in matrix form considering that $\mathbf{A} = \boldsymbol{\alpha}^{-1}(\mathbf{I} - \boldsymbol{\alpha}\mathbf{T})$,

$$\mathbf{A} \cdot \boldsymbol{\mu} = \mathbf{E}^0 \tag{9-21}$$

The diagonal elements of the matrix $\mathbf{A}_{ii}$ are $\boldsymbol{\alpha}_i^{-1}$ and the off-diagonal elements of $\mathbf{A}_{ij}$ are $\mathbf{T}_{ij}$. Equation (9-21) determines how the dipoles are coupled to the static electric field. There are three major methods to determine the dipoles: matrix inversion, iterative methods and predictive methods.

The system in Eq. (9-21) can be solved with direct matrix inversion. Bernardo et al. [53] found the method more robust but significantly slower than the iterative procedure. For a system with $N$ dipoles, solving for each of them involves inverting the $N \times N$ matrix, A – an $O(N^3)$ operation that is typically much more computationally expensive to perform at each step of a $O(N)$ or $O(N^2)$ molecular dynamics simulation and this method has been used very rarely [53]. In the iterative method, the left-hand side of the Eq. (9-3) is calculated by substituting an initial guess for $\boldsymbol{\mu}$ into the right-hand side, and then the cycle is repeated until the desired level of self-consistency is achieved [52, 54, 55]. Both matrix inversion and iterative methods may be used in the SCF calculation to determine the induced dipoles following which the energies, forces as well as second and higher order derivatives acting on the molecular system may be determined. Such information is necessary for energy minimizations and molecular dynamics (MD) simulations. However, the requirement of computational accessibility allowing for MD simulations of large molecular systems requires special considerations when electronic polarizability is included in the model. These will be discussed below in Section 9.4.

### 9.2.2.    Classical Drude Oscillator Model

The models discussed in the previous section treat the polarization on each polarizable center using point dipoles. An alternative approach is to model the polarizable centers using dipoles of finite length, represented by an explicit pair of point charges. A variety of different models of polarizability have used this approach, but especially noteworthy are the Drude models (also known as "shell" or "charge on spring" models) frequently used in simulations of solid state ionic materials and recently extended to water and organic compounds [10, 11, 56–65]. Efforts in our laboratory are aimed at developing a classical Drude model based polarizable force field for biological macromolecules (see Section 9.5).

The Drude model can trace its origins to the work of Paul Drude in 1902 and was developed as a simple way to describe the dispersive properties of materials [66]. In the classical formalism, it represents electronic polarization by introducing a massless charged particle attached to each polarizable atom by a harmonic spring. The positions of these "auxiliary" particles are then adjusted self-consistently to their local energy minima for any given configuration of the atoms in the system. A quantum version of the model has been used in early applications to describe the dipole–dipole dispersion interactions [67–70]. A semiclassical version of the model was used more recently to describe molecular interactions [71], and electron binding [72]. The classical version of the model has been quite useful in statistical mechanical studies of dense systems and in recent decades has seen widespread use in MD and MC simulations. Examples of applications include ionic crystals [73–78], simple liquids of polarizable particles [63, 64, 79–83], liquid water [10, 11, 56, 84–87], and the hydration of small ions [61, 88]. In recent years, the Drude model was extended to interface with QM approaches for use in QM/MM methods [89].

A particularly attractive aspect of the Drude oscillator model is that it preserves the simple charge–charge Coulomb electrostatic interactions to treat polarizability and, therefore, may be implemented in standard biomolecular simulation programs in a relatively straightforward way. Despite this technical advantage, Drude oscillators have not been as widely used as the point dipole or charge transfer models, probably because of the difficulties designing efficient computational schemes to solve the fast motion of the auxiliary particles in MD simulations [56]. In contrast with the point dipole [90] and fluctuating charge models [17, 91], extended Lagrangian algorithms have only recently been implemented [12], as discussed below in Section 9.4. Prior to this applications of the Drude model to liquids usually used computationally costly direct SCF iterative schemes.

In the classical Drude polarizable model polarization is determined by a pair of point charges separated by a variable distance **d**. For a given atom with charge $q$ a mobile Drude particle (also referred to as a Drude oscillator in the text) carrying a charge $q_D$ is introduced. The charge on the parent atom is replaced by $q - q_D$ in order to preserve the net charge of the atom–Drude oscillator pair. The Drude particle is harmonically bound to the atomic particle with a force constant $k_D$. Thus, the Drude model can be described as consisting of an effective nuclear charge and a charge in the valence shell responsible for most of the polarization of the atom. The mathematical formulation of the Drude model is in fact an empirical method of representing the dipolar polarization of the site when $\|\mathbf{d}\| \rightarrow 0$.

In the absence of an electric field, the Drude particle coincides with the position of the atom, **r**, and the atom appears as a point charge of magnitude $q$. In the presence of a uniform electric field **E**, the Drude particle assumes a displaced position $\mathbf{r} + \mathbf{d}$. The Drude separation **d** is related to $k_D$, **E** and $q_D$:

$$\mathbf{d} = \frac{q_D \mathbf{E}}{k_D} \qquad (9\text{-}22)$$

The formula for the induced atomic dipole, $\boldsymbol{\mu}$ as a function of $\mathbf{d}$ is

$$\boldsymbol{\mu} = \frac{q_D^2 \mathbf{E}}{k_D} \tag{9-23}$$

From which results a simple expression for the isotropic atomic polarizability:

$$\alpha = \frac{q_D^2}{k_D} \tag{9-24}$$

As with any model involving inducible dipoles, the potential energy of the system contains terms representing the interaction with the static field, the interaction with other dipoles and the polarization energy (self energy), in addition to the standard representation of the bonding terms (bonds, angles, dihedrals, etc.) and intermolecular interactions, typically represented by a Lennard-Jones term in biological force fields.

$$U(\mathbf{r}, \mathbf{d}) = U_{bond}(\mathbf{r}) - U_{LJ}(\mathbf{r}) - U_{elect}(\mathbf{r}, \mathbf{d}) - U_{self}(\mathbf{d}) \tag{9-25}$$

The dependence on the nuclear positions is indicated by $\mathbf{r}$ and the dependence on the Drude positions is indicated by $\mathbf{d}$. In Eq. (9-25) $U_{bond}(\mathbf{r})$ is the intramolecular energy contribution from, typically, the bond lengths, valence angles, and dihedral angles, $U_{LJ}(\mathbf{r})$ is a Lennard-Jones "6–12" nonpolar contribution, $U_{elect}(\mathbf{r}, \mathbf{d})$ represents all Coulombic interactions, atom–atom, atom–Drude, and Drude–Drude, and $U_{self}(\mathbf{d})$ represents the atom–Drude harmonic bonds. The term $U_{self}(\mathbf{d})$ arises from the harmonic spring separating the two charges and has the simple expression

$$U_{self}(\mathbf{d}) = \frac{1}{2} \sum_{i=1}^{N} k_i d_i^2 \tag{9-26}$$

The electrostatic interaction between independent polarizable atoms is simply the sum of the charge–charge interactions between the four charge sites (i.e. two atoms and their respective Drude particles):

$$U_{elect}(\mathbf{r}, \mathbf{d}) = \frac{1}{2} \sum_{i=1}^{N} \sum_{j \neq i} q_i q_j \left[ \frac{1}{|\mathbf{r}_{ij}|} - \frac{1}{|\mathbf{r}_{ij} - \mathbf{d}_j|} + \frac{1}{|\mathbf{r}_{ij} + \mathbf{d}_i|} + \frac{1}{|\mathbf{r}_{ij} - \mathbf{d}_j + \mathbf{d}_i|} \right] \tag{9-27}$$

For a given $\boldsymbol{\alpha}$ the force constant $k_D$ can be chosen in a way that the displacement $\mathbf{d}$ of the Drude particle remains much smaller than the interatomic distance. This guarantees that the resulting induced dipole $\boldsymbol{\mu}$, is almost equivalent to a point dipole. In the Drude polarizable model the only relevant parameter is the combination $q_D^2/k_D$ which defines the atomic polarizability, $\alpha$. It is

noteworthy that the electrostatic interaction in the Drude model requires only the charge–charge terms already present in MD and Monte-Carlo simulation codes. No new interaction types, such as the dipole field tensor $\mathbf{T}_{ij}$ of Eq. (9-3) are required. The computational advantage of not having to compute the dipole–dipole interactions is balanced by the extra charge–charge calculations. In the case of the CHARMM implementation [10, 11, 58, 60, 63–65] the computational impact of the Drude model is minimized by hydrogen atoms not carrying Drude particles thereby diminishing number charge–charge interactions that need to be calculated.

The Drude oscillators are typically treated as isotropic on the atomic level. However, it is possible to extend the model to include atom-based anisotropic polarizability. When anisotropy is included, the harmonic self-energy of the Drude oscillators becomes

$$
\begin{aligned}
U_{\text{self}} &= \frac{1}{2}\mathbf{d} \cdot \left[ K^D \right] \cdot \mathbf{d} \\
&= \frac{1}{2} \left( \left[ K_{11}^D \right] d_1^2 + \left[ K_{22}^D \right] d_2^2 + \left[ K_{33}^D \right] d_3^2 \right)
\end{aligned}
\tag{9-28}
$$

Where the quantities $d_1^2$, $d_2^2$ and $d_3^2$ are the projections of the Drude displacement vector $\mathbf{d}$ on orthogonal axis defined on a local intramolecular reference frame. The intramolecular reference frame may be defined, for example, by the C$=$O vector and the N—C$=$O plane of an amide bond [65].

The term $U_{\text{elec}}$ of Eq. (9-27) corresponds to the sum over all Coulombic interactions between the core charges $q_i$, located at $\mathbf{r}_i$, and the Drude charges $q_i^D$, located at $\mathbf{r}_i^D = \mathbf{r}_i + \mathbf{d}$. The interactions of the various pairs of charges are treated according to the topological bonding order determined from the atoms in the molecule. As in standard fixed charge force fields, the interactions between core charges corresponding to 1–2 (neighbor) and 1–3 (next-neighbor) pairs are accounted by explicit bonding terms in the potential energy, $U_{\text{internal}}$, and necessarily excluded from the electrostatic energy. Similarly, the interactions of the Drude oscillators with core charges are excluded for 1–2 and 1–3 pairs. The Coulomb interactions between Drude oscillators corresponding to 1–2 and 1–3 atom pairs are present and screened by the function $f_1(r_{ij})$ [50]. The form of the screening function used in the Drude model of CHARMM is

$$
f_1(r_{ij}) = 1 - \left( 1 + \frac{(\alpha_i + \alpha_j)r_{ij}}{2\left(\alpha_i\alpha_j\right)^{1/6}} e^{-(a_i+a_j)r_{ij}/(\alpha_i\alpha_j)^{1/6}} \right)
\tag{9-29}
$$

where $r_{ij}$ is the distance between Drude charges, $\alpha_i$ is the trace of the atomic polarizability tensor, and the Thole damping parameters, $a_i$, modulate the screening strength of $f_1(r_{ij})$. The interactions involving all core charges and all Drude oscillators are included for all 1–4 pairs and beyond without screening. Additional details of the damping function are presented in Section 9.3.

### 9.2.3. Fluctuating Charges Model

Polarizability can also be introduced into standard potentials (Eq. 9-1) by allowing the values of the partial charges to respond to the electric field of their environment. This may be achieved by coupling the charges to their environment using electronegativity equalization (EE) or chemical potential equalization (CPE). This method for treating polarizability has been called the "fluctuating charge" method [17, 92], the "chemical potential (electronegativity) equalization" method [93–108], or "charge equilibration" method [109–115] and has been applied to a variety of systems [17, 116–126]. A practical advantage of this approach is that it introduces polarizability without introducing new interactions. Compared to the Drude model, this can be done using the same number of charge–charge interactions as would be present in a nonpolarizable simulation.

In the fluctuating charge method [109] variable discrete charges are located on atomic sites within the molecule. Their value is computed, for a given molecular geometry, by minimization of the electrostatic energy. The energy of a molecule, as well as a system comprised of a collection of molecules, can be described hierarchically starting from the energy of an isolated atom. Using a neutral atom as a reference point, the energy of an isolated atom can be expanded as second-order Taylor series of the charge [127]:

$$U^{atom}(q_\alpha) = U_\alpha(0) + \chi_\alpha^0 q_\alpha + \frac{1}{2} J_{\alpha\alpha}^0 q_\alpha^2 \tag{9-30}$$

The coefficients $\chi_\alpha^0$ and $J_{\alpha\alpha}^0$ have a clear physical interpretation. $\chi_\alpha^0$ is the "Mulliken electronegativity", namely one-half of the ionization energy and the electron affinity [128].

$$\chi_\alpha^0 = \frac{IP_\alpha + EA_\alpha}{2} \tag{9-31}$$

The value $\frac{1}{2} J_{\alpha\alpha}^0$ is called the "absolute hardness" and can be obtained as half of the difference between the ionization potential and the electron affinity [129].

$$\frac{1}{2} J_{\alpha\alpha}^0 = \frac{IP_\alpha - EA_\alpha}{2} \tag{9-32}$$

When atoms are brought together to form an isolated molecule, the molecular energy consist of contributions from the individual atoms plus intra-atomic interactions, and thus is a function of both charges and coordinates, $U^{molec}(\mathbf{r}_{\mu\nu}, q)$. In the following equations indices $\alpha$ and $\beta$ run over charged sites on a given molecule and $\mu$ and $\nu$ indicate atoms that comprise the molecule. The distinction is relevant since in many models (ex. TIP4P-FQ water model) charged and atomic sites are not coincident. Typically, when there is no external electrostatic field the potential can be partitioned into Coulombic and nonelectrostatic $V^{nonelec}(\mathbf{r}_{\mu\nu})$ types of interactions:

$$U^{molec}(\mathbf{r}, q) = \sum_{\alpha=1}^{N_{site}} \left[ \chi_\alpha^0 q_\alpha + \frac{1}{2} J_{\alpha\alpha}^0 q_\alpha^2 \right]$$

$$+ \sum_{\alpha=1}^{N_{site}} \sum_{\beta>\alpha}^{N_{site}} J_{\alpha\beta}(\mathbf{r}_{\alpha\beta}) q_\alpha q_\beta \qquad (9\text{-}33)$$

$$+ V^{nonelec}(\mathbf{r}_{\mu\nu})$$

In a multi-molecular system with $N_{molec}$ molecules and each molecule consisting of $N_{atom}$ atoms and $N_{site}$ charged sites, the total energy becomes

$$U^{system}(\mathbf{r}, q) = \sum_{i=1}^{N_{molec}} \sum_{\alpha=1}^{N_{site}} \left[ \chi_{i\alpha}^0 q_{i\alpha} + \frac{1}{2} J_{i\alpha i\alpha}^0 q_{i\alpha}^2 \right]$$

$$+ \sum_{i=1}^{N_{molec}} \sum_{\alpha=1}^{N_{site}} \sum_{\beta>\alpha}^{N_{site}} J_{i\alpha i\beta}(\mathbf{r}_{i\alpha i\beta}) q_{i\alpha} q_{i\beta} \qquad (9\text{-}34)$$

$$+ \sum_{i=1}^{N_{molec}} \sum_{j>1}^{N_{molec}} \sum_{\alpha=1}^{N_{site}} \sum_{\beta>\alpha}^{N_{site}} J_{i\alpha j\beta}(\mathbf{r}_{i\alpha j\beta}) q_{i\alpha} q_{j\beta}$$

$$+ V^{nonelec}(\mathbf{r}_{i\mu j\nu})$$

Equation (9-34) without the $V^{nonelec}(\mathbf{r}_{i\mu j\nu})$ term is denoted as $U^{elec}(\mathbf{r}, q)$:

$$U^{elect}(\mathbf{r}, q) = \sum_{i=1}^{N_{molec}} \sum_{\alpha=1}^{N_{site}} \left[ \chi_{i\alpha}^0 q_{i\alpha} + \frac{1}{2} J_{i\alpha i\alpha}^0 q_{i\alpha}^2 \right]$$

$$+ \sum_{i=1}^{N_{molec}} \sum_{\alpha=1}^{N_{site}} \sum_{\beta>\alpha}^{N_{site}} J_{i\alpha i\beta}(\mathbf{r}_{i\alpha i\beta}) q_{i\alpha} q_{i\beta} \qquad (9\text{-}35)$$

$$+ \sum_{i=1}^{N_{molec}} \sum_{j>1}^{N_{molec}} \sum_{\alpha=1}^{N_{site}} \sum_{\beta>\alpha}^{N_{site}} J_{i\alpha j\beta}(\mathbf{r}_{i\alpha j\beta}) q_{i\alpha} q_{j\beta}$$

The energy given by Eq. (9-35) replaces the Coulomb energy $q_i q_j / \mathbf{r}_{ij}$ in Eq. (9-1). The charges $q_i$ are now treated as independent variables, and the polarization response is determined by variations in the charge values. These charges depend on the interactions with other molecules as well as other charge sites on the same molecule, and will change for every time step or configuration sampled during a simulation. The charge values used for each configuration are, in principle, those that minimize the energy given by Eq. (9-35).

Charge conservation can be imposed in two ways. A charge neutrality constraint can be applied to the entire system, thus allowing charge to move from one atomic site to another until the electronegativities are equal on all the atoms of the system.

Alternatively, charge can be constrained independently on each molecule (or other subgroup), so that charge flows only between atoms on the same molecule until the electronegativities are equalized within each molecule, but not between distinct molecules [17]. At each molecular configuration $\mathbf{r}$, the charges are redistributed among the atoms so as to equalize the instantaneous electronegativities. In other words, the ground-state charge distribution $q^{eq}(\mathbf{r})$ satisfies the following equations

$$\left(\frac{\partial U^{system}}{\partial q_{i\alpha}}\right)_{q=q^{eq}} = \left(\frac{\partial U^{system}}{\partial q_{i\beta}}\right)_{q=q^{eq}} \tag{9-36}$$

or equivalently

$$\left(\frac{\partial U^{elect}}{\partial q_{i\alpha}}\right)_{q=q^{eq}} = \left(\frac{\partial U^{elect}}{\partial q_{i\beta}}\right)_{q=q^{eq}} \tag{9-37}$$

and the ground-state energy $U_g(\mathbf{r}) = U^{system}(q^{eq}, \mathbf{r})$. The same set of coupled linear equations can be obtained by variationally minimizing the total energy $U^{system}(q, \mathbf{r})$ with respect to the charge distribution, subject to the constraints of molecular charge conservation,

$$\sum_{\alpha} q_{i\alpha} = 0 \tag{9-38}$$

or the constraint of overall charge conservation,

$$\sum_{i}\sum_{\alpha} q_{i\alpha} = 0 \tag{9-39}$$

An analogy may be drawn between the ground state Born-Oppenheimer potential energy and $U_g(\mathbf{r})$. The charges, $q$, in this case play a role similar to the electronic wave functions, and they represent the electronic degrees of freedom. In fact, Eq. (9-37) is the principle of chemical potential equalization [130] in density functional electronic structure theory [131] reformulated for the point-charge representation. In the adiabatic limit, the nuclei evolve dynamically on the Born-Oppenheimer potential energy surface. In the same fashion, the principle of electronegativity equalization requires the charge density to fluctuate adiabatically as the nuclear coordinates evolve in time.

In most cases, charge is conserved for each molecule, so there is no charge transfer between molecules. In quantum mechanics, charge transfer is an important part of the interaction energy, so there are reasons to remove this constraint [132–136]. Unfortunately, many fluctuating charge methods are known to suffer from a superlinear scaling of the polarizability with increasing molecular size that penalizes model transferability and prevents application of these methods to large molecules of biological importance [107, 137]. The principal cause of this failure is the fact that in the traditional fluctuating charge models, charge can flow between covalent bonds at small energetic cost, thus covering large portions of the molecule. This method may

be suited for small molecules, highly conjugated large size $\pi$ systems like polyenes, but is in general not accurate for bigger systems, as charge will be able to flow to far regions of the extended molecules.

One solution was developed based on the concept of atom–atom charge transfer (AACT) [137]. In this approach the energy is Taylor expanded in terms of charges transferred between atomic pairs within the molecule, rather than in terms of the atomic charges themselves. Similar in spirit is the bond-charge increment (BCI) [120, 138] model, which only allows for charge to flow between two atoms that are directly bonded to each other, guaranteeing that the total charge of each set of bonded atoms is conserved. In a related effort the Atom-Bond Electronegativity Equalization Method (ABEEM) [139–143] has been developed based in concepts from density functional theory. In this model, the total electronic energy of a molecule in the ground state is a complex function of different quantities: (a) valence-state chemical potential of atom $a$, bond $a - b$ and lone-pairs, (b) valence-state hardness of atom $a$, bond $a - b$, and lone-pairs, (c) partial charges of atom $a$, bond $a - b$, and lone-pairs, (d) distances between the different components, atoms, bonds and lone-pairs. The combination of ABEEM and molecular mechanics was performed bringing ABEEM charges of atoms, bonds, and lone-pair electrons into the intermolecular electrostatic interaction term in molecular mechanics [144, 145].

## 9.3.    DAMPING FUNCTIONS IN POLARIZABLE FORCE FIELDS

The development of the methods described in Section 9.2 was an important step in modeling polarization because it led to accurate calculations of molecular polarizability tensors. The most serious issue with those methods is known as the "polarization catastrophe" since they are unable to reproduce the substantial decrease of the total dipole moment at distances close to contact as obtained from ab initio calculations. As noted by Applequist et al. [49], and Thole [50], a property of the unmodified point dipole is that it may originate infinite polarization by the cooperative interaction of the two induced dipoles in the direction of the line connecting the two. The mathematical origins of such singularities are made more evident by considering a simple system consisting of two atoms (A and B) with isotropic polarizabilities, $\alpha_A$ and $\alpha_B$. The molecular polarizability, has two components, one parallel and one perpendicular to the bond axis between A and B,

$$\alpha_{||} = (\alpha_A + \alpha_B + 4\alpha_A\alpha_B/r^3)/(1 - 4\alpha_A\alpha_B/r^6) \qquad (9\text{-}40a)$$

$$\alpha_{\perp} = (\alpha_A + \alpha_B - 2\alpha_A\alpha_B/r^3)/(1 - \alpha_A\alpha_B/r^6) \qquad (9\text{-}40b)$$

When the distance $r$ between the two points approaches $4(\alpha_A\alpha_B)^{1/6}$, the parallel component $\alpha_{||}$ goes to infinity and will be negative for shorter distances.

The singularities can be avoided by making the polarizabilities sufficiently small so that at the typical distances between the atoms ($> 1$ Å) the factor $4\alpha_A\alpha_B/r^6$ is always less than one. The Applequist polarizabilities are in fact small compared to ab

initio values [50] Applequist's atomic polarizabilities were selected to optimize the molecular polarizabilities for a set of 41 molecules. Careful choice of polarizabilities can move the singularities in Eqs. (9-40a) and (9-40b) to small distances, but not eliminate them altogether, thus causing problems for simulation techniques such as MC, which tend to sample these nonphysical regions of configuration space.

While nonbonded atom pairs will typically not come within 1 Å of each other, it is possible for covalently bound pairs, either directly bounds, as in 1–2 pairs, or at the vertices of an angle, as in 1–3 pairs. Accordingly it may be considered desirable to omit the 1–2 and 1–3 dipole–dipole interactions as is commonly performed on additive force fields for the Coulombic and van der Waals terms. However, it has been shown that inclusion of the 1–2 and 1–3 dipole–dipole interactions is required to achieve anistropic molecular polarizabilites when using isotropic atomic polariz-abilites [50]. For example, in a Drude model of benzene in which isotropic polariza-tion was included on the carbons only inclusion of the 1–2 and 1–3 dipole–dipole interactions along with the appropriate damping of those interactions allowed for reproduction of the anisotropic molecular polarizability of the molecule [64]. Thus, it may be considered desirable to include these short range interactions in a polarizable force field.

Another way to deal with this limitation is the inclusion of electrostatic damp-ing, which can be achieved through the method outlined by Thole [50, 146], that uses charge distributions instead of point charges. The screening (attenuation) of the dipole–dipole interaction can be physically interpreted as correcting for the fact that the electronic distribution is not well represented by point charges and point dipoles at small distances [49, 50, 147]. Mathematically, screening avoids the singularities such as those in Eqs. (9-40a) and (9-40b). Basically, Eqs. (9-2), (9-3) and (9-4) retain their significance with the only change being that both the electric field created by fixed charges is damped by the function $f_1(r)$, that effectively changes the contri-bution to **E** from the charge at $j$ [53, 148–151], and, when created by point dipoles, by $f_2(r)$.

$$\mathbf{E} = f_1(\mathbf{r})q\frac{\mathbf{r}}{r^3} \tag{9-41}$$

$$\mathbf{T} = f_2(\mathbf{r})3\frac{\mathbf{r}\cdot\mathbf{r}}{r^5} - f_1(\mathbf{r})\frac{\mathbf{I}}{r^3} \tag{9-42}$$

The method of Thole was developed with the help of the induced dipole formu-lation, when all dipoles interact through the dipole field tensor. The modification introduced by Thole consisted in changing the dipole field tensor:

$$
\begin{aligned}
\left(T_{pq}\right)_{ij} &= \delta_{ij}r^{-3} - 3x_i x_j r^{-5} \\
&= \left(\alpha_p\alpha_q\right)^{-1/2}\left(\delta_{ij}\mu^{-3} - 3u_i u_j u^{-5}\right) \\
&= \left(\alpha_p\alpha_q\right)^{-1/2} t_{ij}\left(\mathbf{u}\right)
\end{aligned}
\tag{9-43}
$$

where $\mathbf{u} = \mathbf{x}/(\alpha_p \alpha_q)^{-1/6}$ and $\delta_{ij}$ is the Kronecker delta. T is a shape function that does not depend on $p$ or $q$, but is related to a model charge distribution $\rho(u)$. Thole originally investigated various forms for the charge distribution and two were considered:

$$\rho(u) = \begin{cases} \dfrac{3}{\pi} \dfrac{(a-u)}{a^4} & u < a \\ 0 & u \geq a \end{cases} \quad \text{(linear)} \tag{9-44a}$$

$$\rho(u) = \left(\dfrac{a^3}{8\pi} e^{-au}\right) \text{(exponential)} \tag{9-44b}$$

with associated dipole–dipole tensors:

$$\begin{aligned} T_{ij} &= \left(4v^3 - 3v^4\right)\delta_{ij}/r^3 - 3v^4\left(r_i r_j/r^5\right) & r < s \\ T_{ij} &= \delta_{ij}/r^3 - 3r_i r_j/r^5 & r \geq s \\ s &= a(\alpha_p \alpha_q)^{1/6} & v = r/s \end{aligned} \quad \text{(linear)} \tag{9-45a}$$

and

$$T_{ij} = \left[1 - \left(a^2 r^2/2 + ar + 1\right)e^{-ar}\right]\delta_{ij}/r^3 - 3\left[1 - \left(a^3 r^3/6 + a^2 r^2/2 + ar + 1\right)e^{-ar}\right]r_i r_j/r^5 \quad \text{(exponential)} \tag{9-45b}$$

Thole's polarizability parameters were selected to optimize the molecular polarizabilities for a set of 16 molecules. The method was later expanded to fit 52 molecules [146]. It must be emphasized that this electric-field damping method is totally independent of the polarization scheme used. For the Drude and fluctuating charge methods only $f_1(r)$ is required, whereas for methods based on induced dipoles both $f_1(r)$ and $f_2(r)$ are necessary. In the context of the induced dipole model other models were proposed since the formula of Thole does not provide enough attenuation. For example, in Ref. [152] the field is evaluated using

$$S(r_{ij}) = 1 - \exp\left[-\gamma\left(\dfrac{r_{ij}}{s_{ij}}\right)^n\right] \tag{9-46}$$

for all values of $r_{ij}$ and the constants $\gamma$ and $n$ determine the extent of attenuation.

## 9.4.    MOLECULAR DYNAMICS WITH POLARIZABLE FORCE FIELDS

In general MD simulations are performed via integrating Newton's equations of motion using a variety of integrators as previously described [153–155]. The concepts for MD simulations are similar to all methods of describing polarization discussed

in Section 9.2. In the following sections details of the implementations for different methodologies will be addressed.

### 9.4.1. Molecular Dynamics Using Induced Dipoles

In the case of the induced dipole methods, computation of the forces acting on the particles of the system requires self consistent determination of the inducible dipoles. A mathematically elegant way of doing it is by performing a SCF calculation to obtain the polarization contribution to the energy and forces from which the new atomic positions are determined. The computational requirement of this step can be greatly decreased by using a good initial guess for the SCF calculation and an initial guess for the electric field is typically obtained from the static field, $\mathbf{E}^0$ and the dipoles from the previous time step of the MD simulation [52, 156]. The iterative SCF calculation is then performed. Convergence limits on the dipoles reported in the literature range from $1 \times 10^{-2}$ D to $1 \times 10^{-6}$ D and have been made more strict over time [14, 157–161]. This is necessary as MD simulations require very strict convergence limits due to poor energy conservation [162]. While iterative methods can achieve the required level of convergence, the procedure is still CPU time intensive in cases that the system is large. (For example, in a simulation of 4000 polarizable sites the matrix $\mathbf{A}$ in Eq. (9-21) will have dimensions of 12,000 * 12,000 = 144,000,000 elements and the CPU time needed to achieve convergence will be quite significant.

Predictive methods that calculate $\mathbf{u}$ for the next time step of a MD simulation based on information from previous timesteps have been developed to minimize the computational cost. Ahlström et al. [13] used a first-order predictor algorithm, in which values of $\mathbf{u}$ from the two previous times steps are used to determine $\mathbf{u}$ at the next time step. A very serious drawback of this method is that it is not stable for long simulation times. However, it has been combined with iterative solutions, either by providing the initial iteration of the electric field values [163, 164], or by performing an iterative SCF step less frequently than every step [13, 165]. Higher-order predictor algorithms have also been described in the literature [13, 163, 166].

Another approach to minimize computational costs was a simplification of the iterative method proposed by Kaminski et al. [167]. The method consists in truncating the iterative SCF process after the second iteration. Equation (9-47) shows the process which is actually the initial iterations of the full iterative process.

$$\boldsymbol{\mu}_i^{1st} = \alpha_i \mathbf{E}_i^0 \tag{9-47a}$$

$$\boldsymbol{\mu}_i^{2nd} = \alpha_i \mathbf{E}_i^0 + \alpha_i \sum_{j \neq i} \mathbf{T}_{ij} \boldsymbol{\mu}_i^{1st} \tag{9-47b}$$

In the first-order approximation (Eq. 9-47a) the magnitude of the inducible dipoles is determined based on the assumption that they cannot interact with each other at all. The second-order approximation from Eq. (9-47b) is able to retain a greater part

of the dipole–dipole interactions than the first-order approximation while providing the benefits of reduced computational cost compared to the full iterative method. It should be noted that the second-order approximation in Eq. (9-47b) does not have a direct physical meaning and can be viewed as introducing a set of induced dipoles with magnitudes calculated on the assumption that each of them perceive all the other dipoles as if those other dipoles were induced by the electrostatic field of the permanent charges only.

An important advance in making explicit polarizable force fields computationally feasible for MD simulation was the development of the extended Lagrangian methods. This extended dynamics approach was first proposed by Sprik and Klein [91], in the sipirit of the work of Car and Parrinello for ab initio MD dynamics [168]. A similar extended system was proposed by van Belle et al. for inducible point dipoles [90, 169]. In this approach each dipole is treated as a dynamical variable in the MD simulation and given a mass, $M_\mu$, and velocity, $\dot{\boldsymbol{\mu}}$. The dipoles thus have a kinetic energy, $\sum_i M_\mu (\dot{\mu})^2/2$, and are propagated using the equations of motion just like the atomic coordinates [90, 91, 170, 171]. The equation of motion for the dipoles is

$$M_\mu \ddot{\boldsymbol{\mu}} = -\nabla \boldsymbol{\mu}_i = \mathbf{E}_i - \boldsymbol{\alpha}_i^{-1} \cdot \boldsymbol{\mu}_i. \qquad (9\text{-}48)$$

The dipole mass does not correspond to any physical mass and is chosen based on the requirement of having approximately adiabatic fictitious dynamics for the duration of the MD simulation. This is satisfied with small values of $M_\mu$. The timestep of the simulation is chosen such that the frequency of the fictitious modes is much larger than the fastest nuclear frequency in order to keep the coupling between fictitious and real dynamics low, which is achieved by using small integration timesteps [172]. It is desirable to keep the kinetic energy of the dipoles small so that the dipole degrees-of-freedom are cold and near the potential energy minimum (corresponding to the exact solution of Eq. 9-3) such that it approximates the SCF condition.

Because this method avoids iterative calculations to attain the SCF condition, the extended Lagrangian method is a more efficient way of calculating the dipoles at every time step. However, polarizable point dipole methods are still more computationally intensive than nonpolarizable simulations. Evaluating the dipole–dipole interactions in Eqs. (9-7) and (9-20) is several times more expensive than evaluating the Coulombic interactions between point charges in Eq. (9-1). In addition, the requirement for a shorter integration timestep as compared to an additive model increases the computational cost.

Most liquid phase molecular simulations with explicit atomic polarizabilities are performed with MD rather than MC techniques. This is due to the fact that, despite its general computational simplicity, MC with explicit polarization [173, 174] requires that Eq. (9-21) be solved every MC step, when even one molecule in the system is moved, and the number of configurations in an average Monte Carlo computation is orders of magnitude greater than in a MD simulation. For nonpolarizable, pairwise-additive models, MC methods can be efficient because only the

interactions involving the moved particle need to be recalculated (while the other $(N-1) \times (N-1)$ interactions are unchanged). For polarizable models, all $N \times N$ interactions must be calculated, in principle, when one particle moves. Consequently, exact polarizable MC calculations can be two to three orders of magnitude slower than comparable nonpolarizable calculations [175]. This is true for all models that account for explicit polarization. Thus, employing MC becomes much less practical for polarizable systems, even though it might be otherwise preferable. The induced dipole polarizable model has, nevertheless, been used in MC simulations with single particle moves [158, 159, 167, 176–181]. To make the update of the induced dipoles at each step more efficient the distances between all the particles may be stored, since most of them are unchanged. However, the memory usage is high. Alternatively, various approximate methods, involving incomplete convergence or updating only a subset of the dipoles, have been suggested [176]. Unfortunately, these methods result in significant errors in computed physical properties [175, 178].

One final point concerns the long-ranged nature of the interactions in induced dipole based models. Dipole–dipole and dipole–charge interactions are termed long-ranged because they do not decrease faster than volume grows – i.e., as $r^3$. If periodic boundary conditions are used, which is common on MD simulations of biomolecular systems, some treatment of the long-ranged interactions is needed. All models, whether polarizable or not, face this problem, but for polarizable models this is a more significant issue. The use of cut-offs or other truncation schemes will change both the static field and the dipole field tensor. These changes to the electric field will modify the value of the induced dipoles, which in turn will change the field at other sites. Accordingly, the treatment of long-ranged forces feeds back on itself in a way that does not occur with nonpolarizable models. Nymand and Linse [182] showed that different boundary conditions (including Ewald sums, spherical cut-off and reaction field methods) lead to more significant differences in equilibrium, dynamical and structural properties for polarizable water models than for nonpolarizable models. It is thus crucial to treat the long-ranged interactions as accurately as possible in polarizable simulations. The most complete treatment of the long-ranged forces is the Ewald summation technique [161, 162, 183, 184]. Faster-scaling methods, such as the fast multipole and particle mesh algorithms, have also been extended to the treatment of point dipoles [161, 185–188].

### 9.4.2. Molecular Dynamics Using the Classical Drude Oscillator

Several different methods exist for treating the motion of the polarizable degrees of freedom in dynamic simulations with the Drude model. Similarly to models based on induced dipoles, there are iterative adiabatic techniques and fully dynamic methods based on extended Lagrangians. In the adiabatic methods, it is assumed that a correspondence between the Drude charge and the electronic degrees of freedom obeying the Born-Oppenheimer approximation exists. The heavier and slowly-moving nuclei and atomic core charges are considered to move adiabatically in the field generated

by the Drude charges. The positions of the Drude particles are assumed to react instantaneously in response to the motion of the nuclei, and thus always remain in positions in which they experience no net force. Those are the positions that minimize the total energy of the system attained via an SCF calculation as discussed above. The forces on the atomic charges are then used to propagate the dynamics, using standard numerical integration methods. The other alternative is to treat the charges dynamically, allowing them to occupy positions away from the minimum-energy position dictated by the nuclei, and thus experience non-zero forces.

When the Drude particles are treated adiabatically, a SCF method must be used to solve for the displacements of the Drude particle, $\mathbf{d}$, similarly to the dipoles $\boldsymbol{\mu}_i$ in the induced dipole model. The implementation of the SCF condition corresponding to the Born-Oppenheimer approximation is straightforward and the real forces acting on the nuclei must be determined after the Drude particles have attained the energy minimum for a particular nuclear configuration. In the case of $N$ polarizable atoms with positions $\mathbf{r}$, the relaxed Drude particle positions $\mathbf{r} + \mathbf{d}^{SCF}$ are found by solving

$$\frac{\partial U}{\partial \mathbf{d}_i} = 0 \tag{9-49}$$

where index $i$ runs from 1 to $N$. $U_{\text{bond}}$ and $U_{\text{LJ}}$ are independent of $\mathbf{d}$, and

$$\frac{\partial U_{self}}{\partial \mathbf{d}_i} + \frac{\partial U_{elec}}{\partial \mathbf{d}_i} = 0. \tag{9-50}$$

These equations define the force equilibria on the Drude particles

$$k_D \mathbf{d}_i - q_D \mathbf{E}_i = 0 \tag{9-51}$$

where $\mathbf{E}_i$ is the total electric field in $\mathbf{r} - \mathbf{d}$, arising from the fixed charges as well as all the induced dipoles (modeled with Drude oscillators). For atomic positions $\mathbf{r}$, the relaxed displacements produce the potential

$$U^{SCF}(\mathbf{r}) = U(\mathbf{r}, \mathbf{d}) \tag{9-52}$$

and the atomic motions in the SCF regime are described by

$$m_i \ddot{\mathbf{r}}_i = -\frac{\partial U\left(\mathbf{r}, \mathbf{d}^{SCF}\right)}{\partial \mathbf{r}_i} \tag{9-53}$$

Integrating Eq. (9-53) in MD simulations requires that the positions of the Drude particles be set at their energy minimum at every integration time step, by solving Eq. (9-51). This simple simulation method has been widely used in MD [75, 76, 87, 189] and to a lesser extent in MC simulations, although examples of application exist [171]. Nonetheless, the SCF procedure is limited and computationally expensive, because any nonconverged SCF calculation introduces systematic drag

forces on the physical atoms that considerably affect energy conservation and the stability of the temperature [56, 76, 190]. Depending on the convergence criterion used, these iterative methods typically require between three and ten iterations [56, 57, 76, 87, 190, 191]. It is interesting to compare the Drude model with models based on induced dipoles. The induced energy of the Drude model system can be written as

$$U_{ind}\left(\mathbf{r}_i, \mathbf{d}\right) = \sum_{i=1}^{N} \left\{ \frac{1}{2}k_i d_i^2 + q_i \left[ \mathbf{r}_{ii} \cdot \mathbf{E}_i^0 - (\mathbf{r}_{ii} + \mathbf{d}_i) \cdot \mathbf{E}_i^0 \right] + \left[ \frac{1}{2} \sum_{i=1}^{n} \sum_{j \neq i} \right. \right.$$
$$\left. \left. \times \left( \frac{1}{\mathbf{r}_{ij}} - \frac{1}{|\mathbf{r}_{ij} - \mathbf{d}_j|} + \frac{1}{|\mathbf{r}_{ij} + \mathbf{d}_i|} + \frac{1}{|\mathbf{r}_{ij} - \mathbf{d}_j + \mathbf{d}_i|} \right) \right] \right\} \quad \text{(9-54)}$$

This equation is the equivalent of Eq. (9-12) for the induced dipole model but has one important difference. Equation (9-13), the derivative of Eq. (9-12), is linear and standard matrix methods can be used to solve for the $\boldsymbol{\mu}_i$ because Eq. (9-12) is a quadratic function of $\boldsymbol{\mu}_i$, while Eq. (9-54) is not a quadratic function of $\mathbf{d}$ and thus matrix methods are usually not used to find the Drude particle displacements that minimize the energy.

An important alternative to SCF is to extend the Lagrangian of the system to consider dipoles as additional dynamical degrees of freedom as discussed above for the induced dipole model. In the Drude model the additional degrees of freedom are the positions of the moving Drude particles. All Drude particles are assigned a small mass $m_{D,i}$, taken from the atomic masses, $m_i$, of their parent atoms and both the motions of atoms and Drude particles (at positions $\mathbf{r}_i$ and $\mathbf{r}_{D,i} \equiv \mathbf{r}_i + \mathbf{d}_i$) are propagated

$$\left(m_i - m_{D,i}\right) \ddot{\mathbf{r}}_i = -\frac{\partial U}{\partial \mathbf{r}_i} \quad \text{(9-55)}$$

$$m_{D,i} \ddot{\mathbf{r}}_{D,i} = -\frac{\partial U}{\partial \mathbf{r}_{D,i}} \quad \text{(9-56)}$$

The motion of Drude particles is expected to be decoupled from the atomic motion if $m_D$ is sufficiently small. The obvious drawback is that a small $m_D$ requires a small integration time step [77, 85]. For a single Drude oscillator, a significant speedup can be attained by using a multi-timestep integration approach [118, 119], but this advantage is lost for a dense system of polarizable atoms, because the long-range $1/r^3$ dipole–dipole interactions include high-frequency oscillations and have to be integrated using very short time steps. However, even if $m_D$ is very small, the Drude particles will eventually reach a thermal equilibrium with the rest of the system. Therefore, simulation approaches relying solely on the kinetic decoupling of the Drude oscillators to maintain a Born-Oppenheimer regimen are inappropriate for long simulation runs. To overcome this the long thermalization time can be exploited to remain close to the SCF energy surface by periodically resetting the positions

of the Drude oscillators to their energy minimum [118], but doing so makes the simulation irreversible.

From this point of view it is of interest to examine the consequences of full thermalization of the classical Drude oscillators on the properties of the system. This is particularly important given the fact that any classical fluctuations of the Drude oscillators are *a priori* unphysical according to the Born-Oppenheimer approximation upon which electronic induction models are based. It has been shown [12] that under the influence of thermalized (hot) fluctuating Drude oscillators the corrected effective energy of the system, truncated to two-body interactions is

$$U^{eff}(\mathbf{r}) = U^{SCF}(\mathbf{r}) - \frac{3}{2}k_B T \sum_{i=1}^{N} \sum_{j \neq i} \frac{\alpha_i \alpha_j}{r_{ij}^6} + \cdots . \tag{9-57}$$

In addition to the static induction effects included in $U_{SCF}$, the hot Drude oscillators give rise to a $1/r^6$, temperature-dependent, attractive term. This $\frac{3}{2}k_B T \alpha^2 / r^6$ term is the classical thermodynamic equivalent of the London quantum dispersive attraction $IE\alpha^2/r^6$. It corresponds to a small perturbation to the London forces, because $k_B T$ is at least two orders of magnitude smaller than the typical ionization energy $IE$. The smaller the temperature of the Drude motion, the closer the effective potential is to the SCF potential, making Eq. (9-57) independent of $m_D$, the mass of the oscillators.

To approximately reproduce the dynamics equivalent to the SCF regimen of Eq. (9-51) a scheme involving low-temperature Drude particles was devised. It involves the use of two Nose'–Hoover thermostats [192]: one to keep the atoms at room temperature $T$ and the second to reduce the thermal fluctuations of the Drude oscillators by imposing a temperature $T_*$ on the Drudes such that $T_* \ll T$. The idea of cooling the polarization degrees of freedom with a separate thermostat was carefully studied by Sprik [15], who showed that, for cold dipoles, both the equilibrium and diffusion properties are independent of the value of the dipole inertia parameter (the analog of $m_D$), as long as it is sufficiently small. For Drude oscillators, the temperature $T_*$ should be small enough to leave almost no kinetic energy in the atom–Drude vibrations, yet large enough to allow the Drude particles to readjust to the room-temperature motion of the atoms. This is attained with the second thermostat by coupling the motion of the Drude particles relative to their nuclei, $\mathbf{d}$ (not to their absolute motion $\dot{\mathbf{r}}_D$). Denoting $\mathbf{R}_i$ the center of mass of each $(\mathbf{r}_i, \mathbf{r}_{D,i})$ pair, $m_i$ the total mass of the pair (as before), and $m_i' = m_D(1 - m_D/m_i)$ the reduced mass, the equations of motion are

$$m_i \ddot{\mathbf{R}}_i = \mathbf{F}_{R,i} - m_i \dot{\mathbf{R}}_i \dot{\eta} \tag{9-58}$$

$$m_i' \ddot{\mathbf{d}}_i = \mathbf{F}_{d,i} - m_i' \dot{\mathbf{d}}_i \dot{\eta}_* \tag{9-59}$$

$$Q\ddot{\eta} = \sum m_j \dot{R}_j^2 - N_f k_B T \tag{9-60}$$

$$Q_* \ddot{\eta}_* = \sum m_j' \dot{d}_j^2 - N_{f*} k_B T_* \tag{9-61}$$

Indices $i$ and $j$ run from 1 to $N$, the total number of atoms. Because not all atoms have to be polarizable, the total number of Drude particles, $N_D$, may be less than $N$. If a given atom $i$ bears no Drude oscillator, $\mathbf{R}_i$ corresponds to $\mathbf{r}_i$, $m'_i$ is zero, and the corresponding Eqs. (9-59) and (9-61) are ignored. $N_f$ is the number of degrees of freedom associated with the atomic motion, accounting for distance constraints imposed by SHAKE [193], and $N_f \equiv 3N_D$ is the number of degrees of freedom associated with the motion of the Drude oscillators. Q and $Q_*$ are the inertia factors of the Nosé–Hoover thermostats. The "velocities" $\dot{\eta}$ and $\dot{\eta}_*$ are acting as friction coefficients, that is, as scaling exponents on the velocities $\dot{\mathbf{r}}$ and $\dot{\mathbf{d}}$, respectively.

### 9.4.3. Molecular Dynamics Using Fluctuating Charges

In most fluctuating charge models, if the energy is quadratic in the charges (as in Eq. 9-35), the minimization condition (Eq. 9-37) leads to a coupled set of linear equations for the charges. As with the polarizable induced dipole and Drude models, solving for the charges can be done by matrix inversion, iteration, or extended Lagrangian methods. And similarly to the other polarizable models, the matrix methods tend to be avoided because of their computational cost, and when they are used, the matrix inversion is typically not performed at every step [194, 195].

Applications of the fluctuating charge model have relied on iterative methods to determine the converged charges [52, 159, 164, 196] and for very large-scale systems, multilevel methods have also been developed [197, 198]. MC methods have also been used with fluctuating-charge models [116, 194].

As with the other polarizable models, proper treatment of long-range electrostatic interactions is essential as found by English [199]. In a comparison of the Lekner [200–202], Ewald [162, 183] and reaction field [203–205] methods to handle the long range electrostatics of several water models, including the additive flexible SPC, rigid SPC, SPC/E, TIP4P and TIP4P-Ew, and the polarizable TIP4P-FQ, English [199] found that the Lekner method gave the best results, while the reaction field method produced the worst agreement with the experimental data. Another interesting conclusion was that the Ewald method was 3.5–5 times faster than the Lekner technique. Despite this variety of available techniques, the most common approach is to use a matrix inversion or iterative method only to obtain the initial energy-minimized charge distribution and then use an extended Lagrangian method to propagate the charges dynamically in order to take advantage of its computational efficiency.

As discussed for the induced dipole and Drude polarizable models the extended Lagrangian method [168, 206, 207] is the most efficient strategy to perform MD simulations. In the extended Lagrangian applied to a fluctuating charge system [17, 92], the charges are given a fictitious mass, $m_q$, and evolved in time according to Newton's equation of motion. The extended Lagrangian corresponding to the energy defined in Eq. (9-34) is

$$L = \frac{1}{2} \sum_{i=1}^{N_{molec}} \left( \sum_{\mu=1}^{N_{atom}} M_\mu \dot{r}_{i\mu}^2 + \sum_{\alpha=1}^{N_{site}} m_q q_{i\alpha}^2 \right) - U^{system}(\mathbf{r}, \mathbf{q}) - \sum_{i=1}^{N_{molec}} \lambda_i \sum_{\alpha=1}^{N_{site}} q_{i\alpha} \quad (9\text{-}62)$$

In Eq. (9-62) $M_\mu$ is the mass of atom $\mu$. The mass $m_q$ does not correspond to any physical mass and is simply set to a value small enough such that the charges follow the atomic coordinates adiabatically. The Lagrangian also includes an $N_{molec}$ number of constraints to ensure that each molecule remains electrostatically neutral.

Based on the Lagrangian of Eq. (9-62), the equations of motion for the positions and charges are

$$M_\mu \ddot{\mathbf{r}}_{i\mu} = -\frac{\partial}{\partial \mathbf{r}_{i\mu}} \left[ U^{system}(\mathbf{r}, \mathbf{q}) \right] \quad (9\text{-}63)$$

and

$$m_q \ddot{\mathbf{q}}_{i\alpha} = -\frac{\partial U^{system}(\mathbf{r}, \mathbf{q})}{\partial \mathbf{r}_{i\mu}} - \lambda_i \quad (9\text{-}64)$$

These equations of motion can be integrated by many standard ensembles: constant energy, constant volume, constant temperature and constant pressure. More complex forms of the extended Lagrangian are possible and readers are referred to Ref. [17] for a Lagrangian that allows intermolecular charge transfer.

Although a direct comparison between the iterative and the extended Lagrangian methods has not been published, the two methods are inferred to have comparable computational speeds based on indirect evidence. The extended Lagrangian method was found to be approximately 20 times faster than the standard matrix inversion procedure [117] and according to the calculation of Bernardo et al. [208] using different polarizable water potentials, the iterative method is roughly 17 times faster than direct matrix inversion to achieve a convergence of $1.0 \times 10^{-8}$ D in the induced dipole.

## 9.5.    POLARIZABLE FORCE FIELDS IN CHARMM

### 9.5.1.    Classical Drude Oscillatory

The polarizable Drude model in CHARMM results from the work of MacKerell, Roux and co-workers [10, 11] and it is geared at developing polarizable force fields for biological macromolcules. Significant progress on the model has been made to date. The algorithm for extended Lagrangian MD has been generically described in Ref. [12] and a water model that is a generalization of the TIP4P model and has been described in Refs. [10] and [11]. The second water model, termed SWM4-NDP, uses a negative charge on the Drude particle (e.g. NDP) and represents the model that acts as the baiss for the rest of the force field. The protocols to determine fixed charges, polarizabilities and Thole factors are published in Refs. [58] and [60].

Results of parameter optimization and MD simulations of small model compounds have been published, including alcohols [63], alkanes [63], aromatic [64] and heteroaromatic [209] compounds and liquid amides [65]. Studies of ions in aqueous solution were also performed [61, 88] and results from an MD simulation on a DPPC lipid monolayer have been reported (Harder, MacKerell, Roux, submitted). Notable from the monolayer study was the reproduction of the dipole potential across the monolayer, a value that cannot be reproduced using non-polarizable models. This exciting, unforeseen observation points to the types of results that may be obtained from polarizable macromolecular force fields that are not accessible to the present additive models.

Details of the implementation of the Drude model in CHARMM follow. The potential energy function is the same of Eq. (9-25). The interactions of the various pairs of charges are treated according to the connectivity of the atoms in the molecule. As in standard additive force fields, the interactions between core charges corresponding to 1–2 (neighbor) and 1–3 (next-neighbor) covalently bound atom pairs are accounted by explicit bonding terms in the potential energy, $U_{internal}$, and excluded from the electrostatic and LJ energy calculations. Similarly, the interactions of the Drude oscillators with core charges are excluded for 1–2 and 1–3 pairs; however, the 1–2 and 1–3 dipole–dipole interactions are included in the model. The Coulomb interactions between Drude oscillators corresponding to 1–2 and 1–3 atom pairs are damped by a Thole like function $S_{ij}$ [50, 146]. An important extension of the damping function was the inclusion of atom-specific damping parameters, $\alpha$ in Eq. (9-27), allowing for improved reproduction of the molecular polarizability tensor in polar neutral species [65, 209]. The interactions involving all core charges and all Drude oscillators are included for all 1–4 pairs and beyond without screening.

An important addition to the model was the inclusion of virtual particles representative of lone pairs on hydrogen bond acceptors [60]. Their inclusion was motivated by the inability of the atom-based electrostatic model to treat interactions with water as a function of orientation. By distributing the atomic charges on to lone pairs it was possible to reproduce QM interaction energies as a function of orientation. The addition of lone pairs may be considered analogous to the use of atomic dipoles on such atoms. In the model, the polarizability is still maintained on the parent atom. In addition, anisotropic atomic polarizability, as described in Eq. (9-28), is included on hydrogen bond acceptors [65]. Its inclusion allows for reproduction of QM polarization response as a function of orientation around S, O and N atoms and it facilitates reproduction of QM interaction energies with ions as a function of orientation.

The parametrization protocol developed for the polarizable Drude model of CHARMM is well defined. A procedure for determining core and Drude charges [58, 60] and Thole damping parameters [65] has been developed and is analogous to work by Friesner and co-workers [117, 120, 138, 210, 211]. A map of the electrostatic potential (ESP) that surrounds the model compound monomer is evaluated on a set of specified grid points using density functional theory computations at the B3LYP/aug-cc-pVDZ level. To measure the electronic response of the molecule, a series of perturbed ESP maps is computed by placing a single $+0.5e$ test charge

at chemically relevant positions around the molecule. The same calculations are repeated using the Drude model, restricting the force constant tensor of each oscillator to be isotropic. Optimal parameters are chosen to minimize the difference between the QM and Drude ESP maps [58, 60]. On atoms bearing lone pairs the atomic charge is moved to the lone pair. It must be emphasized that the lone pair geometry and the force constant tensor are not part of the automated fitting procedure described above. The force constant tensor of atoms with lone pairs required for the anisotropic polarizability is determined by comparing to the local QM polarization response in the vicinity of the oxygen atom. Perturbation charges are placed around the lone pair and the QM and Drude polarization responses are computed. The components of the force constants are then manually adjusted to reproduce the QM polarization responses. Optimization of the components of the potential energy equation (Eq. 9-25) not dependent of the Drude oscillator positions, namely the bonding and Lennard-Jones terms, are adjusted as described previously for the additive CHARMM force field [2, 212, 213].

### 9.5.2. Fluctuating Charge Model

The polarizable fluctuating charge model in CHARMM results from the work of Patel, Brooks and co-workers [92, 214]. The water model is based on the TIP4P-FQ model of Rick, Stuart and Berne [17]. In the development of the force field the electronegativities and hardnesses were treated as empirical parameters and do not have any association with experimental or QM values, for example, from ionization energies and electron affinities of single atoms.

The fitting of electronegativities and hardnesses is done independently of each other with the help of a reformulation of the fluctuating charge model in terms of a linear response model [117, 120, 210]. In the presence of an external potential $\vec{\Phi}$ the electrostatic energy defined in Eq. (9-35) is:

$$
\begin{aligned}
U_{pert}^{ee}\left(\mathbf{r}, \mathbf{q}\right) = & \sum_{i=1}^{N_{molec}} \sum_{\alpha=1}^{N_{site}} \left[ \chi_{i\alpha}^{0} q_{i\alpha}^{pert} + \frac{1}{2} J_{i\alpha i\alpha}^{0} \left( q_{i\alpha}^{pert} \right)^2 \right] \\
& + \sum_{i=1}^{N_{molec}} \sum_{\alpha=1}^{N_{site}} \sum_{\beta > \alpha}^{N_{site}} J_{i\alpha i\beta} \left( \mathbf{r}_{i\alpha i\beta} \right) q_{i\alpha}^{pert} q_{i\beta}^{pert} \\
& + \sum_{i=1}^{N_{molec}} \sum_{j>1}^{N_{molec}} \sum_{\alpha=1}^{N_{site}} \sum_{\beta > \alpha}^{N_{site}} J_{i\alpha j\beta} \left( \mathbf{r}_{i\alpha j\beta} \right) q_{i\alpha}^{pert} q_{j\beta}^{pert} \\
& + \sum_{i=1}^{N_{molec}} \sum_{\alpha=1}^{N_{site}} \Phi_{i\alpha} q_{i\alpha}^{pert}
\end{aligned}
\tag{9-65}
$$

where $q^{pert}$ is a charge in the presence of the external field. At equilibrium, the conditions

$$\frac{\partial U^{ee}}{\partial q_i} = 0, \qquad i = 1, N_{site} \tag{9-66}$$

and

$$\frac{\partial U^{ee}_{pert}}{\partial q_i^{pert}} = 0, \qquad i = 1, N_{site} \tag{9-67}$$

lead to the following equations:

$$\mathbf{J} \cdot \mathbf{q} = -\chi \tag{9-68}$$

and

$$\mathbf{J} \cdot \mathbf{q}^{pert} = -(\chi + \mathbf{\Phi}) \tag{9-69}$$

Taking the difference of Eqs. (9-68) and (9-69) yields an expression for the response of the molecular charge distribution to the external field

$$\mathbf{J} \cdot \left(\mathbf{q}^{pert} - \mathbf{q}\right) = \mathbf{J} \cdot \Delta\mathbf{q} = -\mathbf{\Phi} \Rightarrow \Delta\mathbf{q} = -\left(\mathbf{J}\right)^{-1}\mathbf{\Phi} \tag{9-70}$$

An objective function measuring the deviation from the parameterized model and the target response, determined, for example, from density functional theory based methods can be defined as

$$\varepsilon = \left\| \left(\Delta q^{DFT} - \Delta q^{FQ}\right) \right\| \tag{9-71}$$

Minimization of Eq. (9-71) or other related expression of $\Delta q^{DFT}$ and $\Delta q^{FQ}$ yields optimal model parameters. Details of the choice of model compounds and placement of probes is given in reference [92]. Parameterization of the atom electronegativities is performed next by fitting to charge distributions of the optimized model compounds in the gas phase. Dipole moments are included in the fitting procedure as this was found necessary to accurately reproduce gas-phase moments, since those are very sensitive to small changes in the charge distributions [92]. MD is performed within an extended Lagrangian formalism and optimization of the LJ parameters has been based on the reproduction of small molecular pure solvent properties [92]. Application of the model have been performed on interfaces [215–217], ions in aqueous solution [218, 219] and a simulation study of DMPC have been reported [220]. Recently, the implementation of the fluctuating charge model in CHARMM was extended to allow for calculation of free energies of salvation [221], an important tool for future force field.

### 9.5.3.    Induced Dipole Model

An implementation of the induced dipole method in CHARMM has been reported
[25], based on the polarizable intermolecular potential functions (PIPF) model of
Gao and co-workers [23, 24]. The PIPF potential combined with the CHARMM22
force field as been designed PIPF-CHARMM. The implementation closely follows
the description of the induced dipole model in Section 9.2. The induced dipole at
the $i$th interaction site due to the homogeneous external electric field $\mathbf{E}_i$ is given by
Eq. (9-3) and the dipole field tensor is defined by Eq. (9-4). Infinite polarization[1] is
avoided by using Thole's damping scheme [50, 146]. The damping scheme imple-
mented is equivalent to considering a smeared charge distribution between two inter-
acting sites with a charge distribution of the exponential form, similar to Eq. (9-44b).
The three standard methods to solve Eq. (9-3) discussed in Section 9.2.1, namely the
self-consistent approach (iterative procedure), the direct solution via matrix alge-
bra and the extended Lagrangian method have been implemented. Application of
the PIPF-CHARMM method to liquid amides and alkanes have been reported [25].
When used in conjunction with the CHARMM22 force field the method required
only minor modifications of several torsional parameters, yielding adequate struc-
tural and thermodynamic properties.

### 9.6.    CONCLUSION

In this chapter we have presented an overview of the formalisms used for the in-
clusion of inducible electronic polarizability into empirical forces. All methods are
based on the polarizability relaxing in a self consistent fashion in the electric field
of the static charges and multipoles and in the field of the inducible dipoles. Intro-
duction of dipoles may be based on the explicit introduction of inducible dipoles on
atomic centers or along bonds, introduction of dipoles via the inclusion of additional
charge sites attached directly to polarizable atoms, termed the Drude model, or by
allowing charges to fluctuate in the surrounding electric field. All three approaches
have been implemented and are currently undergoing development, though at the
time of the writing of this chapter, a broad polarizable force field for biological
macromolecules was not yet available. In one limit, the three methods are similar,
yielding molecular polarizabilities as a function of environment; however, differ-
ences do exist. For example, in the fluctuating charge model out of plane polarization
does not occur in planar molecular unless off-atomic sites are included in the model
whereas the Drude model, via the explicit inclusion of different particles allows for
the inclusion of mechanical polarizabilities by including LJ parameters on the Drude
particles [8]. While such aspects could be the subject of additional discussions,
in the end it is the proper implementation of a polarizable model in combination
with an accurate optimization of the parameters in that model that will lead to the

---

[1] See Section 9.5 for a detailed discussion.

ultimate utility of polarizable force fields to the scientific community. We look forward to these future developments in the laboratories of others as well as in our own laboratories.

## ACKNOWLEDGMENTS

## REFERENCES

1. Ponder JW, Case DA (2003) Force fields for protein simulations. In: Daggett V, Eisendberg DS, Richards FM, Kuriyan J (eds) Protein Simulations, vol 66. Elsevier Academic Press, New York, pp 27–86
2. Mackerell AD (2004) Empirical force fields for biological macromolecules: Overview and issues. J Comput Chem 25(13):1584–1604
3. Jorgensen WL, Chandrasekhar J, Madura JD, Impey RW, Klein ML (1983) Comparison of simple potential functions for simulating liquid water. J Chem Phys 79(2):926
4. Cornell WD, Cieplak P, Bayly CI, Gould IR, Merz KM, Ferguson DM, Spellmeyer DC, Fox T, Caldwell JW, Kollman PA (1995) A 2nd generation force-field for the simulation of proteins, nucleic-acids, and organic-molecules. J Am Chem Soc 117(19):5179–5197
5. MacKerell AD, Bashford D, Bellott M, Dunbrack RL, Evanseck JD, Field MJ, Fischer S, Gao J, Guo H, Ha S, Joseph-McCarthy D, Kuchnir L, Kuczera K, Lau FTK, Mattos C, Michnick S, Ngo T, Nguyen DT, Prodhom B, Reiher WE, Roux B, Schlenkrich M, Smith JC, Stote R, Straub J, Watanabe M, Wiorkiewicz-Kuczera J, Yin D, Karplus M (1998) All-atom empirical potential for molecular modeling and dynamics studies of proteins. J Phys Chem B 102(18):3586–3616
6. MacKerell AD Jr, Brooks B, Brooks CL III, Nilsson L, Roux B, Won Y, Karplus M (1998) CHARMM: The energy function and its paramerization with an overview of the program. In: Schleyer PvR, Allinger NL, Clark T, Gasteiger J, Kollman PA, Schaefer HF III, Schreiner PR (eds) Encyclopedia of computational chemistry, vol 1. John Wiley & Sons, Chichester, UK, p 271
7. Halgren TA, Damm W (2001) Polarizable force fields. Curr Opin Struct Biol 11(2):236–242
8. Rick SW, Stuart SJ (2002) Potentials and algorithms for incorporating polarizability in computer simulations. In: Reviews in computational chemistry, vol 18. Wiley-Vch, Inc, New York, pp 89–146
9. Warshel A, Kato M, Pisliakov AV (2007) Polarizable force fields: history, test cases, and prospects. J Chem Theory Comput 3(6):2034–2045
10. Lamoureux G, MacKerell AD, Roux B (2003) A simple polarizable model of water based on classical Drude oscillators. J Chem Phys 119(10):5185–5197
11. Lamoureux G, Harder E, Vorobyov IV, Roux B, MacKerell AD (2006) A polarizable model of water for molecular dynamics simulations of biomolecules. Chem Phys Lett 418(1–3):245–249
12. Lamoureux G, Roux B (2003) Modeling induced polarization with classical Drude oscillators: theory and molecular dynamics simulation algorithm. J Chem Phys 119(6):3025–3039
13. Ahlstrom P, Wallqvist A, Engstrom S, Jonsson B (1989) A molecular-dynamics study of polarizable water. Mol Phys 68(3):563–581
14. Caldwell J, Dang LX, Kollman PA (1990) Implementation of nonadditive intermolecular potentials by use of molecular-dynamics – development of a water water potential and water ion cluster interactions. J Am Chem Soc 112(25):9144–9147

15. Sprik M (1991) Computer-simulation of the dynamics of induced polarization fluctuations in water. J Phys Chem 95(6):2283–2291
16. Millot C, Stone AJ (1992) Towards an accurate intermolecular potential for water. Molcular Phys 77(3):439–462
17. Rick SW, Stuart SJ, Berne BJ (1994) Dynamical fluctuating charge force-fields – application to liquid water. J Chem Phys 101(7):6141–6156
18. Giese TJ, York DM (2004) Many-body force field models based solely on pairwise Coulomb screening do not simultaneously reproduce correct gas-phase and condensed-phase polarizability limits. J Chem Phys 120(21):9903–9906
19. Masia M, Probst M, Rey R (2004) On the performance of molecular polarization methods. I. Water and carbon tetrachloride close to a point charge. J Chem Phys 121(15):7362–7378
20. Masia M, Probst M, Rey R (2005) On the performance of molecular polarization methods. II. Water and carbon tetrachloride close to a cation. J Chem Phys 123(16)
21. Friesner RA (2006) Modeling polarization in proteins and protein-ligand complexes: Methods and preliminary results. Adv Protein Chem 72: 79–104
22. Swart M, van Duijnen PT (2006) DRF90: a polarizable force field. Mol Simul 32(6):471–484
23. Gao JL, Habibollazadeh D, Shao L (1995) A polarizable intermolecular potential function for simulation of liquid alcohols. J Phys Chem 99(44):16460–16467
24. Gao JL, Pavelites JJ, Habibollazadeh D (1996) Simulation of liquid amides using a polarizable intermolecular potential function. J Phys Chem 100(7):2689–2697
25. Xie WS, Pu JZ, MacKerell AD, Gao JL (2007) Development of a polarizable intermolecular potential function (PIPF) for liquid amides and alkanes. J Chem Theory Comput 3(6):1878–1889
26. Ma BY, Lii JH, Allinger NL (2000) Molecular polarizabilities and induced dipole moments in molecular mechanics. J Comput Chem 21(10):813–825
27. Berendsen HJC, Grigera JR, Straatsma TP (1987) The missing term in effective pair potentials. J Phys Chem 91(24):6269–6271
28. Reynolds CA, Ferenczy GG, Richards WG (1992) Methods for determining the reliability of semiempirical electrostatic potentials and potential derived charges. Theochem J Mol Struct 88, 249–269
29. Ferenczy GG, Winn PJ, Reynolds CA (1997) Toward improved force fields: 2. Effective distributed multipoles. J Phys Chem A 101(30):5446–5455
30. Ferenczy GG, Winn PJ, Reynolds CA, Richter G (1997) Effective distributed multipoles for the quantitative description of electrostatics and polarisation in intermolecular interactions. Abs Papers Am Chem Soc 214:38-COMP
31. Winn PJ, Ferenczy GG, Reynolds CA (1997) Toward improved force fields: 1. Multipole-derived atomic charges. J Phys Chem A 101(30):5437–5445
32. Winn PJ, Ferenczy GG, Reynolds CA (1999) Towards improved force fields: III. Polarization through modified atomic charges. J Comput Chem 20(7):704–712
33. Wu JH, Winn PJ, Ferenczy GG, Reynolds CA (1999) Solute polarization and the design of cobalt complexes as redox-active therapeutic agents. Int J Quant Chem 73(2):229–236
34. Gooding SR, Winn PJ, Maurer RI, Ferenczy GG, Miller JR, Harris JE, Griffiths DV, Reynolds CA (2000) Fully polarizable QM/MM calculations: an application to the nonbonded iodine–oxygen interaction in dimethyl-2-iodobenzoylphosphonate. J Comput Chem 21(6):478–482
35. Ferenczy GG, Reynolds CA (2001) Modeling polarization through induced atomic charges. J Phys Chem A 105(51):11470–11479
36. Illingworth CJR, Gooding SR, Winn PJ, Jones GA, Ferenczy GG, Reynolds CA (2006) Classical polarization in hybrid QM/MM methods. J Phys Chem A 110(20):6487–6497

37. Ferenczy GG (1991) Charges derived from distributed multipole series. J Comput Chem 12(8): 913–917
38. Chipot C, Angyan JG, Ferenczy GG, Scheraga HA (1993) Transferable net atomic charges from a distributed multipole analysis for the description of electrostatic properties – a case-study of saturated-hydrocarbons. J Phys Chem 97(25):6628–6636
39. Zhu SB, Yao S, Zhu JB, Singh S, Robinson GW (1991) A flexible polarizable simple point-charge water model. J Phys Chem 95(16):6211–6217
40. Ren PY, Ponder JW (2002) Consistent treatment of inter- and intramolecular polarization in molecular mechanics calculations. J Comput Chem 23(16):1497–1506
41. Grossfield A, Ren PY, Ponder JW (2003) Ion solvation thermodynamics from simulation with a polarizable force field. J Am Chem Soc 125(50):15671–15682
42. Grossfield A, Ren PY, Ponder JW (2003) Single ion solvation thermodynamics from simulations. Biophys J 84(2):94A-94A
43. Ren PY, Ponder JW (2003) Polarizable atomic multipole water model for molecular mechanics simulation. J Phys Chem B 107(24):5933–5947
44. Ren PY, Ponder JW (2004) Temperature and pressure dependence of the AMOEBA water model. J Phys Chem B 108(35):13427–13437
45. Grossfield A (2005) Dependence of ion hydration on the sign of the ion's charge. J Chem Phys 122(2)
46. Jiao D, King C, Grossfield A, Darden TA, Ren PY (2006) Simulation of Ca2+ and Mg2+ solvation using polarizable atomic multipole potential. J Phys Chem B 110(37):18553–18559
47. Rasmussen TD, Ren PY, Ponder JW, Jensen F (2007) Force field modeling of conformational energies: importance of multipole moments and intramolecular polarization. Int J Quant Chem 107(6):1390–1395
48. Piquemal JP, Perera L, Cisneros GA, Ren PY, Pedersen LG, Darden TA (2006) Towards accurate solvation dynamics of divalent cations in water using the polarizable amoeba force field: from energetics to structure. J Chem Phys 125(5):054511
49. Applequist J, Carl JR, Fung K-K (1972) Atom dipole interaction model for molecular polarizability. Application to polyatomic molecules and determination of atom polarizabilities. J Am Chem Soc 94(9):2952–2960
50. Thole BT (1981) Molecular polarizabilities calculated with a modified dipole interaction. Chem Phys 59(3):341–350
51. Stone AJ (1997) The theory of intermolecular forces. Oxford University Press, Oxford
52. Vesely FJ (1977) N-particle dynamics of polarizable Stockmayer-type molecules. J Comput Phys 24(4):361–371
53. Bernardo DN, Ding YB, Kroghjespersen K, Levy RM (1994) An anisotropic polarizable water model – incorporation of all-atom polarizabilities into molecular mechanics force-fields. J Phys Chem 98(15):4180–4187
54. Pollock EL, Alder BJ, Patey GN (1981) Static dielectric properties of polarizable Stockmayer fluids. Physica A Stat Theor Phys 108(1):14–26
55. van Belle D, Couplet I, Prevost M, Wodak SJ (1987) Calculations of electrostatic properties in proteins – analysis of contributions from induced protein dipoles. J Mol Biol 198(4): 721–735
56. Yu HB, Hansson T, van Gunsteren WF (2003) Development of a simple, self-consistent polarizable model for liquid water. J Chem Phys 118(1):221–234
57. Yu HB, van Gunsteren WF (2004) Charge-on-spring polarizable water models revisited: from water clusters to liquid water to ice. J Chem Phys 121(19):9549–9564

58. Anisimov VM, Lamoureux G, Vorobyov IV, Huang N, Roux B, MacKerell AD (2005) Determination of electrostatic parameters for a polarizable force field based on the classical Drude oscillator. J Chem Theory Comput 1(1):153–168

59. Yu HB, van Gunsteren WF (2005) Accounting for polarization in molecular simulation. Comput Phys Commun 172(2):69–85

60. Harder E, Anisimov VM, Vorobyov IV, Lopes PEM, Noskov SY, MacKerell AD, Roux B (2006) Atomic level anisotropy in the electrostatic modeling of lone pairs for a polarizable force field based on the classical Drude oscillator. J Chem Theory Comput 2(6):1587–1597

61. Lamoureux G, Roux B (2006) Absolute hydration free energy scale for alkali and halide ions established from simulations with a polarizable force field. J Phys Chem B 110(7):3308–3322

62. Yu HB, Geerke DP, Liu HY, van Gunsteren WE (2006) Molecular dynamics simulations of liquid methanol and methanol–water mixtures with polarizable models. J Comput Chem 27(13): 1494–1504

63. Anisimov VM, Vorobyov IV, Roux B, MacKerell AD (2007) Polarizable empirical force field for the primary and secondary alcohol series based on the classical drude model. J Chem Theory Comput 3(6):1927–1946

64. Lopes PEM, Lamoureux G, Roux B, MacKerell AD (2007) Polarizable empirical force field for aromatic compounds based on the classical Drude oscillator. J Phys Chem B 111(11):2873–2885

65. Harder E, Anisimov VM, Whitfield TW, MacKerell AD, Roux B (2008) Understanding the dielectric properties of liquid amides from a polarizable force field. J Phys Chem B 112(11):3509–3521

66. Drude P (2008) The theory of optics (1902). Kessinger Publishing Company

67. London F (1937) The general theory of molecular forces. Trans Faraday Soc 33:8–26

68. Bade WL (1957) Drude-model calculation of dispersion forces. I. General theory. J Chem Phys 27(6):1280–1284

69. Bade WL, Kirkwood JG (1957) Drude-model calculation of dispersion forces. II. The linear lattice. J Chem Phys 27(6):1284–1288

70. Bade WL (1958) Drude-model calculation of dispersion forces. III. The fourth-order contribution. J Chem Phys 28(2):282–284

71. Amos AT (1996) Bond properties using a modern version of the Drude model. Int J Quant Chem 60(1):67–74

72. Wang F, Jordan KD (2002) Application of a Drude model to the binding of excess electrons to water clusters. J Chem Phys 116(16):6973–6981

73. Dick BG, Overhauser AW (1958) Theory of the dielectric constants of alkali halide crystals. Phys Rev 112(1):90–103

74. Hanlon JE, Lawson AW (1959) Effective ionic charge in alkali halides. Phys Rev 113(2):472–478

75. Jacucci G, McDonald IR, Singer K (1974) Introduction of the shell model of ionic polarizability into molecular dynamics calculations. Phys Lett A 50(2):141–143

76. Lindan PJD, Gillan MJ (1993) Shell-model molecular-dynamics simulation of superionic conduction in CAF2. J Phys Conden Matter 5(8):1019–1030

77. Mitchell PJ, Fincham D (1993) Shell-model simulations by adiabatic dynamics. J Phys Conden Matter 5(8):1031–1038

78. Lindan PJD (1995) Dynamics with the shell-model. Mol Simul 14(4–5):303–312

79. Noskov SY, Lamoureux G, Roux B (2005) Molecular dynamics study of hydration in ethanol-water mixtures using a polarizable force field. J Phys Chem B 109(14):6705–6713

80. Hoye JS, Stell G (1980) Dielectric theory for polar molecules with fluctuating polarizability. J Chem Phys 73(1):461–468

81. Pratt LR (1980) Effective field of a dipole in non-polar polarizable fluids. Mol Phys 40(2):347–360

82. Lado F (1997) Molecular theory of a charged particle in a polarizable nonpolar liquid. J Chem Phys 106(11):4707–4713

83. Cao J, Berne BJ (1993) Theory of polarizable liquid crystals: optical birefringence. J Chem Phys 99(3):2213–2220

84. Saint-Martin H, Medina-Llanos C, Ortega-Blake I (1990) Nonadditivity in an analytical intermolecular potential – the water–water interaction. J Chem Phys 93(9):6448–6452

85. de Leeuw NH, Parker SC (1998) Molecular-dynamics simulation of MgO surfaces in liquid water using a shell-model potential for water. Phys Rev B 58(20):13901–13908

86. Saint-Martin H, Hernandez-Cobos J, Bernal-Uruchurtu MI, Ortega-Blake I, Berendsen HJC (2000) A mobile charge densities in harmonic oscillators (MCDHO) molecular model for numerical simulations: the water–water interaction. J Chem Phys 113(24):10899–10912

87. van Maaren PJ, van der Spoel D (2001) Molecular dynamics simulations of water with novel shell-model potentials. J Phys Chem B 105(13):2618–2626

88. Whitfield TW, Varma S, Harder E, Lamoureux G, Rempe SB, Roux B (2007) Theoretical study of aqueous solvation of $K^+$ comparing ab initio, polarizable, and fixed-charge models. J Chem Theory Comput 3(6):2068–2082

89. Lu ZY, Zhang YK (2008) Interfacing ab initio quantum mechanical method with classical Drude oscillator polarizable model for molecular dynamics simulation of chemical reactions. J Chem Theory Comput 4(8):1237–1248

90. van Belle D, Froeyen M, Lippens G, Wodak SJ (1992) Molecular-dynamics simulation of polarizable water by an extended lagrangian method. Mol Phys 77(2):239–255

91. Sprik M, Klein ML (1988) A polarizable model for water using distributed charge sites. J Chem Phys 89(12):7556–7560

92. Patel S, Brooks CL (2004) CHARMM fluctuating charge force field for proteins: I parameterization and application to bulk organic liquid simulations. J Comput Chem 25(1):1–15

93. Nalewajski RF (1991) Normal (decoupled) representation of electronegativity equalization equations in a molecule. Int J Quant Chem 40(2):265–285

94. Chelli R, Procacci P (2002) A transferable polarizable electrostatic force field for molecular mechanics based on the chemical potential equalization principle. J Chem Phys 117(20):9175–9189

95. Itskowitz P, Berkowitz ML (1997) Chemical potential equalization principle: direct approach from density functional theory. J Phys Chem A 101(31):5687–5691

96. Nalewajski RF (1998) On the chemical potential/electronegativity equalization in density functional theory. Polish J Chem 72(7):1763–1778

97. Chelli R, Ciabatti S, Cardini G, Righini R, Procacci P (1999) Calculation of optical spectra in liquid methanol using molecular dynamics and the chemical potential equalization method. J Chem Phys 111(9):4218–4229

98. Bret C, Field MJ, Hemmingsen L (2000) A chemical potential equalization model for treating polarization in molecular mechanical force fields. Mol Phys 98(11):751–763

99. Chelli R, Ciabatti S, Cardini G, Righini R, Procacci P (2000) Calculation of optical spectra in liquid methanol using molecular dynamics and the chemical potential equalization method (vol 111, p 4218, 1999). J Chem Phys 112(12):5515–5515

100. Nalewajski RF (2000) Charge sensitivities of the externally interacting open reactants. Int J Quant Chem 78(3):168–178

101. Llanta E, Ando K, Rey R (2001) Fluctuating charge study of polarization effects in chlorinated organic liquids. J Phys Chem B 105(32):7783–7791

102. York DM (2002) Chemical potential equalization: a many-body force field for molecular simulations. Abs Papers Am Chem Soc 224:U472–U472

103. Smith PE (2004) Local chemical potential equalization model for cosolvent effects on biomolecular equilibria. J Phys Chem B 108(41):16271–16278
104. Chelli R, Barducci A, Bellucci L, Schettino V, Procacci P (2005) Behavior of polarizable models in presence of strong electric fields. I. Origin of nonlinear effects in water point-charge systems. J Chem Phys 123(19):194109
105. Medeiros M (2005) Monte Carlo simulation of polarizable systems: early rejection scheme for improving the performance of adiabatic nuclear and electronic sampling Monte Carlo simulations. Theor Chem Acc 113(3):178–182
106. Piquemal JP, Chelli R, Procacci P, Gresh N (2007) Key role of the polarization anisotropy of water in modeling classical polarizable force fields. J Phys Chem A 111(33):8170–8176
107. Warren GL, Davis JE, Patel S (2008) Origin and control of superlinear polarizability scaling in chemical potential equalization methods. J Chem Phys 128(14):144110
108. Zhang Y, Lin H (2008) Flexible-boundary quantum-mechanical/molecular-mechanical calculations: partial charge transfer between the quantum-mechanical and molecular-mechanical subsystems. J Chem Theory Comput 4(3):414–425
109. Rappe AK, Goddard WA (1991) Charge equilibration for molecular-dynamics simulations. J Phys Chem 95(8):3358–3363
110. Kitao O, Ogawa T (2003) Consistent charge equilibration (CQEq). Mol Phys 101(1–2):3–17
111. Nistor RA, Polihronov JG, Muser MH, Mosey NJ (2006) A generalization of the charge equilibration method for nonmetallic materials. J Chem Phys 125(9):094108
112. Ogawa T, Kurita N, Sekino H, Kitao O, Tanaka S (2004) Consistent charge equilibration (CQEq) method: application to amino acids and crambin protein. Chem Phys Lett 397(4–6):382–387
113. Sefcik J, Demiralp E, Cagin T, Goddard WA (2002) Dynamic charge equilibration-morse stretch force field: application to energetics of pure silica zeolites. J Comput Chem 23(16):1507–1514
114. Tanaka M, Siehl HU (2008) An application of the consistent charge equilibration (CQEq) method to guanidinium ionic liquid systems. Chem Phys Lett 457(1–3):263–266
115. Brodersen S, Wilke S, Leusen FJJ, Engel GE (2005) Comparison of static and fluctuating charge models for force-field methods applied to organic crystals. Cryst Growth Des 5(3):925–933
116. Chen B, Xing JH, Siepmann JI (2000) Development of polarizable water force fields for phase equilibrium calculations. J Phys Chem B 104(10):2391–2401
117. Liu YP, Kim K, Berne BJ, Friesner RA, Rick SW (1998) Constructing ab initio force fields for molecular dynamics simulations. J Chem Phys 108(12):4739–4755
118. Stuart SJ, Berne BJ (1996) Effects of polarizability on the hydration of the chloride ion. J Phys Chem 100(29):11934–11943
119. Stuart SJ, Berne BJ (1999) Surface Curvature Effects in the Aqueous Ionic Solvation of the Chloride Ion. J Phys Chem A 103(49):10300–10307
120. Banks JL, Kaminski GA, Zhou RH, Mainz DT, Berne BJ, Friesner RA (1999) Parametrizing a polarizable force field from ab initio data. I. The fluctuating point charge model. J Chem Phys 110(2):741–754
121. Rick SW, Berne BJ (1996) Dynamical fluctuating charge force fields: The aqueous solvation of amides. J Am Chem Soc 118(3):672–679
122. Toufar H, Baekelandt BG, Janssens GOA, Mortier WJ, Schoonheydt RA (1995) Investigation of supramolecular systems by a combination of the electronegativity equalization method and a Monte-Carlo simulation technique. J Phys Chem 99(38):13876–13885
123. Perng BC, Newton MD, Raineri FO, Friedman HL (1996) Energetics of charge transfer reactions in solvents of dipolar and higher order multipolar character. 1. Theory. J Chem Phys 104(18):7153–7176

124. Perng BC, Newton MD, Raineri FO, Friedman HL (1996) Energetics of charge transfer reactions in solvents of dipolar and higher order multipolar character. 2. Results. J Chem Phys 104(18): 7177–7204

125. Field MJ (1997) Hybrid quantum mechanical molecular mechanical fluctuating charge models for condensed phase simulations. Mol Phys 91(5):835–845

126. Ribeiro MCC, Almeida LCJ (1999) Fluctuating charge model for polyatomic ionic systems: a test case with diatomic anions. J Chem Phys 110(23):11445–11448

127. Iczkowski RP, Margrave JL (1961) Electronegativity. J Am Chem Soc 83(17):3547–3551

128. Mulliken RS (1934) A new electroaffinity scale, together with data on valence states and on valence ionization potentials and electron affinities. J Chem Phys 2(11):782–793

129. Parr RG, Pearson RG (1983) Absolute hardness: companion parameter to absolute electronegativity. J Am Chem Soc 105(26):7512–7516

130. Parr RG, Donnelly RA, Levy M, Palke WE (1978) Electronegativity: the density functional viewpoint. J Chem Phys 68(8):3801–3807

131. Hohenberg P, Kohn W (1964) Inhomogeneous electron gas. Phys Rev 136(3B):B864–B871

132. Kitaura K, Morokuma K (1976) A new energy decomposition scheme for molecular interactions within the Hartree-Fock approximation. Int J Quant Chem 10(2):325–340

133. Weinhold F (1997) Nature of H-bonding in clusters, liquids, and enzymes: an ab initio, natural bond orbital perspective. J Mol Struct Theochem 398–399:181–197

134. van der Vaart A, Merz KM (1999) The role of polarization and charge transfer in the solvation of biomolecules. J Am Chem Soc 121(39):9182–9190

135. Korchowiec J, Uchimaru T (2000) New energy partitioning scheme based on the self-consistent charge and configuration method for subsystems: application to water dimer system. J Chem Phys 112(4):1623–1633

136. Jeziorski B, Moszynski R, Szalewicz K (1994) Perturbation-theory approach to intermolecular potential-energy surfaces of Van-Der-Waals complexes. Chem Rev 94(7):1887–1930

137. Chelli R, Procacci P, Righini R, Califano S (1999) Electrical response in chemical potential equalization schemes. J Chem Phys 111(18):8569–8575

138. Stern HA, Kaminski GA, Banks JL, Zhou RH, Berne BJ, Friesner RA (1999) Fluctuating charge, polarizable dipole, and combined models: parameterization from ab initio quantum chemistry. J Phys Chem B 103(22):4730–4737

139. Yang ZZ, Wang CS (1997) Atom-bond electronegativity equalization method. 1. Calculation of the charge distribution in large molecules. J Phys Chem A 101(35):6315–6321

140. Wang CS, Li SM, Yang ZZ (1998) Calculation of molecular energies by atom-bond electronegativity equalization method. Theochem J Mol Struct 430, 191–199

141. Wang CS, Yang ZZ (1999) Atom-bond electronegativity equalization method. II. Lone-pair electron model. J Chem Phys 110(13):6189–6197

142. Cong Y, Yang ZZ (2000) General atom-bond electronegativity equalization method and its application in prediction of charge distributions in polypeptide. Chem Phys Lett 316(3–4):324–329

143. Yang ZZ, Wang CS (2003) Atom-bond electronegativity equalization method and its applications based on density functional theory. J Theor Comput Chem 2(2):273–299

144. Yang ZZ, Wu Y, Zhao DX (2004) Atom-bond electronegativity equalization method fused into molecular mechanics. I. A seven-site fluctuating charge and flexible body water potential function for water clusters. J Chem Phys 120(6):2541–2557

145. Wu Y, Yang ZZ (2004) Atom-bond electronegativity equalization method fused into molecular mechanics. II. A seven-site fluctuating charge and flexible body water potential function for liquid water. J Phys Chem A 108(37):7563–7576

146. van Duijnen PT, Swart M (1998) Molecular and atomic polarizabilities: Thole's model revisited. J Phys Chem A 102(14):2399–2407

147. Stillinger FH (1979) Dynamics and ensemble averages for the polarization models of molecular interactions. J Chem Phys 71(4):1647

148. Wallqvist A, Berne BJ (1993) Effective potentials for liquid water using polarizable and nonpolarizable models. J Phys Chem 97(51):13841–13851

149. Stillinger FH, David CW (1978) Polarization model for water and its ionic dissociation products. J Chem Phys 69(4):1473

150. Kuwajima S, Warshel A (1990) Incorporating electric polarizabilities in water water interaction potentials. J Phys Chem 94(1):460–466

151. Ojamae L, Shavitt I, Singer SJ (1998) Potential models for simulations of the solvated proton in water. J Chem Phys 109(13):5547–5564

152. Ding YB, Bernardo DN, Kroghjespersen K, Levy RM (1995) Solvation free-energies of small amides and amines from molecular-dynamics free-energy perturbation simulations using pairwise additive and many-body polarizable potentials. J Phys Chem 99(29):11575–11583

153. Tuckerman M, Berne BJ, Martyna GJ (1992) Reversible multiple time scale molecular-dynamics. J Chem Phys 97(3):1990–2001

154. Tuckerman M, Berne BJ, Martyna GJ (1993) Reversible multiple time-scale molecular-dynamics – reply. J Chem Phys 99(3):2278–2279

155. Martyna GJ, Tuckerman ME, Tobias DJ, Klein ML (1996) Explicit reversible integrators for extended systems dynamics. Mol Phys 87(5):1117–1157

156. Pollock EL, Alder BJ (1977) Effective field of a dipole in polarizable fluids. Phys Rev Lett 39(5):299–302

157. Chialvo AA, Cummings PT (1996) Engineering a simple polarizable model for the molecular simulation of water applicable over wide ranges of state conditions. J Chem Phys 105(18):8274–8281

158. Dang LX, Chang TM (1997) Molecular dynamics study of water clusters, liquid, and liquid–vapor interface of water with many-body potentials. J Chem Phys 106(19):8149–8159

159. Barnes P, Finney JL, Nicholas JD, Quinn JE (1979) Cooperative effects in simulated water. Nature 282(5738):459–464

160. Chang TM, Peterson KA, Dang LX (1995) Molecular-dynamics simulations of liquid, interface, and ionic solvation of polarizable carbon-tetrachloride. J Chem Phys 103(17):7502–7513

161. Toukmaji A, Sagui C, Board J, Darden T (2000) Efficient particle-mesh Ewald based approach to fixed and induced dipolar interactions. J Chem Phys 113(24):10913–10927

162. Allen MP, Tildesley DJ (1989) Computer simulation of liquids. Clarendon Press, New York

163. Brodholt J, Sampoli M, Vallauri R (1995) Parameterizing polarizable intermolecular potentials for water with the ice 1 h phase. Mol Phys 85(1):81–90

164. Svishchev IM, Kusalik PG, Wang J, Boyd RJ (1996) Polarizable point-charge model for water: results under normal and extreme conditions. J Chem Phys 105(11):4742–4750

165. Wallqvist A, Ahlstrom P, Karlstrom G (1990) A new intermolecular energy calculation scheme – applications to potential surface and liquid properties of water. J Phys Chem 94(4):1649–1656

166. Ruocco G, Sampoli M (1994) Computer-simulation of polarizable fluids – a consistent and fast way for dealing with polarizability and hyperpolarizability. Mol Phys 82(5):875–886

167. Kaminski GA, Friesner RA, Zhou RH (2003) A computationally inexpensive modification of the point dipole electrostatic polarization model for molecular simulations. J Comput Chem 24(3):267–276

168. Car R, Parrinello M(1985) Unified approach for molecular dynamics and density-functional theory. Phys Rev Lett 55(22):2471–2474

169. van Belle D, Wodak SJ (1995) Extended Lagrangian formalism applied to temperature control and electronic polarization effects in molecular dynamics simulations. Comput Phys Commun 91 (1–3):253–262

170. Halley JW, Rustad JR, Rahman A (1993) A polarizable, dissociating molecular-dynamics model for liquid water. J Chem Phys 98(5):4110–4119

171. Saboungi ML, Rahman A, Halley JW, Blander M (1988) Molecular-dynamics studies of complexing in binary molten-salts with polarizable anions – Max4. J Chem Phys 88(9):5818–5823

172. Harder E, Kim BC, Friesner RA, Berne BJ (2005) Efficient simulation method for polarizable protein force fields: application to the simulation of BPTI in liquid. J Chem Theory Comput 1(1): 169–180

173. Kiyohara K, Gubbins KE, Panagiotopoulos AZ (1998) Phase coexistence properties of polarizable water models. Mol Phys 94(5):803–808

174. Jedlovszky P, Vallauri R (1999) Temperature dependence of thermodynamic properties of a polarizable potential model of water. Mol Phys 97(11):1157–1163

175. Mahoney MW, Jorgensen WL (2001) Rapid estimation of electronic degrees of freedom in Monte Carlo calculations for polarizable models of liquid water. J Chem Phys 114(21):9337–9349

176. Goodfellow JM (1982) Cooperative effects in water-biomolecule crystal systems. Proc Natl Acad Sci USA 79(16):4977–4979

177. Cabral BJC, Rivail JL, Bigot B (1987) A Monte-Carlo simulation study of a polarizable liquid – influence of the electrostatic induction on its thermodynamic and structural-properties. J Chem Phys 86(3):1467–1473

178. Rullmann JAC, van Duijnen PT (1988) A polarizable water model for calculation of hydration energies. Mol Phys 63(3):451–475

179. Cieplak P, Kollman P, Lybrand T (1990) A new water potential including polarization – application to gas-phase, liquid, and crystal properties of water. J Chem Phys 92(11):6755–6760

180. Jedlovszky P, Vallauri R (2005) Liquid–vapor and liquid–liquid phase equilibria of the Brodholt-Sampoli-Vallauri polarizable water model. J Chem Phys 122(8):081101

181. Valdez-Gonzales M, Sanit-Martin H, Hernandez-Cobos J, Ayala R, Sanchez-Marcos E, Ortega-Blake I (2007) Liquid methanol Monte Carlo simulations with a refined potential which includes polarizability, nonadditivity, and intramolecular relaxation. J Chem Phys 127(22):224507

182. Nymand TM, Linse P (2000) Molecular dynamics simulations of polarizable water at different boundary conditions. J Chem Phys 112(14):6386–6395

183. Ewald PP (1921) Die Berechnung optischer und elektrostatischer Gitterpotentiale. Annalen der Physik 369(3):253–287

184. Darden T, Perera L, Li LP, Pedersen L (1999) New tricks for modelers from the crystallography toolkit: the particle mesh Ewald algorithm and its use in nucleic acid simulations. Struct Fold Des 7(3):R55–R60

185. Kutteh R, Nicholas JB (1995) Efficient dipole iteration in polarizable charged systems using the cell multipole method and application to polarizable water. Comput Phys Commun 86(3):227–235

186. Kutteh R, Nicholas JB (1995) Implementing the cell multipole method for dipolar and charged dipolar systems. Comput Phys Commun 86(3):236–254

187. Toukmaji AY, Board JA (1996) Ewald summation techniques in perspective: a survey. Comput Phys Commun 95(2–3):73–92

188. Sandak B (2001) Multiscale fast summation of long-range charge and dipolar interactions. J Comput Chem 22(7):717–731

189. Jacucci G, McDonald IR, Rahman A (1976) Effects of polarization on equilibrium and dynamic properties of ionic systems. Phys Rev A 13(4):1581–1592

190. Sangster MJL, Dixon M (1976) Interionic potentials in alkali halides and their use in simulations of the molten salts. Adv Phys 25(3):247–342

191. Dixon M, Sangster MJL (1975) Simulation of molten NaI including polarization effects. J Phys C Solid State Phys 8(1):L8–L11

192. Hoover WG (1985) Canonical dynamics: equilibrium phase-space distributions. Phys Rev A 31(3):1695–1697

193. Ryckaert JP, Ciccotti G, Berendsen HJC (1977) Numerical integration of the Cartesian equations of motion of a system with constraints: molecular dynamics of *n*-alkanes. J Comput Phys 23(3): 327–341

194. Medeiros M, Costas ME (1997) Gibbs ensemble Monte Carlo simulation of the properties of water with a fluctuating charges model. J Chem Phys 107(6):2012–2019

195. Campbell T, Kalia RK, Nakano A, Vashishta P, Ogata S, Rodgers S (1999) Dynamics of oxidation of aluminum nanoclusters using variable charge molecular-dynamics simulations on parallel computers. Phys Rev Lett 82(24):4866–4869

196. Streitz FH, Mintmire JW (1994) Electrostatic potentials for metal-oxide surfaces and interfaces. Phys Rev B 50(16):11996–12003

197. Keffer DJ, Mintmire JW (2000) Efficient parallel algorithms for molecular dynamics simulations using variable charge transfer electrostatic potentials. Int J Quant Chem 80(4–5):733–742

198. Nakano A (1997) Parallel multilevel preconditioned conjugate-gradient approach to variable-charge molecular dynamics. Comput Phys Commun 104(1–3):59–69

199. English NJ (2005) Molecular dynamics simulations of liquid water using various long range electrostatics techniques. Mol Phys 103(14):1945–1960

200. Gronbech-Jensen N (1997) Lekner summation of long range interactions in periodic systems. Int J Mod Phys C 8(6):1287–1297

201. Lekner J (1989) Summation of dipolar fields in simulated liquid vapor interfaces. Physica A 157(2):826–838

202. Lekner J (1991) Summation of coulomb fields in computer-simulated disordered-systems. Physica A 176(3):485–498

203. Barker JA, Watts RO (1973) Monte Carlo studies of the dielectric properties of water-like models. Mol Phys 26(3):789–792

204. Neumann M (1983) Dipole moment fluctuation formulas in computer simulations of polar systems. Mol Phys 50(4):841–858

205. Neumann M (1985) The dielectric constant of water. Computer simulations with the MCY potential. J Chem Phys 82(12):5663–5672

206. Andersen HC (1980) Molecular dynamics simulations at constant pressure and/or temperature. J Chem Phys 72(4):2384–2393

207. Parrinello M, Rahman A (1980) Crystal structure and pair potentials: a molecular-dynamics study. Phys Rev Lett 45(14):1196–1199

208. Bernardo DN, Ding YB, Kroghjespersen K, Levy RM (1995) Evaluating polarizable potentials on distributed-memory parallel computers – program-development and applications. J Comput Chem 16(9):1141–1152

209. Lopes PEM, Lamoureux G, MacKerell AD, Polarizable Empirical Force Field for Nitrogen-containing Heteroaromatic Compounds Based on the Classical Drude Oscillator. Accepted for publication on J Comput Chem

210. Stern HA, Rittner F, Berne BJ, Friesner RA (2001) Combined fluctuating charge and polarizable dipole models: application to a five-site water potential function. J Chem Phys 115(5):2237–2251

211. Kaminski GA, Stern HA, Berne BJ, Friesner RA, Cao YXX, Murphy RB, Zhou RH, Halgren TA (2002) Development of a polarizable force field for proteins via ab initio quantum chemistry: first generation model and gas phase tests. J Comput Chem 23(16):1515–1531

212. Foloppe N, MacKerell AD (2000) All-atom empirical force field for nucleic acids: I. Parameter optimization based on small molecule and condensed phase macromolecular target data. J Comput Chem 21(2):86–104

213. Yin DX, Mackerell AD (1998) Combined ab initio empirical approach for optimization of Lennard-Jones parameters. J Comput Chem 19(3):334–348

214. Patel S, Mackerell AD, Brooks CL (2004) CHARMM fluctuating charge force field for proteins: II – Protein/solvent properties from molecular dynamics simulations using a nonadditive electrostatic model. J Comput Chem 25(12):1504–1514

215. Patel S, Brooks CL (2005) Structure, thermodynamics, and liquid–vapor equilibrium of ethanolfrom molecular-dynamics simulations using nonadditive interactions. J Chem Phys 123(16):164502

216. Patel S, Brooks CL (2005) A nonadditive methanol force field: bulk liquid and liquid–vapor interfacial properties via molecular dynamics simulations using a fluctuating charge model. J Chem Phys 122(2)

217. Zhong Y, Warren GL, Patel S (2008) Thermodynamic and structural properties of methanol-water solutions using nonadditive interaction models. J Comput Chem 29(7):1142–1152

218. Warren GL, Patel S (2007) Hydration free energies of monovalent ions in transferable intermolecular potential four point fluctuating charge water: an assessment of simulation methodology and force field performance and transferability. J Chem Phys 127(6):064509

219. Warren GL, Patel S (2008) Comparison of the solvation structure of polarizable and nonpolarizable ions in bulk water and near the aqueous liquid–vapor interface. J Phys Chem C 112(19):7455–7467

220. Patel S, Brooks CL (2006) Fluctuating charge force fields: recent developments and applications from small molecules to macromolecular biological systems. Mol Simul 32(3–4), 231–249

221. Patel S, Brooks CL (2006) Revisiting the hexane-water interface via molecular dynamics simulations using nonadditive alkane-water potentials. J Chem Phys 124(20):204706

# Part III
# Biocatalysis Applications

CHAPTER 10

# MODELING PROTONATION EQUILIBRIA IN BIOLOGICAL MACROMOLECULES

JANA KHANDOGIN

*Department of Chemistry & Biochemistry, University of Oklahoma, Norman, OK 73019, USA,*
*e-mail: Jana.K.Shen@ou.edu*

**Abstract:**     The stability and function of proteins are dependent on the charge states. For more than a decade, theoretical methods for the prediction of protonation equilibria in proteins have been based on a macroscopic description in which the dielectric response of protein to the fluctuating environment is modeled implicitly through an effective dielectric constant. Recently, constant pH molecular dynamics methods have been developed, which allow for an explicit coupling between the conformational dynamics and protonation equilibria in proteins. Of particular interest is the continuous constant pH method based on $\lambda$ dynamics and GB implicit models. This method has enabled accurate and robust $pK_a$ predictions for proteins, and simulations of pH-coupled protein folding from first principles

**Keywords:**     Molecular dynamics, $pK_a$, Protein folding, pH-dependent conformational change

## 10.1.     INTRODUCTION

Protein can gain or loose protons in response to variations in solution acidity or pH. Protonation or deprotonation leads to a change in the charge properties, which in turn affects the stability and function of the protein. Protein folding occurs in a specific pH range, a deviation from which may facilitate misfolding or aggregation. A well-known example is the tetrameric transthyretin, natively folded at pH 7.5. At pH 3, it undergoes dissociation into a monomeric intermediate, which proceeds to form amyloid fibrils that are linked to clinical syndromes, senile systemic amyloidosis and familial amyloid polyneuropathy I [39]. The ion-exchange function of an integral membrane protein, the $Na^+/H^+$ antiporter is regulated by pH. This antiporter protein comprises several helices which form cytoplasmic and periplamic funnels. At alkaline pH, a conformational change is induced such that two of the helices are reoriented, exposing the $Na^+$-binding site to the cytoplasmic funnel. Binding of $Na^+$-ions in turn results in opening of the periplasmic funnel leading to release of cations [36]. The catalytic function of enzymes is dependent on pH. For instance,

the pH for maximum protein degradtion activity of the enzyme pepsin is 2 while pancreatic lipase which catalyzes the hydrolysis of fat molecules has an optimum pH of 8.

Most frequently, protonation or deprotonation of protein occurs at ionizable, also called titratable, side chains, such as aspartate, glutamate or histidine. The ionization equilibrium of a titratable site,

$$HA \rightleftharpoons A^- + H^+ \tag{10-1}$$

is characterized with the equilibrium constant $K_a$ defined as

$$K_a = \frac{[A^-][H^+]}{[HA]} \tag{10-2}$$

The p$K_a$ value of a side chain is related to $K_a$ by

$$pK_a = -\log K_a \tag{10-3}$$

Combining Eq. (10-2) and Eq. (10-3) gives the Henderson-Hasselbach equation:

$$pH = pK_a + \log \frac{[A^-]}{[HA]} \tag{10-4}$$

Thus, the p$K_a$ value equals pH if the concentrations of protonated and deprotonated states are identical. For this reason, p$K_a$ is also called p$K_{1/2}$. Equation (10-4) can be rearranged and written in a generalized form in the presence of multiple titration sites,

$$S^{deprot} = \frac{1}{1 + 10^{n(pK_a - pH)}} \tag{10-5}$$

where $S^{deprot}$ denotes the deprotonated fraction and $n$ is the Hill coefficient, reflecting the extent of cooperativity in protonation between multiple titration sites. If $n$ is greater than 1, protonation is cooperative. If $n$ is smaller than 1, protonation is anti-cooperative. If $n$ equals 1, protonation is independent. The Hill coefficient was originally derived for characterizing cooperative binding of ligands to a receptor [105]. A plot with $S^{deprot}$ vs. pH is often known as the titration curve.

Comparison of solution pH with the p$K_a$ of a side chain informs about the protonation state. A unique p$K_a$, termed the standard or model p$K_a$, can be experimentally determined for each ionizable side chain in solution when it is incorporated in a model compound, often a blocked amino acid residue [73] (Table 10-1). In a protein environment, however, the p$K_a$ value of an ionizable side chain can substantially deviate from the standard value, due to desolvation effects, hydrogen bonding, charge-charge, charge-dipole, and other electrostatic interactions with the

*Table 10-1.* Standard p$K_a$ values for titratable side chains in a protein [85, 97]

| Side chain | p$K_a$ | Side chain | p$K_a$ |
|---|---|---|---|
| Asp | 4.0 | Cys | 10.5 |
| Glu | 4.4 | Tyr | 10.1 |
| His$^\delta$ | 6.6 | Lys | 10.8 |
| His$^\varepsilon$ | 7.0 | Arg | 12.5 |
| C-ter | 3.8 | N-ter | 7.5 |

surrounding polar groups. Experimental determination of abnormally shifted p$K_a$'s is tedious and sometimes even impossible. For instance, variation in pH may induce large conformational changes or even unfolding of the protein, which makes it impossible to obtain extremely shifted p$K_a$'s by NMR titrations. This is because the measurement relies on the observation of chemical shifts for both protonated and deprotonated states, but the baseline representing the fully protonated state can not be obtained if a protein unfolds at a pH that is higher than the expected p$K_a$ [75]. Experimental p$K_a$ determination of integral membrane proteins is also nontrivial because standard techniques for water soluble proteins can not be applied [32]. In addition to being highly desirable in these cases, theoretical p$K_a$ studies can offer the structural and energetic basis for p$K_a$ shifts that are often not accessible by experimental techniques [32].

## 10.2. ABSOLUTE p$K_a$ FOR SMALL MOLECULES

In principle, one can calculate p$K_a$ by making use of the thermodynamic cycle as shown in Figure 10-1 and a relationship between the free energy of deprotonation in aqueous phase and p$K_a$,

$$\Delta G_{aq}^{\text{deprot}} = \ln(10) RT \, pK_a \qquad (10\text{-}6)$$

The p$K_a$ obtained this way is often referred to as the absolute p$K_a$. Thus, p$K_a$ can be computed if the gas-phase deprotonation free energy and solvation energies of AH, A$^-$ and H$^+$ are known [60, 81].



*Figure 10-1.* Thermodynamic cycle for the absolute p$K_a$ calculation of small molecules. Subscripts gas and aq denote gas- and aqueous-phase, respectively

$$\Delta G_{aq}^{dep} = \Delta G_{gas}^{dep} - \Delta G^{solv}(HA) + \Delta G^{solv}(A^-) + \Delta G^{solv}(H^+) \qquad (10\text{-}7)$$

This approach may be possible for small molecules because the gas-phase proton affinity can be obtained quantum mechanically with an accuracy of 1–2 kcal/mol [19]. However, the solvation free energy of $H^+$ cannot be calculated and the experimental value is only known approximately, from 259.5 to 262.5 kcal/mol [60]. Also, because the proton affinity and solvation free energies in Eq. (10-7) are on the order of hundreds of kcal/mol, small percentage errors in the calculation can give rise to large error in $\Delta G_{aq}^{dep}$ and p$K_a$. Thus, this method for calculation of absolute p$K_a$'s remains impractical at the present time [6].

Another method to obtain absolute p$K_a$'s for small molecules is to compute deprotonation free energy directly from the free energies of species in the reaction using quantum mechanical and continuum solvation methods (Eq. 10-1),

$$\Delta G_{aq}^{deprot} = G(A^-) - G(HA) + [G(H_3O^+) - G(H_2O)] \qquad (10\text{-}8)$$

where the free energy difference between $H_3O$ and $H_2O$ is treated as a constant [48]. Klamt et al. applied density functional theory [77] and a continuum solvent model, the COSMO solvation method [47], with the TURBOMOLE/COSMO program [86, 101] to compute the free energy of deprotonation for a large variety of organic and inorganic molecules. They found that the deprotonation free energies were linearly correlated with the experimental p$K_a$'s,

$$pK_a = A\frac{\Delta G_{aq}^{deprot}}{\ln(10)RT} + B \qquad (10\text{-}9)$$

with a slope of about 0.58. This result was surprising because the expected slope is 1 if calculation errors were minimum.

## 10.3.    p$K_a$ CALCULATION FOR BIOLOGICAL MACROMOLECULES

Over the past two decades, significant progress has been made in the development of theoretical methods for the calculation of p$K_a$ for proteins. Unlike in the calculation of absolute p$K_a$'s, p$K_a$ for a protein side chain is obtained by computing the difference in deprotonation free energy ($\Delta\Delta G$) in the protein environment and that in the corresponding model compound (Figure 10-2). A fundamental assumption made in these methods is that quantum mechanical effects in the deprotonation energy, which include bond breaking and other electronic effects are approximately identical for the protein and the model compound. Consequently, $\Delta\Delta G$ can be obtained from computing the free energy difference due to nonbonded interactions of the ionizable site with the protein environment *vs.* that in solution. Below we will discuss

$$HA_{(mod)} \underset{\Longleftarrow}{\overset{\Delta G_{(mod)}}{\Longrightarrow}} A^-_{(mod)} + H^+$$

$$\Bigg| \quad \Delta\Delta G = \Delta G_{(prot)} - \Delta G_{(mod)} \quad \Bigg|$$

$$HA_{(prot)} \underset{\Longleftarrow}{\overset{\Delta G_{(prot)}}{\Longrightarrow}} A^-_{(prot)} + H^+$$

*Figure 10-2.* Thermodynamic cycle for the p$K_a$ calculation of proteins. Subscripts mod and prot refer to the model compound and protein, respectively

the development of two types of widely used methods, which differ in the way the dielectric environment of protein is modeled.

### 10.3.1.    Methods Based on a Macroscopic Description of Protein

Since electrostatic effects dominate the thermodynamic cycle as shown in Figure 10-2, major development efforts have focused on the calculation of electrostatic energy for transferring the neutral and charged forms of the ionizable group from water with dielectric constant of about 80 to the protein with a low dielectric constant (see later discussions). This led to the development of continuum based models, where water and protein are described as uniform dielectric media, and enter into the linearized Poisson-Boltzmann (PB) electrostatic equation,

$$\nabla \cdot [\varepsilon(\mathbf{r})\nabla\phi(\mathbf{r})] - \kappa^2(\mathbf{r})\varepsilon(\mathbf{r})\phi(\mathbf{r}) = -4\pi\rho_0(\mathbf{r}) \qquad (10\text{-}10)$$

where $\varepsilon(\mathbf{r})$ is the dielectric constant, $\phi(\mathbf{r})$ is the electrostatic potential, and $\rho_0(\mathbf{r})$ is the fixed solute charge density. $\kappa(\mathbf{r})$ is related to the ionic strength $I$ of the solution by $\kappa^2(\mathbf{r}) = 8\pi e I/\varepsilon(\mathbf{r})kT$. This second-order differential equation can be solved analytically for a spherical protein as in the Tanford-Kirkwood model [96]. Over the past two decades, efficient algorithms based on finite-difference approach have been developed for solving the PB equation with an arbitrarily shaped dielectric boundary, such as the molecular or van der Waals surface [7, 22, 49, 72]. These algorithms were implemented in several computer programs such as DELPHI [72], MEAD [5], and UHBD [63] that have been widely used for solving the PB equation and computing p$K_a$'s for proteins [2, 8, 106].

In PB based methods, the p$K_a$ for a protein ionizable site $i$ is given by [6],

$$pK_a^{prot} = pK_{int} - \frac{1}{\ln(10)RT}\Delta G_{Coul} \qquad (10\text{-}11)$$

where $pK_{int}$ is the intrinsic p$K_a$ according to Tanford. This is the p$K_a$ an ionizable group would have if all other groups are held in the neutral state [6]. $\Delta G_{Coul}$

represents the Coulomb interaction between other charged sites in the protein. $pK_{int}$ can be written in reference to the model $pK_a$ ($pK_a^{mod}$) as

$$pK_{int} = pK_a^{mod} - \frac{1}{\ln(10)RT}(\Delta G_{Born} + \Delta G_{bg}) \qquad (10\text{-}12)$$

where $\Delta G_{Born}$ is the Born or reaction field energy describing the difference between the self energy of each charge in the protein and in water, and $\Delta G_{bg}$ represents the interaction energy between the titrating group and other charge sites in the protein. The division in energetic contributions allows one to analyze the molecular determinants of the calculated $pK_a$ values [32].

Another advantage of PB based $pK_a$ calculations is that effects of electrolytes are readily accounted for in the PB equation. The Coulombic contribution in conjunction with salt dependence to the abnormally depressed $pK_a$'s of histidine in staphylococcal nuclease has been experimentally tested [56]. Recently, the methodology used in the PB calculations (Eqs. 10-11 and 10-12) has been combined with the generalized Born (GB) implicit solvent model [94] to offer $pK_a$ predictions at a reduced computational cost [52].

An inherent problem to methods that rely on a dielectric continuum description of protein is the lack of representation of the dielectric relaxation in protein, which arises from dipole movement, charge fluctuations, local rearrangement of polar side groups, and in some extreme case, global unfolding. Another mechanism for dielectric relaxation in protein is through conformational reorganization upon ionization. To compensate for the lack of explicit treatment of dielectric response, an effective dielectric constant $\varepsilon_{eff}$ for protein is used in the electrostatic calculation. The value of $\varepsilon_{eff}$ can be empirically adjusted, typically between 4 and 20 [2]. However, as noted by García-Moreno and Fitch, no single dielectric constant is able to reproduce experimental $pK_a$'s for both internal and surface residues [32]. This can be explained by the heterogeneous nature of the dielectric property as well as the conformational flexibility of the protein [90]. In fact, using molecular simulations, Simonson and Brooks showed that the dielectric constant in a protein changes abruptly from 2 to 3 in the interior, similar to the experimental value for dry protein powders [12, 35], to a significantly higher value, 14–25, in the outer region [90]. Mehler suggested that in polar solvent, an effective dielectric function, $\varepsilon_{eff}(\mathbf{r})$, that has a sigmodal shape, is better suited for describing charge-charge interactions [64]. For the interested reader, the meaning of protein dielectric constant has been extensively discussed by Warshel and coworkers [87, 104].

Several remedies have been suggested for improving the PB based $pK_a$ prediction methods. Most of them are based on strategies that combine conformational flexibility with the PB calculation. You and Bashford included multiple conformers by systematically scanning the side chain torsion angles [107]. Alexov and Gunner used Monte-Carlo protocol to sample positions of hydroxyl and other polar protons [1]. This method, referred to as the multi-conformation continuum electrostatic (MCCE), was later extended to include rotamers for residues that have strong electrostatic

interactions, and gave a root mean square error below 1 p$K$ unit for Asp, Glu, His, Lys, Tyr, N-ter and C-ter residues in 12 proteins using a single $\varepsilon_{eff}$ of 4 [33]. Other ways to account for protein dielectric response due to conformational variations are to average p$K_a$'s and to use an average structure from a short-time molecular dynamics trajectory [51, 99]. Finally, Mehler and Guarnieri addressed the problem with the heterogeneity of protein dielectric environment by devising a sigmoidally screened Coulomb potential parameterized to characterize the hydrophobicity and hydrophilicity around titratable residues [65]. This method was able to reproduce the divergent p$K_a$'s for two buried residues in lysozyme [65] and some of the largest p$K_a$ shifts (up to 4.7 unit) in staphylococcal nuclease [66]. Despite the encouraging results demonstrated by the above mentioned development, a fundamental issue remains, namely, the coupling between conformational dynamics and ionization equilibria is missing in the macroscopic approaches based on static structures.

### 10.3.2. Methods Based on a Microscopic Description of Protein

The linkage of protein conformation to ionization of titratable residues on a local level, such as the ionization induced reorientation of nearby polar groups [24, 53, 54] can be addressed by p$K_a$ calculation using free energy simulation techniques such as the free energy perturbation (FEP) method [68, 80, 91, 103]. In this approach, the free energy of the charging process in protein environment is obtained using molecular dynamics with either explicit or continuum solvent models. In the same spirit, another method is to apply the Gaussian fluctuation formula derived from linear response theory [59] in molecular dynamics simulations with explicit solvent [23]. Warshel et al. developed a method [88] base on the protein dipoles Langevin dipoles (PDLD) model, which considers protein to be surrounded by point dipoles in a grid representing solvent molecules [55]. In this method, the p$K_a$ calculation follows Eqs. (10-11) and (10-12) as in PB based methods. However, none of the above mentioned methods can be applied to situations, where ionization is coupled to a large conformational transition, such as unfolding of protein [6]. In this case, methods have to developed that can simultaneously describe the conformational dynamics of protein and the protonation process of titratable residues. Below we will outline recent progress in the development of these methods and the application to p$K_a$ calculations and the study of pH-coupled conformational processes in proteins.

### 10.3.2.1. *Constant pH Molecular Dynamics Based on Discrete Protonation States*

In recent years, a class of methods has been developed for molecular dynamics simulations to be performed with an external pH parameter, like temperature or pressure [18, 43, 44, 70]. These methods treat the solution as an infinite proton bath, and are thus referred to as constant pH molecular dynamics (PHMD). In PHMD, conformational dynamics of a protein is sampled simultaneously with the protonation states as a function of pH. As a result, protein dielectric response to the

correlated event of charging and conformational dynamics is explicitly taken into account. Also, a physical description for the coexistence of the charged and neutral states is enabled. Currently, there are two types of PHMD methods. In the methods based on discrete protonation states, also known as stochastic titration, molecular dynamics (MD) is periodically interrupted by Monte-Carlo (MC) sampling of protonation states [4, 16, 71, 102]. After the MD sampling of conformational states, a MC step generates a Boltzmann distribution of protonation states at a given pH, which in turn, affects the generation of new conformational states. The procedure repeats itself until convergence, at which point the coupling between conformational and protonation equilibria is established. The discrete protonation states methods developed so far differ mainly in the solvent model used for the MD and MC sampling but also in the protocol for the update of protonation states as well as the method for evaluating the free energy of deprotonation (relative to that of the model compound value) in the MC step.

Baptista et al. developed a mixed solvent scheme [4], in which an explicit solvent MD simulation was combined with MC sampling using the free energy obtained from the PB calculation (Eq. 10-11). After a switch in protonation states, explicit solvent molecules are allowed to relax at the fixed solute conformation. This step is necessary to alleviate a problem due to the abrupt change in the charge states, which is a potential pitfall for all discrete states methods. An instantaneous switch in charge states may lead to a large increase in unfavorable energy, which may result in a low acceptance ratio in the MC move (see later discussions) [93]. The effect due to mobile counter ions can be trivially included in the PB calculation. However, it is more problematic in the sampling of conformational states using explicit solvent, because the total charge of the solute is changing and unknown at the start of the simulation. Machuqueiro and Baptista tested two approaches [61], the generalized reaction field (GRF) method [98] which can account for the effects due to ionic strength, and the Particle mesh Ewald (PME) method [21] with inclusion of an approximate number of counter ions. A most recent study showed that the GRF method gives better p$K_a$ prediction for hen egg lysozyme [62].

To explore the explicit solvent scheme for both conformational and protonation states sampling, Bürgi et al. employed a free energy simulation technique, thermodynamic integration (TI), to estimate the deprotonation free energy for the MC trial moves [16]. However, since free energy calculation is prohibitively expensive, the TI sampling could be performed only for a very short period of time (20 ps in the work of Bürgi et al. [16]), giving rise to a large uncertainly in the estimated free energy [70]. Most recently, Stern proposed a method, in which MC trial moves consist of short molecular dynamics runs with a time-dependent Hamiltonian that interpolates between the old and new protonation states [93]. This method may offer reduction in computational cost as compared to the one based on the free energy calculation. Also, because the trial MC move is not instantaneous, the problem with discontinuous force is circumvented.

The discrete protonation states methods employing implicit solvent models in both MD and MC steps have significantly lower computational cost. Dlugosz and

Antosiewicz combined the implicit solvent model, Analytical Continuum Electrostatics [84], for conformational sampling with the PB calculation [25]. To speed up simulations even further, Mongan et al. developed a protocol that employs a GB implicit solvent model for both conformational and protonation states sampling [71]. Another advantage of using implicit solvent model for conformational sampling is that the problem due to discontinuous energy may be largely circumvented because of the instantaneous adjustment of the solvent to new protonation states. Another trick that may ease the problem is to change the protonation state of one site at each MC step although this may slow the convergence for strongly coupled titrations [70, 71]. Besides the pitfall due to discontinuous energy, a discrete protonation states simulation is significantly more expensive than a standard MD simulation, because of the extra computational time spent on energy evaluations in the MC trial moves.

The accuracy of PHMD methods and their feasibility for studying pH-dependent conformational phenomena of proteins can be assessed by $pK_a$ calculations. In this case, PHMD simulations are performed with several pH values. The resulting occupancy values for deprotonated states ($S^{\text{de prot}}$) are plotted against pH (Figure 10-3). A titration curve and $pK_a$ values (Figure 10-3) can be obtained by fitting the data to the generalized HH equation (Eq. 10-5).

The discrete protonation states methods have been tested in $pK_a$ calculations for several small molecules and peptides, including succinic acid [4, 25], acetic acid [93], a heptapeptide derived from ovomucoid third domain [27], and decalysine [61]. However, these methods have sofar been tested on only one protein, the hen egg lysozyme [16, 61, 71]. While the method using explicit solvent for both MD and MC sampling did not give quantitative agreement with experiment due to convergence difficulty [16], the results using a GB model [71] and the mixed PB/explicit



*Figure 10-3.* Theoretical titration curves for the model compounds of Asp and His obtained from REX-CPHMD simulations [41]. *Solid curves* are the obtained by fitted the computed deprotonated fraction to the generalized Henderson-Hasselbach equation. The dashed lines indicate the computed $pK_a$ values

solvent scheme [61] are encouraging. Machuqueiro and Baptista performed 30 ns titration simulations with a MC trial move every 2 ps followed by 0.2 ps solvent relaxation [61]. The overall root-mean-square deviation between the computed and experimental p$K_a$ values was below 1 p$K$ unit.

### 10.3.2.2.   *Constant pH Molecular Dynamics Based on Continuous Protonation States*

In an early work by Mertz and Pettitt, an open system was devised, in which an extended variable, representing the extent of protonation, was used to couple the system to a chemical potential reservoir [67]. This method was demonstrated in the simulation of the acid-base reaction of acetic acid with water [67]. Recently, PHMD methods based on continuous protonation states have been developed, in which a set of continuous titration coordinates, $\lambda_i$, bound between 0 and 1, is propagated simultaneously with the conformational degrees of freedom in explicit or continuum solvent MD simulations. In the acidostat method developed by Börjesson and Hünenberger for explicit solvent simulations [13], $\lambda_i$ is relaxed towards the equilibrium value via a first-order coupling scheme in analogy to Berendsen's thermostat [10]. However, the theoretical basis for the equilibrium condition used in the derivation seems unclear [3]. A test using the p$K_a$ calculation for several small amines did not yield HH titration behavior [13].

   An alternative formalism to the acidostat method, continuous constant pH molecular dynamics (CPHMD), has been developed by Brooks and coworkers [42, 58] drawing on the flavor of the early work by Mertz and Pettitt [67]. In CPHMD, each titration residue carries a titration coordinate, $\lambda_j$, which is a function of an unbounded variable $\theta_j$,

$$\lambda_j = \sin^2(\theta_j) \tag{10-13}$$

The titration coordinates evolve along with the dynamics of the conformational degrees of freedom, $r_i$, in simulations with GB implicit solvent models [37, 57]. An extended Hamiltonian formalism, in analogy to the $\lambda$ dynamics technique developed for free energy calculations [50], is used to propagate the titration coordinates. The deprotonated and protonated states are those, for which the $\lambda$ value is approximately 1 or 0 (end-point states), respectively. Thus, in contrast to the acidostat method, where $\lambda$ represents the extent of deprotonation, $S^{\text{deprot}}$ is estimated from the relative occupancy of the states with $\lambda \approx 1$ (see later discussions). The extended Hamiltonian in the CPHMD method is a sum of the following terms [42].

$$H^{\text{ext}}(\{r_i\}, \{\theta_j\}) = H^{\text{hyb}}(\{\mathbf{r}_j\}, \{\theta_j\}) + \sum_i \frac{m_j}{2} \dot{\theta}_j^2 + U^*(\{\theta_j\}) \tag{10-14}$$

The hybrid Hamiltonian, which depends on both spatial and titration coordinates, can be written in terms of van der Waals, Coulomb, and GB electrostatic energies,

$$H^{\mathrm{hyb}}(\{\mathbf{r}_i\}, \{\theta_j\}) = U^{\mathrm{vdW}}(\{\mathbf{r}_i\}, \{\theta_j\}) + U^{\mathrm{elec}}(\{\mathbf{r}_i\}, \{\theta_j\}) + U^{\mathrm{GB}}(\{\mathbf{r}_j\}, \{\theta_j\}) \quad (10\text{-}15)$$

where the van der Waals energy for the interaction between the titrating proton $j$ and another atom $i$ is, given by

$$U^{\mathrm{vdW}}(j, i) = \begin{cases} (1 - \lambda_j)\tilde{U}^{\mathrm{vdW}}(j, i) & i: \text{other atom} \\ (1 - \lambda_j)(1 - \lambda_i)\tilde{U}^{\mathrm{vdW}}(j, i) & i: \text{titrating proton} \end{cases} \quad (10\text{-}16)$$

where $\tilde{U}^{\mathrm{vdW}}(j, i)$ represents the protonation independent (full-strength) van der Waals interaction energy. The dependence of the Coulomb and GB electrostatic energies on $\lambda$ is incorporated through a linear attenuation of partial charges.

$$q_{\alpha, j} = \lambda_j q_{\alpha, j}^{\mathrm{deprot}} + (1 - \lambda_j)q_{\alpha, j}^{\mathrm{prot}} \quad (10\text{-}17)$$

where $q_{\alpha, j}^{\mathrm{deprot}}$ and $q_{\alpha, j}^{\mathrm{prot}}$ are the partial charge of atom $\alpha$ from the titrating residue $j$ in the deprotonated and protonated states, respectively.

The second term in Eq. (10-14) represents the kinetic energy of $\theta_j$, which has a fictitious mass similar to that of a heavy atom. The last term represents the sum of the biasing potential for each titrating group, defined as

$$U^*(\{\theta_j\}) = \sum_j [-U^{\mathrm{bar}}(\theta_j) - U^{\mathrm{mod}}(\theta_j) + U^{\mathrm{pH}}(\theta_j)] \quad (10\text{-}18)$$

For the sake of convenience, we will express the biasing potentials as a function of $\lambda$ below. $U^{\mathrm{bar}}(\theta_j)$ is a harmonic potential,

$$U^{\mathrm{barr}}(\lambda_j) = 4\beta_i \left( \lambda_j - \frac{1}{2} \right)^2 \quad (10\text{-}19)$$

that equally lowers the energy of the end-point states. This term is necessary for suppressing the population of mixed states $(0 < \lambda < 1)$, which are unphysical and could introduce distortion in the presence of strongly coupled titration groups. Previous studies [41–46] indicated that the barrier height $\beta_j$ can be set to 1.5 kcal/mol (or 2.5 kcal/mol for double-site titration, see later discussions), which offers a reasonable tradeoff between the transition rate of protonation states and a low population (below 15%) of mixed states.

$U^{\mathrm{mod}}(\theta_j)$ represents the potential of mean force (PMF) for deprotonating a model compound along the titration coordinate. This ensures that the titration simulation of a model compound at pH $= \mathrm{p}K_a^{\mathrm{mod}}$, yields approximately 50% protonated and 50% deprotonated states. In other words, the PMF along the titration coordinate for a model compound is flattened out at pH $= \mathrm{p}K_a^{\mathrm{mod}}$, thus allowing us to model only the difference between the free energy of deprotonation in protein and that in

solution (Figure 10-2). Due to the pair-wise form of the GB energy function, $U^{\text{mod}}$ is a quadratic function of $\lambda$,

$$U^{\text{mod}}(\lambda_j) = A_j(\lambda_i - B_j)^2 \tag{10-20}$$

The parameters, $A_j$ and $B_j$, can be obtained via thermodynamic integration using titration simulations of the model compound at different $\theta$ values [58]. Finally, $U^{\text{pH}}(\theta_i)$ models the free energy dependence on the external pH by

$$U^{\text{pH}}(\lambda_j) = \ln 10 \cdot k_{\text{B}} T \cdot \lambda_j (pK_j^{\text{mod}} - pH) \tag{10-21}$$

where $pK_a(j)$ is the model $pK_a$. Thus, the ratio between the protonated and deprotonated states (the protonation equilibrium, Eq. (10-1)) is controlled by the difference between the simulation pH and the standard $pK_a$'s, and between the nonbonded environment in a protein and that in solution.

The formalism introduced so far treats the protonation equilibria of different sites independently, and can be referred to as the single-site model. This is, however, not realistic for histidine and carboxyl side chains of Asp, Glu and C-terminal residues. Histidine can loose or gain a proton from either Nδ or Nε, while protonation/deprotonation can occur at either of the carboxylate oxygens in Asp. Another way of viewing the competitive titration problem is to consider the titration products as tautomers, although in the case of a carboxylate group, the two protonated forms are chemically equivalent. Nevertheless, we can treat them as tautomers because rotation of the C—O bond in a carboxylate group is slow on a MD time scale (Figure 10-4). In order to couple the protonation equilibria of the tautomers, a double-site model was developed, in which a new extended variable, $x_j$, describing tautomeric interconversion, is introduced and is treated in the same manner as the



*Figure 10-4.* The double- and single-site titration models for His and Asp groups [42]. (**A**) In the double site model, only one $\lambda$ is used for describing the equilibrium between the protonated and deprotonated forms, while the tautomer interversion process is represented by the variable $x$. (**B**) In the single-site model, protonation at different sites is represented by different $\lambda$ variables. HSP refers to the doubly protonated form of histidine. HSD and HSE refer to the singly protonated histidine with a proton on the δ and ε nitrogens, respectively. ASP1 and ASP2 refer to the protonated carboxylic acid with a proton on either of the carboxlate oxygens

fictitious $\lambda$ particles (Figure 10-4) [42]. The double-site model can be viewed as a special case in the two-dimensional $\lambda$ dynamics formalism [42]. Thus, the energy terms for nonbonded interactions need to be reformulated to include the dependence of $x$. Also, the PMF for titrating the model compound is no longer a simple quadratic function of $\lambda$, but rather, it becomes a bivariate polynomial, quadratic in either $\lambda$ or $x$ [42]. The exact functional form can be obtained by thermodynamic integration at different $\lambda$ or $x$ values [42]. In principle, the two-dimensional $\lambda$ dynamics technique can be extended to devise a multiple-site titration model for the titration of the amino group in lysine. However, the necessity for developing such a sophisticated model is unclear. The ionization equilibrium of lysine is typically not considered in studies under cellular (pH=4.5–8) or physiological (pH≈7.5) conditions because its known $pK_a$'s in proteins are mostly above 10 [65]. Currently, the CPHMD method is implemented in the CHARMM molecular dynamics program package (version c33a1 and up) [15].

### 10.3.2.3.    *De Novo Prediction of $pK_a$ Values in Proteins*

One major bottleneck in the development of methods based on a microscopic description of protein is convergence [70]. Given the current CPU power, de novo prediction of protein $pK_a$'s in explicit solvent is practically unattainable using methods based on either discrete or continuous protonation states due to poor convergence. The current CPHMD implementation utilizes the recently developed implicit solvent GB models, GBMV [57] and GBSW [37] with a simple solvent accessible surface area approximation for modeling the nonpolar solvation effects. These GB models offer high accuracy (less than 1% error with respect to PB solvation energies for a large set of proteins [30]) but are significantly less computationally demanding relative to PB calculations or explicit solvent simulations. Nevertheless, random errors in the CPHMD titration simulations of model compounds were as large as 0.5 p$K$ units, indicating the lack of convergence [42]. Considering the fact that protonation states are very sensitive to local conformational variations, the CPHMD method was combined with the replica-exchange (REX) conformational sampling technique [95], in order to enhance the protonation states sampling. In a REX algorithm, multiple independent copies of the system are simulated in parallel at different temperatures. Conformational states at adjacent temperatures are allowed to swap periodically based on the Monte Carlo criterion, thus enhancing barrier crossing in the potential energy surface. The REX technique has been applied to significantly speed up the first principles protein folding simulations [74]. Indeed, by incorporating the REX algorithm, the magnitude of random errors was reduced to below 0.16 p$K$ units in model compound titrations with CPHMD using the same amount of total simulation time [41]. The replica-exchange protocol was implemented in a PERL based package, MMTSB tool set [29] (http://blue11.bch.msu.edu/mmtsb), which was interfaced with the CHARMM program to enable REX-CPHMD simulations.

The accuracy of p$K_a$ predictions is intimately linked to that of the underlying solvent model. Overstabilization of attractive electrostatic interactions in GB models [38] led to a systematic underestimation of p$K_a$'s for carboxyl residues that interact with positively charged groups [42]. By employing an optimized set of atomic input radii for the GBSW model and the improved torsion energetics for the implicit solvent force field [17], significant improvement was demonstrated [41]. In PB based methods, effects due to mobile ions are naturally included in solving the PB equation (Eq. 10-10). In GB based methods, however, one has to resort to some approximation. One way is to apply a scaling function, $e^{\kappa r}$, to the solvent dielectric constant $\varepsilon$, where $\kappa$ depends on the ionic strength of the solution and is defined earlier [92]. Although the extent of salt dependence recovered by this crude treatment is still unclear; and may be somewhat larger than that from the PB calculation [92], REX-CPHMD simulations that incorporate the approximate screening function yield p$K_a$ results in better agreement with experiments that were conducted in 100–300 mM salt solution [45].

By combining with the REX sampling protocol, employing an improved GB model, and accounting for salt effects, REX-CPHMD simulations are able to offer quantitative prediction of protein p$K_a$'s. A benchmark study demonstrated that 1-ns REX-CPHMD simulations give results with an RMS error below 1 p$K_a$ unit for a stringent test set of proteins, in which anomalously large p$K_a$ shifts are observed in carboxyl and histidine residues [45]. The ability of predicting the ionization equilibria for these residues is important, because they are often involved in the catalytic activity of enzymes [28, 31] and functional processes driven by proton gradients [9, 36, 40]. One of the remaining challenges is to improve the p$K_a$ prediction for deeply buried sites. In fact, buried side chains gave the largest error in the benchmark study using the REX-CPHMD technique [45]. This may be attributed to the underestimation of desolvation energies of buried residues as well as buried charge-charge interactions, because the employed GB model uses a dielectric boundary (van der Waals surface), which neglects solvent excluded volumes [41]. Furthermore, the efficiency of REX-CPHMD simulations can be greatly enhanced by coupling both temperature and pH in a two-dimensional version of the REX protocol. Nonetheless, this [41] and other on-going studies (Khandogin and Brooks III, unpublished results) clearly indicate that de novo p$K_a$ prediction by REX-CPHMD simulations is soon to become a routine task.

## 10.4.    pH-COUPLED MOLECULAR DYNAMICS SIMULATIONS OF PROTEINS

pH effects have been traditionally included in molecular simulations through the fixed charge state approach. The underlying assumptions are (1) p$K_a$'s of titratable side chains are known or can be approximated with the standard p$K_a$'s; (2) titrable side chains assume one of the charge states (charged or neutral); (3) their charge states remain fixed in the course of the simulation. Obviously, these assumptions

often do not reflect the physical reality. First, p$K_a$'s in a protein can be significantly shifted from the corresponding standard values. Second, when solution pH is close to the p$K_a$ of the titrable group, both charge states are populated in the protonation equilibrium. Lastly, the charge states of titratable groups fluctuate in response to environment. PHMD simulations offer a rigorous way of including solvent pH effects in molecular simulations, thus allowing for studies of pH-coupled conformational phenomena. Although the ability to model pH-coupled conformational transitions is limited by the bottleneck in conformational sampling and the accuracy of force field and solvent model, recent applications of PHMD to explore pH-dependent conformational properties and folding of proteins are encouraging.

### 10.4.1.    pH-Dependent Structural Properties in Peptides

Using stochastic titration simulations with implicit solvent, Dlugosz and Antosiewicz investigated the question regarding the significance of solvent-solute proton exchange for the structure of polypeptides under pH conditions such that both protonated and deprotonated states are occupied [26]. They showed that the distance distributions between Asp and Lys residues in a pentapeptide can not be reproduced by a linear combination of those from the fixed-protonation state simulations [26].

The origin of the pH-dependent helicity of peptides and proteins has been a subject for both experimental [11] and theoretical studies in the past two decades. An early work by Ripoll et al., which allows a coupling between ionization and conformational states using Monte Carlo simulation and the Multigrid Boundary Element method for solving PB equation, examined the pH-dependent helicity of a 17-residue peptide [83]. This study showed that conformations containing right-handed α helical segments are energetically more favorable at low pH, in contrast to simulations using fixed charge states. The difference was attributed to the restrictions imposed on the conformational space in the fixed-charge approach, which led to the exclusion of low-energy conformations [83]. A similar methodology was later used to study the pH-dependent helicity of a penta-peptide consisting of lysine and alanine residues, and resulted in a quantitative agreement with experiment [100]. More recently, the pH dependent helicity of decalysine was studied by Börjesson and Hünenberger using the acidostat method in explicit solvent [14], and by Machuqueiro and Baptista using the stochastic titration method with the mixed solvent scheme [61]. In the latter study [61], the authors compared the calculated helix profile as a function of pH with experiment. The overall shape of the helix profile vs. pH was reproduced although the helix-coil transition obtained from the simulation was not as abrupt as that from experiment. The above-mentioned simulation studies focused on the comparison of theoretical and experimental p$K_a$'s and total helix content. However, convergence of conformational properties was not demonstrated due to the lack of conformational sampling. This, in turn, made it difficult to gain detailed molecular mechanisms of the structural transitions.

### 10.4.2.     De Novo Simulations of pH-Dependent Proten Folding

Because of the quantitative accuracy in predicting protein $pK_a$'s and the advantage of being only marginally slower than standard GB simulations, which are significantly faster compared to explicit solvent simulations, CPHMD has become a powerful tool for gaining atomistic insights into a host of pH-dependent conformational phenomena. Khandogin and Brooks applied REX-CPHMD simulations to investigate the pH-dependent folding of the C-peptide from ribonuclease A, starting from an extended structure [45]. Although molecular mechanisms governing the pH-dependent folding behavior can be probed by spectroscopic measurements combined with amino acid mutations, much detailed information is inaccessible. For example, a bell-shaped pH profile of helix content for the C-peptide has been known for a long time [11, 76, 89]. However, questions regarding the composition of the conformational equilibrium and specific interactions responsible for the pH dependence remained elusive. The study by Khandogin and Brooks was able to reveal, not only a pH-dependent total helix content in agreement with experiment, but also a pH-dependent conformational equilibrium consisting of unfolded and partially folded states of various helical lengths [45]. This study was able to provide a direct correlation between the specific electrostatic interactions and stability of different helical conformations, thus offering atomic insights into the mechanism of pH-modulated helix-coil transitions [45]. The latter information is both consistent with and complementary to the existing experimental data.

De novo folding simulations based on REX-CPHMD were also applied to elucidate the folding behavior of β amyloid peptides (Aβ) from Alzheimers disease, Aβ(1–28), Aβ(10–42) [43, 44]. Unlike the C-peptide and many globular proteins, which are maximally folded at an intermediate pH that minimizes the total charge, Aβ is helical with an inverse bell-shaped pH profile in aqueous TFE solution but it is largely unfolded in water [43, 44]. The study of Khandogin and Brooks reconciled this discrepancy in folding by showing that Aβ is mainly coil-like but contains several short, nascent helical segments that have increased helix propensity under low and high pH conditions (Figure 10-5). Furthermore, this study revealed a pH-dependent solvent exposure in the central hydrophobic cluster, which forms the minimum fibril forming sequence, and a pH-dependent β-turn formation in residues 23–26. Taken together, these data predicted that, at pH 6, Aβ adopts conformational states that are most prone for the formation of β-sheet based aggregates. Finally, this study suggested that, although minimum charge-charge repulsion at the isoelectric point provides an initial driving force for aggregation, the folding conformational landscape of Aβ is strongly modulated by pH leading to an enhanced intermolecular hydrophobic association which stabilizes the β-sheet based oligomers.

REX-CPHMD simulations have also been applied to understand the mechanism of the formation of protein intermediate states. Recent solution NMR data revealed a sparsely populated intermediate in the villin headpiece domain, in which the N-terminal subdomain is largely random but the C-terminal subdomain adopts a native-like fold [34]. Interestingly, H41 in this intermediate state titrates at a pH value of

*Figure 10-5.* Representative conformations of the β amyloid peptide (10–42) under different pH conditions. The conformations were obtained as centroids of the most populated clusters from the replica-exchange CPHMD folding simulations [43, 44]. The N-terminal residues 10–28 are shown in *blue*; the C-terminal residues 29–42 are shown in *red*. In the most aggregation-prone state (pH 6), the side chains of the central hydrophobic cluster Leu-17, Val-18, Phe-19, Phe-20 and Ala-21 are shown as van der Waals spheres in *pink*, *grey*, *cyan*, *purple* and *green*, respectively

5.6, more than 1.5 units higher than that of the native state [34]. REX-CPHMD simulations initiated from an X-ray crystal and a minimized average solution NMR structure gave rise to two conformational states that have distinct titration behavior for H41, corresponding the measurements for the native and intermediate states [46]. Supported by the dynamical, structural and titration properties, the simulation data suggested that the state derived from the solution NMR structure resembles a putative intermediate that has a largely unfolded N-terminal subdomain. Moreover, the formation of the putative intermediate was thought to be the result of the loss of a hydrogen-bonded network centered at H41. Thus, this work put forth a proposal that the hydrogen-bonded network and not the protonation state of H41 is a prerequisite for folding in the villin headpiece domain.

Another interesting application area of PHMD simulations is to investigate electrostatic interactions in the unfolded states of proteins. A traditional view that unfolded proteins adopt random conformational states that are devoid of electrostatic and hydrophobic interactions, are recently challenged by experimental data [20, 69]. REX-CPHMD folding simulations of the 35 residue C-terminal subdomain of the villin headpiece domain revealed a significant deviation from the standard $pK_a$ values for several titratable residues. Additional simulations, in which a charged group is neutralized confirmed the existence of specific electrostatic interactions in the unfolded states (JK and CLB, manuscript in preparation).

## 10.5.    SUMMARY AND OUTLOOK

Many biological and chemical processes are modulated by solution salt and pH conditions. Over the past decade, significant progress has been made in the development of theoretical methods that explicitly account for the microscopic coupling between conformational dynamics and protonation equilibria in proteins. These methods, known as constant pH molecular dynamics, eliminate the need for the ad hoc assignment of protein dielectric constant in the $pK_a$ calculation as in the traditional Poisson-Boltzmann (PB) based approaches, and it can be applied to study pH-dependent conformational phenomena of proteins. Of particular interest are the stochastic titration method that combines the PB calculation sampling of protonation states and explicit solvent molecular dynamics for sampling of conformational states, and the continuous constant pH molecular dynamics (CPHMD) method based on λ dynamics and GB implicit solvent models. The accuracy and speed of the CPHMD method has enabled, for the first time, robust and quantitative $pK_a$ prediction for proteins on the first-principles level. By employing the replica-exchange algorithm to enhance the coupled conformational and protonation states sampling, the CPHMD method can compute $pK_a$ during protein folding, thus providing a powerful tool for the simulation of pH-dependent protein folding.

   Despite the above-mentioned success, the accuracy of the CPHMD method can be further improved through continued development of the underlying GB implicit solvent model and its parameterization. The underestimation of $pK_a$ shifts for deeply buried residues may be reduced by using a molecular surface that accounts for solvent excluded volumes. A more elaborate scheme relative to the current approximate function [41] may be considered to better capture salt screening effects. Another effect that may play an important role in further pushing the level of accuracy for $pK_a$ determinations is polarization. In this regard, the CPHMD technique can be combined with methods of treating polarization, such as the fluctuating charge model [82], for which an all-atom force field for protein simulations has been recently developed [78, 79]. Since the solvation energy is coupled with titration coordinates through modulation of partial charges, the CPHMD formalism can be combined, in a similar manner, with other implicit solvent models such as PB. However, since the PB approach requires numerical solutions, the potential of mean force function for model compound titration is no longer analytic and has to be approximated by a continuous function, through, for example, cubic spline fitting to grid data. The CPHMD formalism can also be extended to explicit solvent simulations although, with the current CPU and sampling capability, convergence is expected to be poor.

   The first series of applications of CPHMD to pH-dependent conformational processes has provided novel insights, inaccessible by traditional molecular simulations and current experiments, into the pH-modulated peptide folding, formation of intermediate, and the aggregation properties of amyloidogenic peptides. Ongoing studies (JK and CLB, unpublished data) indicate that the CPHMD method can be applied to investigate a variety of pH-dependent phenomena, such as peptide insertion into the membrane, and the backbone hydrogen bond registry shift in amyloid fibrils.

The development and application of the CPHMD method demonstrate that simulations are capturing physical reality at increasing resolution. With the explosion in computing technologies, we are just at the beginning of a new era, where in silico experimentation becomes an indispensable complement to wet lab experiments in exploring unanswered questions related to a wide variety of biological and chemical processes.

## REFERENCES

1. Alexov EG, Gunner MR (1997) Incorporating protein conformational flexibility into the calculation of pH-dependent protein properties. *Biophys J* 72:2075–2093.
2. Antosiewicz J, McCammon JA, Gilson MK (1994) Prediction of pH-dependent properties of proteins. *J Mol Biol* 238:415–436.
3. Baptista M (2002) Comment on "Explicit-solvent molecular dynamics simulation at constant pH: Methodology and application to small amines". *J Chem Phys* 116:7766–7768.
4. Baptista M, Teixeira VH, Soares CM (2002) Constant-pH molecular dynamics using stochastic titration. *J Chem Phys* 117:4184–4200.
5. Bashford D (1997) *Scientific Computing in Object-Oriented Parallel Environments: Lecture Notes in Computer Science*, volume 1343, chapter An object-oriented programming suite for electrostatic effects in biological molecules. Springer, Berlin, pp 233–240.
6. Bashford D (2004) Macroscopic electrostatic models for protonation states in proteins. *Front Bioscience* 9:1082–1099.
7. Bashford D, Gerwert K (1992) Electrostatic calculations of the p$K_a$ values of ionizable groups in bacteriorhodopsin. *J Mol Biol* 224:473–486.
8. Bashford D, Karplus M (1990) pK$_a$ s of ionizable groups in proteins: Atomic detail from a continuum electrostatic model. *Biochemistry* 29:10219–10225.
9. Belevich I, Verkhovsky MI, Wikström M (2006) Proton-coupled electron transfer drives the proton pump of cytochrome c oxidase. *Nature* 440:829–832.
10. Berendsen HJC, Postma JPM, van Gunsteren WF, Dinola A, Haak JR (1984) Molecular dynamics with coupling to an external bath. *J Chem Phys* 81:3684–3690.
11. Bierzynski, Kim PS, Baldwin RL (1982) A salt bridge stabilizes the helix formed by isolated C-peptide of RNase A. *Proc Natl Acad Sci USA* 79:2470–2474.
12. Bone S, Pethig R (1985) Dielectric studies of protein hydration and hydration-induced flexibility. *J Mol Biol* 181:323–326.
13. Börjesson U, Hünenberger PH (2001) Explicit-solvent molecular dynamics simulation at constant pH: Methodology and application to small amines. *J Chem Phys* 114(22):9706–9719.
14. Börjesson U, Hünenberger PH (2004) pH-dependent stability of a decalysine α-helix studied by explicit-solvent molecular dynamics simulations at constant pH. *J Phys Chem B* 108:13551–13559.
15. Brooks R, Bruccoleri RE, Olafson BD, States DJ, Swaminathan S, Karplus M (1983) Charmm: A program for macromolecular energy minimization and dynamics calculations. *J Comput Chem* 4:187–217.
16. Bürgi R, Kollman PA, van Gunsteren WF (2002) Simulating proteins at constant pH: An approach combining molecular dynamics and Monte Carlo simulation. *Proteins* 47:469–480.
17. Chen J, Im W, Brooks CL III (2006) Balancing solvation and intramolecular interactions: Toward a consistent generalized Born force field. *J Am Chem Soc* 128:3728–3736.

18. Chen J, Brooks CL III, Khandogin J (2008) Recent advances in implicit solvent based methods for biomolecular simulations. *Curr Opin Struct Biol* 18:140–148.

19. Chipman M (2002) Computation of pKa from Dielectric Continuum Theory. *J Phys Chem A* 106: 7413–7422.

20. Cho J-H, Raleigh DP (2005) Mutational analysis demonstrates that specific electrostatic interactions can play a key role in the denatured state ensemble of proteins. *J Mol Biol* 353:174–185.

21. Darden T, York D, Pedersen L (1993) Particle mesh Ewald: An N log(N) method for Ewald sums in large systems. *J Chem Phys* 98:10089–10092.

22. Davis ME, Madura JD, Luty BA, McCammon JA (1991) Electrostatics and diffusion of molecules in solution: Simulations with the University of Houston Brownian dynamics program. *Comput Phys Commun* 62:187–197.

23. Del Buono GS, Figueirido FE, Levy RM (1994) Intrinsic pKas of ionizable residues in proteins: An explicit solvent calculation for lysozyme. *Proteins* 20:85–97.

24. Dillet V, Dyson HJ, Bashford D (1998) Calculations of electrostatic interactions and pKas in the active site of *Escherichia coli* thioredoxin. *Biochemistry* 37:10298–10306.

25. Dlugosz M, Antosiewicz JM (2004) Constant-pH molecular dynamics simulations: A test case of succinic acid. *Chem Phys* 302:161–170.

26. Dlugosz M, Antosiewicz JM (2005) Effects of solute-solvent proton exchange on polypeptide chain dynamics: A constant-pH molecular dynamics study. *J Phys Chem B* 109:13777–13784.

27. Dlugosz M, Antosiewicz JM, Robertson AD (2004) Constant-pH molecular dynamics study of protonation-structure relationship in a heptapeptide derived from ovomucoid third domain. *Phys Rev E* 69:021915.

28. Edgcomb SP, Murphy KP (2002) Variability in the pKa of histidine side-chains correlates with burial within proteins. *Proteins* 49:1–6.

29. Feig M, Karanicolas J, Brooks CL, III (2004) MMTSB tool set: Enhanced sampling and multiscale modeling methods for applications in structure biology. *J Mol Graph Model* 22:377–395.

30. Feig M, Onufriev A, Lee MS, Im W, Case DA, Brooks CL III (2004) Performance comparison of generalized born and poisson methods in the calculation of electrostatic solvation energies for protein structures. *J Comput Chem* 25:265–284.

31. Forsyth WR, Antosiewicz JM, Robertson AD (2002) Empirical relationships between protein structure and carboxyl p$K_a$ values in proteins. *Proteins* 48:388–403.

32. García-Moreno EB, Fitch CA (2004) Structural interpretation of pH and salt-dependent processes in proteins with computational methods. *Methods Enzymol* 380:20–51.

33. Georgescu RE, Alexov EG, Gunner MR (2002) Combining conformational flexibility and continuum electrostatics for calculating pKas in proteins. *Biophys J* 83:1731–1748.

34. Grey MJ,f Tang Y, Alexov E, McKnight CJ, Raleigh DP, Palmer AG III (2006) Characterizing a partially folded intermediate of the villin headpiece domain under non-denaturing conditions: Contribution of His41 to the pH-dependent stability of the N-terminal subdomain. *J Mol Biol* 355: 1078–1094.

35. Harvey SC, Hoekstra P (1972) Dielectric relaxation spectra of water adsorbed on lysozyme. *J Phys Chem* 76:2987–2994.

36. Hunte C, Screpanti E, Venturi M, Rimon A, Padan E, Michel H (2005) Structure of a $Na^+/H^+$ antiporter and insights into mechanism of action and regulation by pH. *Nature* 435:1197–1202.

37. Im W, Lee MS, Brooks CL III (2003) Generalized Born model with a simple smoothing function. *J Comput Chem* 24:1691–1702.

38. Im W, Chen J, Brooks CL III (2006) Peptide and protein folding and conformational equilibria: Theoretical treatment of electrostatics and hydrogen bonding with implicit solvent models. *Adv Protein Chem* 72:173–198.

39. Kelly JW (1996) Alternative conformations of amyloidogenic proteins govern their behavior. *Curr Opin Struct Biol* 6:11–17.

40. Kelly JW (1998) The environmental dependency of protein folding best explains prion and amyloid diseases. *Proc Natl Acad Sci USA* 95:930–932.

41. Khandogin J, Brooks CL III (2006) Toward the accurate first-principles prediction of ionization equilibria in proteins. *Biochemistry* 45:9363–9373.

42. Khandogin J, Brooks CL III (2005) Constant pH molecular dynamics with proton tautomerism. *Biophys J* 89:141–157.

43. Khandogin J, Brooks CL III (2007) Linking folding with aggregation in Alzheimer's beta amyloid peptides. *Proc Natl Acad Sci USA* 104:16880–16885.

44. Khandogin J, Brooks CL III (2007) *Annual report of computational chemistry*, volume 3, chapter Molecular Simulations of pH-Mediated Biological Processes, Elsevier, Amsterdam, pp 3–11.

45. Khandogin J, Chen J, Brooks CL III (2006) Exploring atomistic details of pH-dependent peptide folding. *Proc Natl Acad Sci USA* 103:18546–18550.

46. Khandogin J, Raleigh DP, Brooks CL III (2007) Folding intermediate in the villin headpiece domain arises from disruption of a N-terminal hydrogen-bonded network. *J Am Chem Soc* 129:3056–3057.

47. Klamt A, Schüürmann G (1993) COSMO: A new approach to dielectric screening in solvents with explicit expressions for the screening energy and its gradient. *J Chem Soc Perkin Trans 2*, 2: 799–805.

48. Klamt A, Eckert F, Diedenhofen M, Beck ME (2003) First Principles Calculations of Aqueous $pK_a$ Values for Organic and Inorganic Acids Using COSMO-RS Reveal an Inconsistency in the Slope of the $pK_a$ Scale. *J Phys Chem A* 107:9380–9386.

49. Klapper I, Hagstrom R, Fine R, Sharp K, Honig B (1986) Focusing of Electric Fields in the Active Site of Cu-Zn Superoxide Dismutase: Effects of Ionic Strength and Amino-Acid Modification. *Proteins* 1:47–59.

50. Kong X, Brooks CL, III (1996) λ-dynamics: A new approach to free energy calculations. *J Chem Phys* 105:2414–2423.

51. Koumanov A, Karshikoff A, Friis EP, Borchert TV (2001) Conformational averaging in pK calculations: Improvement and limitations in prediction of ionization properties of proteins. *J Phys Chem B* 105:9339–9344.

52. Kuhn B, Kollman PA, Stahl M (2004) Prediction of $pK_a$ shifts in proteins using a combination of molecular mechanical and continuum solvent calculations. *J Comput Chem* 25:1865–1872.

53. Langsetmo K, Fuchs JA, Woodward C (1991) The conserved, buried aspartic acid in oxidized *Escherichia coli* thioredoxin has a $pK_a$ of 7.5. its titration produces a related shift in global stability? *Biochemistry* 30:7603–7609.

54. Langsetmo K, Fuchs JA, Woodward C, Sharp KA (1991) Linkage of thioredoxin stability to titration of ionizable groups with perturbed $pK_a$. *Biochemistry* 30:7609–7614.

55. Lee S, Chu ZT, Warshel A (1993) Microscopic and semimicroscopic calculations of electrostatic energies in proteins by the POLARIS and ENZYMIX programs. *J Comput Chem* 14: 161–185.

56. Lee K, Fitch CA, Lecomte JT, García-Moreno EB (2002) Electrostatic effects in highly charged proteins: Salt sensitivity of $pK_a$ values of histidines in staphylococcal nuclease. *Biochemistry* 41:5656–5667.

57. Lee S, Feig M, Salsbury FR Jr, Brooks CL III (2003) New analytic approximation to the standard molecular volume definition and its application to generalized Born calculations. *J Comput Chem* 24:1348–1356.

58. Lee S, Salsbury FR Jr, Brooks CL III (2004) Constant-pH molecular dynamics using continuous titration coordinates. *Proteins* 56:738–752.

59. Levy RM, Belhadj M, Kitchen DB (1991) Gaussian fluctuation formula for electrostatic free energy changes in solution. *J Chem Phys* 95:3627–3633.

60. Lim C, Bashford D, Karplus M (1991) Absolute p$K_a$ calculations with continuum dielectric methods. *J Phys Chem* 95:5610–5620.

61. Machuqueiro M, Baptista AM (2006) Constant-pH molecular dynamics with ionic strength effects: Protonation-conformation coupling in decalysine. *J Phys Chem B* 110:2927–2933.

62. Machuqueiro M, Baptista AM (2007). Stochastic titration study of hen egg lysozyme. *Proteins* 72:289–298.

63. Madura JD, Briggs JM, Wade RC, Davis ME, Luty BA, Ilin A, Antosiewicz J, Gilson MK, Bagheri B, Scott LR, McCammon JA (1995) Electrostatics and diffusion of molecules in solution: Simulations with the University of Houston Brownian Dynamics program. *Comput Phys Commun* 91:57–95.

64. Mehler EL, Eichele G (1984) Electrostatic effects in water-accessible regions of proteinst. *Biochemistry* 23:3887–3891.

65. Mehler EL, Guarnieri F (1999) A self-consistent, microenvironment modulated screened coulomb potential approximation to calculate pH-dependent electrostatic effects in proteins. *Biophys J* 75:3–22.

66. Mehler EL, Fuxreiter M, Simon I, Garcia B-Moreno (2002) The Role of hydrophobic microenvironments in modulating p$K_a$ shifts in proteins. *Proteins* 48:283–292.

67. Mertz JE, Pettitt BM (1994) Molecular dynamics at a constant pH. *Int J Supercomput Appl High Perform Comput* 8:47–53.

68. Merz KM Jr (1991) Determination of p$K_a$'s of ionizable groups in proteins: The p$K_a$ of Glu 7 and 35 in hen egg white lysozyme and Glu 106 in human carbonic anhydrase II. *J Am Chem Soc* 113:3572–2575.

69. Mok KH, Kuhn LT, Goez M, Day IJ, Lin JC, Andersen NH, Hore PJ (2007) A pre-existing hydrophobic collapse in the unfolded state of an ultrafast folding protein. *Nature* 447:106–109.

70. Mongan J, Case DA (2005) Biomolecular simulations at constant pH. *Curr Opin Struct Biol* 15:157–163.

71. Mongan J, Case DA, McCammon JA (2004) Constant pH molecular dynamics in generalized Born implicit solvent. *J Comput Chem* 25:2038–2048.

72. Nicholls A, Honig B (1991) A rapid finite difference algorithm, utilizing successive over-relaxation to solve the Poisson-Boltzmann equation. *J Comput Chem* 12:435–445.

73. Nozaki Y, Tanford C (1967) Examination of titration behavior. *Methods Enzymol* 11:715–734.

74. Nymeyer H, Gnanakaran S, García AE (2004) Atomic simulations of protein folding using the replica exchange algorithm. *Methods Enzymol* 383:119–149.

75. Oliveberg M, Arcus VL, Fersht AR (1995) p$K_a$ values of carboxyl groups in the native and denatured states of barnase: The p$K_a$ values of the denatured state are on average 0.4 units lower than those of model compounds. *Biochemistry* 34:9424–9433.

76. Osterhout JJ Jr, Baldwin RL, York EJ, Stewart JM, Dyson HJ, Wright PE (1989) [1]H NMR studies of the solution conformations of an analogue of the C-peptide of ribonuclease A. *Biochemistry* 28:7059–7064.

77. Parr R, Yang W (1989) *Density-Functional Theory of Atoms and Molecules*. Oxford University Press, New York.

78. Patel S, Brooks CL, III (2003) CHARMM fluctuating charge force field for proteins: I parameterization and application to bulk organic liquid simulations. *J Comput Chem* 25:1–16.

79. Patel S, Mackerell AD Jr, Brooks CL III (2004) CHARMM fluctuating charge force field for proteins: II Protein/solvent properties from molecular dynamics simulations using a nonadditive electrostatic model. *J Comput Chem* 25:1504–1514.

80. Riccardi D, Schaefer P, Cui Q (2005) pKa Calculations in Solution and Proteins with QM/MM Free Energy Perturbation Simulations: A Quantitative Test of QM/MM Protocols. *J Phys Chem B* 109:17715–17733.

81. Richardson WH, Peng C, Bashford D, Noodleman L, Case DA (1997) Incorporating Solvation Effects into Density Functional Theory: Calculation of Absolute Acidities. *Int J Quantum Chem* 61:207–217.

82. Rick SW, Lynch DL, Doll JD (1991) A variational Monte Carlo study of argon, neon, and helium clusters. *J Chem Phys* 95:3506–3520.

83. Ripoll DR, Vorobjev YN, Liwo A, Vila JA, Scheraga HA (1996) Coupling between folding and ionization equilibria: Effects of ph on the conformational preferences of polypeptides. *J Mol Biol* 264:770–783.

84. Schaefer M, Karplus M (1996) A Comprehensive Analytical Treatment of Continuum Electrostatics. *J Phys Chem* 100(5):1578–1600.

85. Schaefer M, van Vlijmen HWT, Karplus M (1998) Electrostatic contributions to molecular free energies in solution. *Adv Protein Chem* 51:1–57.

86. Schäfer A, Klamt A, Sattel D, Lohrenz JCW, Eckert F (2000) Cosmo implementation in turbomole: Extension of an efficient quantum chemical code towards liquid systems. *Phys Chem Chem Phys* 2:2187–2193.

87. Schutz CN, Warshel A (2001) What are the dielectric constants of proteins and how to validate electrostatic models? *Proteins* 44:400–417.

88. Sham YY, Chu ZT, Warshel A (1997) Consistent calculations of p$K_a$s of ionizable residues in proteins: Semi-microscopic and microscopic approaches. *J Phys Chem B* 101:4458–4472.

89. Shoemaker KR, Fairman R, Schultz DA, Robertson AD, York EJ, Stewart JM, Baldwin RL (1990) Side-chain interactions in the C-peptide helix: Phe8-His12$^+$. *Biopolymers* 29:1–11.

90. Simonson T, Brooks CL III (1996) Charge screening and the dielectric constant of proteins: Insights from molecular dynamics. *J Am Chem Soc* 118:8452–8458.

91. Simonson T, Carlsson J, Case DA (2004) Proton binding to proteins: p$K_a$ calculations with explicit and implicit solvent models. *J Am Chem Soc* 126:4167–4180.

92. Srinivasan J, Trevathan MW, Beroza P, Case DA (1999) Application of a pairwise generalized Born model to proteins and nucleic acids: Inclusion of salt effects. *Theor Chem Acc* 101:426–434.

93. Stern HA (2007) Molecular simulation with variable protonation states at constant pH. *J Chem Phys* 126:164112.

94. Still WC, Tempczyk A, Hawley RC, Hendrickson T (1990) Semianalytical treatment of salvation for molecular mechanics and dynamics. *J Am Chem Soc* 112:6127–6129.

95. Sugita Y, Okamoto Y (1999) Replica-exchange molecular dynamics method for protein folding. *Chem Phys Lett* 314:141–151.

96. Tanford C, Kirkwood JG (1957) Theory of protein titration curves. I. general equations for impenetrable spheres. *J Am Chem Soc* 79:5333–5339.

97. Tanokura M (1983) $^1$H-NMstudy R on the tautomerism of the imidazole ring of histidine residues: I. Microscopic pK values and molar ratios of tautomers in histidine-containing peptides. *Biochim Biophys Acta* 742:576–585.

98. Tironi G, Sperb R, Smith PE, van Gunsteren WF (1995) A generalized reaction field method for molecular dynamics simulations. *J Chem Phys* 102:5451–5459.

99. van Vlijmen HW, Schaefer M, Karplus M (1998) Improving the accuracy of protein p$K_a$ calculations: Conformational averaging versus the average structure. *Proteins* 33:145–158.

100. Vila JA, Ripoll DR, Scheraga HA (2001) Influence of lysine content and pH on the stability of alanine-based copolypeptides. *Biopolymers* 58:235–246.

101. Von Arnim M, Ahlrichs R (1998) Performance of Parallel TURBOMOLE for Density Functional Calculations. *J Comput Chem* 19(15):1746–1757.
102. Walczak M, Antosiewicz JM (2002) Langevin dynamics of proteins at constant pH. *Phys Rev E* 66:051911.
103. Warshel A, Sussman F, King G (1986) Free energy of charges in solvated proteins: Microscopic calculations using a reversible charging process. *Biochemistry* 25:8368–8372.
104. Warshel A, Sharmaa PK, Katoa M, Parson WW (2006) Modeling electrostatic effects in proteins. *Biochim Biophys Acta* 1764:1647–1676.
105. Weiss N (1997) The Hill equation revisited: Uses and misuses. *FASEB J* 11:835–841.
106. Yang A-S, Gunner MR, Sampogna R, Sharp K, Honig B (1993) On the calculation of p$K_a$'s in proteins. *Proteins* 15:252–265.
107. You TJ, Bashford D (1995) Conformation and hydrogen ion titration of proteins: A continuum electrostatic model with conformational flexibility. *Biophys J* 69:1721–1733.

# CHAPTER 11

# QUANTUM MECHANICAL STUDIES OF THE PHOTOPHYSICS OF DNA AND RNA BASES

KURT A. KISTLER AND SPIRIDOULA MATSIKA

*Department of Chemistry, Temple University 1901 N.13th Street, Philadelphia, PA, 19122, USA,*
*e-mail: smatsika@temple.edu*

**Abstract:**     This chapter reviews current computational studies of the photophysical behavior of DNA and RNA bases. The theoretical concepts and electronic structure methods appropriate for computational photophysical and photochemical studies are presented. The natural nucleobases have ultrashort excited state lifetimes and very low quantum yields for fluorescence. Quantum chemical calculations revealing radiationless decay pathways that quench fluorescence are discussed

**Keywords:**     Fluorescence quenching, Nucleobase, Photophysics, Photochemistry, Quantum Chemistry, Computational Photochemistry, Conical intersections, Radiationless decay, Ab initio

## 11.1.     INTRODUCTION

In this chapter we will review current computational studies of the photophysical behavior of DNA and RNA bases. The photophysics of nucleic acids is of extreme importance since their interaction with UV radiation can lead to photodamage in DNA [1, 2]. During the last decade there has been great progress in understanding the fundamentals of the photophysical behavior of nucleic acids based on femtosecond spectroscopy and high level quantum chemical studies. Earlier work has been reviewed several times [3–5] and the more recent studies were reviewed in great detail in 2004 [6]. The current chapter focuses primarily on recent quantum chemical calculations on nucleobases and their radiationless relaxation through singlet electronic states. The review is mainly divided into two parts. Section 11.2 reviews the fundamental theoretical tools for studying photophysical and photochemical processes in molecular systems while Section 11.3 presents the computational studies on photophysical behavior of nucleobases.

## 11.2.    COMPUTATIONAL PHOTOCHEMISTRY

The interaction of light with matter has many applications in chemistry and biology. Biological processes using light are essential in life, and vital processes, such as photosynthesis, light harvesting, vision, photochemical damage and repair in DNA, utilize light to proceed. These biological functions are based on fundamental photophysical and photochemical processes in molecules present in biological systems. Understanding better these processes requires extensive studies by physical and biological chemists.

The contribution of computational chemistry to the area of photochemistry and photobiology has been slow over the years because of the difficulty of the theory involved, but progress has been accelerated in recent years. The main difficulty arises from the fact that photochemistry and photobiology involve processes taking place on the excited states of molecules. Excited states only exist in a quantum world, and consequently all the methods based on classical mechanics that have been developed and are being used for ground state processes in biochemistry cannot be used, at least not without some incorporation of quantum effects. So a quantum mechanical model is needed even for the simplest qualitative description of photochemical events. When localized events are described a quantum mechanical model can be used for the local excitation in combination with classical mechanics to incorporate the much larger biological environment. The situation is even more challenging since the requirements to calculate excited states, even within quantum mechanical models, are a lot more demanding than the requirements for ground state reactions. So although quantum mechanical methods for ground state problems have been available for many years and are now being used routinely by many chemists, even non-experts, the same is not true for excited states. Accurate methods exist for small systems but the accuracy deteriorates rapidly as the size of the system increases. This area of theoretical chemistry is an active field of research today, and methods are being developed and improved so that we can apply them and expand our knowledge of photochemical and photobiological events.

Computational photochemistry requires knowledge of the excited states beyond initial vertical excitations and calculations of the potential energy surfaces (PES), energies as a function of the geometrical distortions, are essential. Thus, the requirements of electronic structure methods are more than just an accurate description of excited states at a given equilibrium geometry of a molecule, but, rather, an accurate description at different distorted geometries, and possibly the ability to calculate the forces along the PES. When we move away from the initial absorption to reaction mechanisms we need to consider the framework that enables modeling of these excited state processes. Figure 11-1 shows a cartoon of ground and excited state PES and possible processes that can occur upon photoexcitation. Ground state reactions occur through transition states and computational chemists study these reactions by locating transition states between the reactants and products and connecting them through reaction coordinates. Excited states also have transition states, but ultimately a system cannot remain on an excited state with the extra electronic energy, rather,

*Figure 11-1.* Cartoon of ground and excited state potential energy surfaces, indicating points where nona-diabatic transitions can occur

it will return to the ground state, either at the equilibrium structure where absorption occurred, or it will go to a different product. In the former case the overall process is photophysical while in the latter it is photochemical. The important point here is that any photophysical or photochemical process involves transitions between different electronic states. Transitions between different electronic states can occur either with absorption or emission of photons (radiative transitions) or radiationlessly.

In order to computationally study these processes it is crucial to be able to model the radiationless nonadiabatic transitions. It is growth in this area that has catalyzed computational photochemistry and photobiology. The basics of the theory describing radiationless transitions are outlined below.

## 11.2.1.  Nonadiabatic Processes

Based on the Born-Oppenheimer approximation [7] the behavior of molecules is described by the dynamics of the nuclei moving along a single PES generated by the electrons. Nonadiabatic phenomena occur when at least two potential energy surfaces approach each other, and the coupling between them becomes important. As two PES approach each other, the rate of nonadiabatic processes depends on the energy separating these two surfaces. The Landau-Zener formula gives the rate for nonadiabatic transitions $P_{LZ}$ between two states

$$P_{LZ} = e^{-\frac{2\pi\lambda^2}{v\Delta F}}, \tag{11-1}$$

where $\lambda$ is the coupling, $v$ is the velocity at the crossing point, and $\Delta F$ is the difference of the slopes of the two diabatic surfaces at that point. This formula is a one-dimensional scattering solution to the problem derived by assuming the the nuclei follow a classical trajectory [8–10]. Traditionally nonadiabatic transitions are described as internal conversion through avoided crossings. Fermi's Golden Rule [11] can also be used to calculate nonadiabatic transition rates.

In recent years ultrafast experimental techniques have allowed the observation of nonadiabatic processes that take place in femtoseconds [12], and these ultrafast rates cannot be explained with the traditional theories. Conical intersections, which are the actual crossings of two potential energy surfaces, however, can facilitate rapid nonadiabatic transitions. Conical intersections were known mathematically since the 1930s [13, 14] but they were regarded as mathematical curiosities rather than a useful concept for explaining photochemistry. Developments in the 1990s however changed that view when algorithms were developed that allowed for the location of conical intersections without the presence of symmetry [15, 16]. These algorithms have since revealed that conical intersections occur in the excited states of many molecules and are far from uncommon [10, 17–25]. In fact, our understanding of photochemistry has undergone a transformation based on the concept of conical intersections, which has become the standard in photochemistry textbooks [26, 27]. Conical intersections and nonadiabatic dynamics are expected to participate actively in photobiology as studies reveal [28–35]. Modern computational photophysics and photochemistry uses high level electronic structure methods to locate conical intersections and map the PES of ground and excited states in order to understand the underlying mechanisms involved.

### 11.2.1.1. Conical Intersections

Neumann and Wigner proved, in their seminal work in 1929 [13], that for a molecular system with $N^{int}$ internal nuclear coordinates two electronic surfaces become degenerate in a subspace of dimension $N^{int} - 2$. To illustrate this dimensionality rule, consider two intersecting adiabatic electronic states, $\psi_1$ and $\psi_2$. These two states are expanded in terms of two diabatic states $\phi_1$ and $\phi_2$, which are orthogonal to all the remaining electronic states and to each other [36]. The electronic energies are the eigenvalues of the Hamiltonian matrix

$$\mathbf{H} = \begin{pmatrix} H_{11} & H_{12} \\ H_{21} & H_{22} \end{pmatrix} \tag{11-2}$$

where $H_{ij} = \langle \phi_i | H | \phi_j \rangle$. Diagonalizing $\mathbf{H}$ gives the adiabatic energies, $E_{1,2} = \frac{H_{11}+H_{22}}{2} \pm \frac{1}{2}\sqrt{(H_{11} - H_{22})^2 + 4H_{12}^2}$. For the eigenvalues to be degenerate two conditions must be satisfied,

$$H_{11} - H_{22} = 0 \tag{11-3}$$

$$H_{12} = 0 \tag{11-4}$$

Thus the degeneracy occurs in the subspace where these two equations are satisfied, which is $N^{int} - 2$.

A consequence of the noncrossing rule is that the degeneracy is lifted linearly in a two-dimensional space, called the branching plane, and the energies of the two states form a double cone. Conical intersections are characterized by the topography of the potential energy surfaces in their vicinity [20, 36]. This topography plays a significant role in the efficacy of a conical intersection's ability to promote a nonadiabatic transition [20, 30, 32, 36–39].

Location of conical intersections is based on minimizations constrained to satisfy the above Eqs. (11-3) and (11-4). The development of analytic gradient techniques for ab initio methods [40] enables the efficient characterization of PES by finding optimized structures for molecules, locating transition states, and establishing reaction pathways. Optimizations in the excited states and conical intersections are less widespread than those of the ground state but progress is being made, as discussed in the next section. To locate conical intersections efficiently the nonadiabatic couplings are needed in addition to the gradients of the surfaces (because of the second constraint), although algorithms that do not require the coupling can be used. Even though conical intersections are multidimensional seams (of dimension $N^{int} - 2$), usually the minimum energy point on that seam is located and reported.

The derivative (nonadiabatic) coupling, $\mathbf{f}_{IJ}$, is the term neglected in the Born-Oppenheimer approximation that is responsible for nonadiabatic transitions between different states $I$ and $J$. It originates from the nuclear kinetic energy operator operating on the electronic wavefunctions $\psi_I$ and is given by

$$\mathbf{f}_{IJ} = \langle \psi_I | \nabla \psi_J \rangle \tag{11-5}$$

where differentiation, $\nabla$, is with respect to nuclear coordinates and the integration is over electronic coordinates. The derivative coupling $\mathbf{f}_{IJ}$ is a measure of the variation of the electronic wavefunction with nuclear coordinates, and depends on the energy difference between states $I$ and $J$. When the states are well separated, the coupling is small and the Born-Oppenheimer approximation is valid. If, however, the electronic eigenvalues are close, a small change in the nuclear coordinates may cause a large change in the electronic wavefunctions, a situation where the coupling becomes important.

In order to be able to characterize the PES of excited states, locate conical intersections, and derive mechanisms for photophysics and photochemistry, efficient electronic structure methods for excited states are required. In the following section we give a brief overview of the current state of methodological developments in electronic structure methods applicable to excited states.

### 11.2.2. Electronic Structure Methods for Excited States

The electronic structure method used to provide the energies and gradients of the states is crucial in photochemistry and photophysics. Ab initio electronic structure methods have been used for many years. Treating closed shell systems in their ground state is a problem that, in many cases, can now be solved routinely by chemists using standardized methods and computer packages. In order to obtain quantitative results, electron correlation (also referred to as dynamical correlation) should be included in the model and there are many methods available for doing this based on either variational or perturbation principles [41].

However, when treating excited states the situation is more complicated. The simplest method that can be used to study excited states is configuration interaction with single excitations (CIS), which is equivalent to Hartree-Fock (HF) for the ground state. In this method all Slater determinants with promotion of a single electron from the HF wavefunction to the virtual orbitals are constructed and the optimized linear combinations of those are obtained by diagonalizing the Hamiltonian matrix in the basis of the Slater determinants. This has been the most frequently used method in the past to study excited states and it still is the only applicable method for very large systems. Because it does not include dynamical correlation it predicts excitation energies that are too high. Often the energies are scaled by an empirical factor so that they match the experimental absorption maxima. It has also been combined with semiempirical methods giving the opportunity to be used to even larger systems.

In order to improve the quality of excited state description, correlation needs to be included. The methods having been developed can be categorized into two groups, the multireference methods, and the single-reference based methods.

### 11.2.2.1. *Multireference Methods*

Multireference methods provide a straightforward way to treat excited states, since studying excited states requires the equivalent treatment of these states. Multireference methods are extensions of the single reference Hartree-Fock or configuration interaction (CI) methods, where many configurations are used instead of a single configuration,

$$\psi_I = \sum_{a=1}^{N^{CSF}} c_a^I \psi_a \tag{11-6}$$

The basis of the expansion, $\psi_a$, are configuration state functions (CSF), which are linear combinations of Slater determinants that are eigenfunctions of the spin operator and have the correct spatial symmetry and total spin of the electronic state under investigation [42].

The Multi-Configuration Self-Consistent Field (MCSCF) method includes configurations created by excitations of electrons within an active space. Both the coefficients $c_a$ of the expansion in terms of CSFs and the expansion coefficients of the

molecular orbitals setting up the Slater determinants are optimized. The choice of which configurations and active orbitals to be included is critical and depends on the chemical nature of the problem. This selection process is the most difficult step in setting up an MCSCF calculation, and the flexibility in the choice of the active space and its dependence on the investigator's intuition have been criticisms of these methods. The most useful approach is the complete active space MCSCF designated as CASSCF [43, 44]. Nevertheless, one must still choose the orbitals to be included in the active space. For small systems a full valence active space can be used, where all the valence orbitals of all atoms are included. This is not possible for larger systems, and compromises between rigor and computing time have to be made.

To include dynamical correlation into a quantum calculation one must go beyond the MCSCF approach. A multireference configuration interaction model (MRCI) is a CI expansion where many electronic configurations are used as references instead of using a single Hartree-Fock reference. The final expansion is a linear combination of all the references and of the configurations generated from single and double excitations out of these references to the virtual orbitals. The MRCI method is very accurate provided all the important configurations are included in the expansion. This requirement can be satisfied for small systems, but as the size of the system increases the expansion becomes prohibitively large and truncations are necessary. If one is interested in excitation energies at a single geometry (most often vertical excitations) different expansions for the different states can be used to reduce the size of the calculation. Alternatively the configurations can be truncated based on some selection criterion. If, however, one is interested in the excited states over a range of geometries, with the possibility of states crossing, then the importance of consistency of the expansion for all states dictates the approach. Analytic gradients have been developed for MRCI wavefunctions [45–48], which is a huge advantage in computing cost when using these wavefunctions for studying conical intersections and excited state stationary points. The COLUMBUS suite of programs [49] for example has algorithms for studying conical intersections and derivative couplings [50, 51] that rely upon analytic gradient techniques [45–47]. An efficient, internally contracted MRCI method has been developed by Werner and Knowles [52] and is implemented in the MOLPRO suite of programs [53]. This program is widely used, but it has the disadvantage of not including analytic gradients.

Efficient methods to calculate excited states based on perturbation theory have also been developed. Roos and coworkers [54, 55] developed the Complete Active Space Second-Order Perturbation Theory (CASPT2), which has been implemented in the ab initio package MOLCAS [56]. The CASPT2 method [54, 55] perturbatively computes, through second order, the dynamical correlation using a single CASSCF reference state and a non-diagonal zeroth-order Hamiltonian $H_0$. This method has been used for the study of many systems of various sizes and it reproduces experimental excitation energies with high accuracy [57]. The CASPT2 method, however, cannot treat near degeneracies efficiently, since in these cases the CASSCF wavefunction is an insufficient reference state for the perturbation calculation. Multistate perturbative methods have been developed [58, 59] (MS-CASPT2) that avoid

this problem, and have been found to perform well at avoided crossings and when valence-Rydberg mixing occurs [59]. Serrano-Andrés et al. [60] recently investigated the possibility of using CASPT2 and multi-state CASPT2 for locating actual conical intersections [60]. They concluded that these methods can lead to nonphysical results when small active spaces are used. CASPT2 is frequently used to obtain refined energies at selected points (stationary points or conical intersections) optimized at the CASSCF level. This approach has been very useful in studying nonadiabatic problems in systems of moderate size, and, as will be seen in Section 11.3, it is the most frequent protocol used in the studies of radiationless decay in nucleobases.

Other implementations of perturbation theory for excited states have also been developed [58, 61, 62] and exist in other computational packages such as in GAMESS [63]. MCQDPT [64, 65] is a method developed by Nakano that includes perturbation theory for simultaneous treatment of many states. This method calculates perturbation based corrections to both diagonal and off-diagonal elements to give an effective Hamiltonian which is then diagonalized. Diagonalization after inclusion of the off-diagonal perturbation ensures that avoided crossings of states of the same symmetry are treated correctly.

### 11.2.2.2.     *Single-Reference Methods*

A different approach for calculating excited states is based on indirect methods that allow one to calculate excitation energies based on a single reference wavefunction. Single reference methods for the calcualtion of excited states of large molecules have been reviewed recently [66].

Configuration interaction with single substitutions (CIS) is the simplest method that can describe electronically-excited states, but it does not include dynamical correlation and thus it cannot predict accurate energies. Furthermore it is unable to predict the correct topography of the intersecting surfaces because of the identically zero Hamiltonian matrix elements between the HF reference and the singly-excited determinants (Brillouin's theorem). One has to go beyond the CIS approximation and to include the dynamic correlation consistently in order to recover the correct conical topography. The quasidegenerate second-order perturbation corrections to CIS developed by Head-Gordon and coworkers, denoted CIS(D) [67], where the (D) signifies double excitations treated perturbatively, is the cheapest excited-state method that includes dynamical correlation. This method has the accuracy of the second-order Møller-Plesset perturbation theory (MP2) method for the ground state and scales as $N^5$ where $N$ is the number of basis functions.

The CIS(D) method should not be able to resolve the problem of describing conical intersections because it still separates the calculation of the ground and the excited state energies. A new electronic structure model, termed CIS(2), has been developed in which the energies of both the ground and singly-excited states are eigenvalues of a dressed symmetric CI matrix in the space of the reference and singly-excited determinants. The effects of double and triple substitutions are approximately included into the CI matrix elements in the spirit of quasidegenerate second-order perturbation theory [68]. This method can describe conical intersections.

Starting from a coupled cluster representation for the ground state, linear-response and equation-of-motion coupled cluster (LR-CC and EOM-CC respectively) [69, 70] can be used to provide accurate excitation energies. Variations of the method can allow for more extended problems such as bond-breaking [71, 72]. Coupled cluster methods represent the most sophisticated methods to account for dynamical correlation when a single electronic configuration is a good first order description of the chemical system but they are computationally expensive and thus limited to small systems. The symmetry-adapted cluster configuration interactions approach (SAC-CI) is a related approach that can be used for applications to open-shell systems, as well as closed-shell systems [73].

The CC2 method [74] is an approximation to coupled cluster with singles and doubles (CCSD), and the excited state energies calculated have MP2 quality. An implementation that employs the resolution of identity (RI) approximation for two-electron integrals to reduce the CPU time is also available, RI-CC2 [75], which is suitable for large scale integral-direct calculations. This method has been implemented in TURBOMOLE [76].

Alternatively, methods based on the density rather than the wavefunction [77], have proven quite succesful for the ground state and have been extended to an excited state formalism. Time-dependent density functional theory (TDDFT) [78] provides excitation energies at a relatively low cost. TDDFT has gained unprecedented popularity in recent years. This method can predict vertical excitation energies efficiently, but its extension to the description of excited state properties and potential energy surfaces is more complicated and is currently under development. Methods based on TDDFT appear to be inadequate for describing conical intersections [79] due to a qualitatively incorrect description of the energy surfaces in the vicinity of a conical intersection. A method that combines density functional theory (DFT) with MRCI (DFT/MRCI) has also been developed and has been useful in describing medium size systems [80].

## 11.3.    EXCITED STATES IN DNA

The interaction of UV radiation with nucleic acids is of great importance since it can lead to UV-induced damage in DNA with profound consequences, including photocarcinogenesis [1, 2]. The nucleobases are the primary chromophores in DNA and RNA, and consequently, the photophysical and photochemical behavior of the nucleobases has been the focus of extensive theoretical and experimental work over the years [4, 6, 81, 82].

Upon absorption of UV radiation from sunlight the bases can proceed through photochemical reactions that can lead to photodamage in the nucleic acids. Photochemical reactions do occur in the bases, with thymidine dimerization being a primary result, but at low rates. The bases are quite stable to photochemical damage, having efficient ways to dissipate the harmful electronic energy, as indicated by their ultrashort excited state lifetimes. It had been known for years that the excited states were short lived, and that fluorescence quantum yields are very low for all bases [4, 81, 82]. Femtosecond laser spectroscopy has, in recent years, enabled a much

more accurate measurement of lifetimes and detailed studies of excited state dynamics [6, 83–120]. Transient absorption spectroscopy and fluorescence up-conversion techniques have been used to measure the lifetimes and to study the dynamics of nucleosides, nucleotides, and isolated bases, in solution. These studies report lifetimes on the order of hundreds of femtoseconds and suggest that nonradiative relaxation proceeds on an ultrafast timescale to the ground state with the excess energy being transformed into heat [83, 84]. Experiments in the gas phase have been performed as well, reporting spectra and lifetimes [95–115]. A review was published in 2004 summarizing the experimental and theoretical work up to that point of the excited states in DNA and RNA [6]. Much progress has occurred since then in this exciting field. Experiments have moved beyond single bases, nucleotides and nucleosides, to the much more complicated and more biological relevant area of polynucleotides, and single and double stranded DNA. Theoretically, there had been very limited studies of excited state photophysics in the monomers in 2004, but there has been an explosion since then. Many studies have appeared since then on mechanisms for radiationless relaxation in all nucleobases. Computationally, only the bare bases have been studied. It is not clear what effect the phosphate or sugar has on the excited state dynamics. Initially the effect was believed to be very small, but more recent experimental results suggest that the presence of phosphates in nucleotides may play a role in the lifetimes [116], although that has not been explicitly studied theoretically. Recently theoretical attention has also shifted to systems with more than one base, following the lead of experimental work. Computationally, however, it is much more difficult to study these bigger systems, so progress is slower. Here we will focus mainly on the computational studies of radiationless relaxation in monomers, starting by the character and energies of excited states upon initial absorption.

### 11.3.1.  Vertical Excitation Energies

Figure 11-2 shows the structures of the five nucleobases considered here. Only the tautomers present in the nuclei acids will be considered. The first step in understanding the photochemistry and photophysics of nucleic acid bases is to have a good knowledge of their excited states produced by photon absorption. Here only singlet excited states are considered. Most computational studies have focused on the singlet states, since any mechanism that involves singlet states will be much faster than mechanisms involving triplet states. The reader should not infer from the absence of triplet states in this review that they are not contributing to any radiationless processes in the nucleobases. The valence orbitals likely to be excited in these heteroatomic aromatic molecules are $\pi$ orbitals and lone pair orbitals localized on nitrogen or oxygen. Figures 11-3 and 11-4 show the orbitals calculated to be involved in excitation for the pyrimidine and purine bases, respectively. The lowest excited states then are either $\pi\pi^*$ excitations or $n\pi^*$. At higher energies Rydberg states also exist, although these are more difficult to detect experimentally because of their low oscillator strengths. They are also more challenging to compute in a balanced way

*Figure 11-2.* Structures of the DNA/RNA bases



*Figure 11-3.* Molecular orbitals for the pyrimidine bases cytosine and uracil

relative to the valence states due to their different requirements for basis sets and dynamical correlation, so the exact location of these states is less certain.

The ground state minimum of the nucleobases is planar or almost planar. In the bases that have an amino group, adenine, cytosine and guanine, the nitrogen on the animo group is pyramidalized breaking the planarity slightly. So the point group

*Figure 11-4.* Molecular orbitals for the purine bases adenine and guanine

symmetry of the bases is $C_s$, or very close to it, and the excited states have either
$A'$ symmetry (i.e. $\pi\pi^*$ states) or $A''$ symmetry (i.e. $n\pi^*$, $\pi\sigma^*$). Vertical excitation
energies for uracil, thymine, cytosine, adenine and guanine calculated using var-
ious methods are given in Tables 11-1, 11-2, 11-3, 11-4, and 11-5. Only excitation
energies for the canonical tautomers of the nucleobases present in DNA are shown in
these tables and discussed further. This is a representative list of calculated energies
and does not include all published values. The most recent results obtained using
high level ab initio methods including dynamical correlation are included but older
values based on semiempirical methods and scaled CIS results are omitted. The ta-
bles list the methods used, but there is no specific information about the basis sets
used or other related information. The reader is directed to the actual publication for
the computational details.

### 11.3.1.1.   Uracil–Thymine

Uracil has eight $\pi$ orbitals and two lone pair $n_O$ orbitals on the two oxygen atoms.
The lowest excited states, if we neglect Rydberg states, originate from excitations
from these valence orbitals. Figure 11-3 shows the orbitals mainly participating in the
first four excited states. The excited states of uracil have been studied theoretically

*Table 11-1.* Vertical excitation energies in eV relative to the ground state minimum of the singlet electronic excited states in the uracil

|  | $S_1$ | $S_2$ | $S_3$ | $S_4$ | $S_5$ |
|---|---|---|---|---|---|
| CASSCF [125] | 4.78[a] | 6.31[a] | 6.88[b] | 7.03[b] |  |
| CASSCF [126] | 5.11[a] | 7.36[b] |  |  |  |
| CAS-MSPT2 [126] | 4.83[a] | 5.20[b] |  |  |  |
| CASPT2 [125] | 4.54 (0.00018)[a] | 5.00 (0.19)[b] | 5.82 (0.08)[b] | 6.00 (0.00)[a] |  |
| MRCI [147] | 4.80 (0.00015)[a] | 5.79 (0.19)[b] | 6.31 (0.000)[a] | 6.57 (0.035)[b] |  |
| CCSD [128] | 5.428[a] | 6.015[b] |  |  |  |
| QDPT2 [128] | 4.988[a] | 5.883[b] |  |  |  |
| QCCSD [128] | 4.988[a] |  |  |  |  |
| DFT/MRCI [130] | 4.61 (0.0002)[a] | 5.44 (0.2626)[b] | 5.95 (0.000)[a] | 6.15 (0.0501)[b] |  |
| CC2 [124] | 4.79 (0.000)[a] | 5.34 (0.182)[b] | 6.02 (0.003)[c] | 6.09 (0.000)[a] | 6.25 (0.032)[b] |
| CIS(2) [68] | 5.107 | 5.891 |  |  |  |
| TDDFT [129] | 4.64[a] | 5.11[b] | 5.64[a] | 5.85[b] |  |
| TDDFT [131] | 4.66 (0.0000)[a] | 5.17 (0.1343)[b] | 5.67 (0.0021)[c] | 5.79 (0.0000)[a] | 5.89 (0.0366)[b] |

Oscillator strengths are given in parenthesis.
[a] $n\pi^*$ state.
[b] $\pi\pi^*$ state.
[c] Rydberg state.

*Table 11-2.* Vertical excitation energies in eV relative to the ground state minimum of the singlet electronic excited states in thymine

| CASSCF [125] | 5.22[a] | 6.75[b] | 6.77[a] | 7.15[b] |  |
|---|---|---|---|---|---|
| CASSCF [126] | 5.30[a] | 7.49[b] |  |  |  |
| CAS-MSPT2 [126] | 4.88[a] | 5.17[b] |  |  |  |
| CASPT2 [125] | 4.39 (0.00019)[a] | 4.88 (0.17)[b] | 5.88 (0.17)[b] | 5.91 (0.00091)[a] |  |
| CASPT2 [152] | 4.81[b] | 4.97[a] | 5.99[b] | 6.51[a] |  |
| CC2 [152] | 4.88[a] | 5.29[b] | 6.39[b] | 6.25[a] |  |
| CC2 [124] | 4.82 (0.000)[a] | 5.20 (0.182)[b] | 5.74 (0.000)[c] | 6.16 (0.000)[a] | 6.27 (0.037)[b] |
| TDDFT [131] | 4.69 (0.0000)[a] | 4.96 (0.1388)[b] | 5.43 (0.0003)[c] | 5.81 (0.0001)[a] | 5.95 (0.0668)[b] |

Oscillator strengths are given in parenthesis.
[a] $n\pi^*$ state.
[b] $\pi\pi^*$ state.
[c] Rydberg state.

previously with ab initio methods ranging from CIS [121–123] to highly correlated CC2 [124], CASPT2 methods [125, 126], MRCI methods [127], coupled cluster [128], and density functional approaches for excited states [129–131]. The lowest singlet excited states using different methods are shown in Table 11-1. Although the different methods give different vertical excitation energies, all of them agree that the first excited state $S_1$ is an $n\pi^*$ and the second is a $\pi\pi^*$ more than 0.4 eV above. The next two states, $S_3$ and $S_4$, are $n\pi^*$ and $\pi\pi^*$ respectively, but their predicted ordering may change depending on the method used. The correlated methods MRCI,

*Table 11-3.* Vertical excitation energies in eV relative to the ground state minimum of the singlet electronic excited states in cytosine

|                      | $S_1$          | $S_2$            | $S_3$            | $S_4$            | $S_5$           |
| -------------------- | -------------- | ---------------- | ---------------- | ---------------- | --------------- |
| CASSCF [153]         | 5.21[a]        | 5.24[b]          | 6.00             |                  |                 |
| MRCI [157]           | 5.14 (0.067)[a] | 5.29 (0.002)[b] | 5.93 (0.001)[b]  |                  |                 |
| CASPT2 [140]         | 4.50 (0.065)[a] | 4.88 (0.001)[b] | 5.23 (0.003)[b]  |                  |                 |
| CASPT2 [134]         | 4.39 (0.061)[a] | 5.00 (0.005)[b] | 6.53 (0.001)[b]  |                  |                 |
| CR-EOMCCSD(T) [263]  | 4.76[a]        | 5.24[b]          |                  |                  |                 |
| CC2 [124]            | 4.66 (0.052)[a] | 4.87 (0.002)[b] | 5.26 (0.002)[b]  | 5.53 (0.005)[c]  | 5.61 (0.138)[a] |
| CIS(2) [68]          | 5.026          | 5.330            |                  |                  |                 |
| DFT/MRCI [133]       | 4.83 (0.0803)[a] | 5.02 (0.0022)[b] | 5.50 (0.0014)[b] | 5.67 (0.1807)[a] | 5.91 (0.000)[b] |
| TDDFT [153]          | 4.71 (0.036)[a] | 4.76 (0.002)[b] | 5.15 (0.001)[b]  |                  |                 |
| TDDFT [131]          | 4.62 (0.0455)[a] | 4.70 (0.0033)[b] | 5.10 (0.0008)[b] | 5.26 (0.0051)[c] | 5.44 (0.0744)[a] |

Oscillator strengths are given in parenthesis.
[a] $\pi\pi^*$ state.
[b] $n\pi^*$ state.
[c] Rydberg state.

*Table 11-4.* Vertical excitation energies in eV relative to the ground state minimum of the singlet electronic excited states in adenine

|                | $S_1$            | $S_2$            | $S_3$            | $S_4$            | $S_5$           |
| -------------- | ---------------- | ---------------- | ---------------- | ---------------- | --------------- |
| CASSCF [138]   | 5.226[a]         | 6.193[b]         | 6.674[b]         | 6.859[a]         | 7.232[a]        |
| CASPT2 [138]   | 4.852 (0.0064)[a] | 4.902 (0.1416)[a] | 5.503 (0.0039)[b] | 5.685 (0.0061)[b] | 5.68[c] [160]  |
| CASPT2 [143]   | 5.13 (0.070)[a]  | 5.20 (0.370)[a]  | 6.15 (0.001)[b]  | 6.86 (0.001)[b]  |                 |
| CASPT2 [164]   | 4.96 (0.004)[b]  | 5.16 (0.004)[a]  | 5.35 (0.175)[a]  |                  |                 |
| CASPT2 [141]   | 5.01 (0.006)[b]  | 5.05 (0.002)[a]  | 5.16 (0.21)[a]   | 5.72 (0.005)[b]  |                 |
| MCQDPT [264]   | 4.92[a]          |                  |                  |                  |                 |
| MRCI [210]     | 5.67[a]          | 5.68[b]          | 6.11[a]          | 6.25[b]          |                 |
| CIPSI [219]    | 4.96 (0.001)[b]  | 4.97 (0.010)[a]  | 5.34 (0.359)[a]  |                  |                 |
| CC2 [124]      | 5.12 (0.007)[b]  | 5.25 (0.302)[ac] | 5.53 (0.011)[c]  | 5.75 (0.003)[bc] | 5.86 (0.004)[c] |
| DFT/MRCI [139] | 4.90 (0.034)[a]  | 5.01 (0.001)[b]  | 5.04 (0.307)[a]  | 5.31 (0.006)[c]  | 5.45 (0.003)[b] |
| TDDFT [219]    | 4.97 (0.000)[b]  | 5.08 (0.167)[a]  | 5.35 (0.065)[a]  |                  |                 |
| TDDFT [131]    | 4.86 (0.0095)[b] | 4.97 (0.2114)[a] | 5.20 (0.0178)[a] | 5.23 (0.0114)[c] | 5.53 (0.0056)[b] |

Oscillator strengths are given in parenthesis.
[a] $\pi\pi^*$ state.
[b] $n\pi^*$ state.
[c] Rydberg state.

CASPT2, and DFT/MRCI agree in the ordering of the states with the exception that the $S_3$, $S_4$ states are switched in CASPT2. A comparison of the different methods demonstrates the difficulty in calculating these states. $S_1$ energies range from 4.5 to 5.4 eV, $S_2$ from 5.0 to 7.3 eV, $S_3$ from 5.6 to 6.9 eV and $S_4$ from 5.8 to 7.0 eV. It should also be emphasized that the energies calculated with multireference methods, like MRCI and CASPT2, depend on the choice of active space. This will be even more apparent on comparisons for some of the other nucleobases where calcula-

*Table 11-5.* Vertical excitation energies in eV relative to the ground state minimum of the singlet electronic excited states in guanine

|  | $S_1$ | $S_2$ | $S_3$ | $S_4$ | $S_5$ |
|---|---|---|---|---|---|
| CASSCF [145] | 5.207[a] | 5.694[b] | 5.982[c] | 6.930[c] | 7.607[b] |
| CASPT2 [145] | 4.653 (0.188)[c] | 4.896 (0.0007)[a] | 5.422 (0.115)[c] | 5.510 (0.001)[b] | 6.711 (0.026)[b] |
| CASPT2 [143] | 4.76 (0.133)[c] | 5.09 (0.231)[c] | 5.79 ($10^{-4}$)[b] | 6.60 (0.008)[b] |  |
| CIPSI [220] | 4.76 (0.135)[c] | 5.64 (0.184)[c] |  |  |  |
| CC2 [124] | 4.98 (0.132)[c] | 5.08 (0.028)[a] | 5.38 (0.003)[ba] | 5.43 (0.141)[a] | 5.47 (0.179)[ca] |
| DFT/MRCI [171] | 4.53 (0.22)[c] |  |  |  |  |
| TDDFT [131] | 4.64 (0.0312)[c] | 4.86 (0.1275)[c] | 5.04 (0.0100)[a] | 5.17 (0.2193)[c] | 5.28 (0.0023)[b] |

Oscillator strengths are given in parenthesis.
[a] Rydberg state.
[b] $n\pi^*$ state.
[c] $\pi\pi^*$ state.

tions for different active spaces are present. Here the most reasonable active space for non-Rydberg low energy states is an active space of 14 electrons in 10 orbitals ($8 \pi + 2 n$), denoted (14,10). This active space has been used for the MRCI calculations. Roos and coworkers examined the effect of different active spaces on their CASPT2 results [125]. The results shown in Table 11-1 are using an active space of 12 electrons in 12 orbitals. A recent multistate CASPT2 (CAS-MSPT2) [126] used an active space of eight electrons in seven orbitals averaged over five states. In this type of multistate calculation even the number of states included can have a big effect on excitation energies. The effect of dynamical correlation can be seen by comparison between a CASSCF calculation and the corresponding CASPT2 that uses the same active space. Dynamical correlation has small importance on the $n\pi^*$ state, lowering the state by 0.3 eV, but it has a huge effect on the $\pi\pi^*$ states, lowering the energy by about 2 eV for the first one and more than 1 eV for the second. The basis set is another parameter that can affect the energies. For uracil the studies that have used a variety of basis sets have not shown a big effect [130], but further work is needed to examine this in greater detail. The Rydberg states have not been studied extensively by many methods. CC2 calculates that the first Rydberg state is around 6 eV [124], and TDDFT predicts that state at 5.67 eV [131].

Thymine is very similar to uracil since it differs only by a methyl group. Vertical excitation energies are given in Table 11-2. The first excited state is again an $n\pi^*$ state, and $S_2$ is a $\pi\pi^*$. The shift of the energies compared to those of uracil is generally small, but it also varies depending on the method used to calculate the excitation energies for both bases.

Experimentally the lowest peak in uracil corresponding to the bright $\pi\pi^*$ states is observed in the gas phase at 5.1 eV [132]. In thymine the first band maximum is at 4.8 eV [132]. An extended table of available experimental results taken from absorption spectra and circular dichroism is given by Roos and coworkers [125] for both molecules.

### 11.3.1.2.   Cytosine

The final pyrimidine base, cytosine, has an amino group and a carbonyl group as substituents instead of two carbonyl groups present in uracil and thymine. As a consequence the excited states are different than those in uracil and thymine. The orbitals involved in the lowest valence excited states are shown in Figure 11-3. There are three states that are quite close energetically, making it harder for the theoretical methods to predict even the right ordering of the states. Excitation energies are given in Table 11-3. Contrary to uracil and thymine, the first excited state in cytosine is a $\pi\pi^*$ state with calculated vertical excitation energy 4.4–5.1 eV. $S_2$ and $S_3$ are $n\pi^*$ states where the lone pair orbitals involved are one on oxygen and the other on nitrogen. $S_1$ and $S_2$ are calculated to be very close energetically by most methods. Like uracil and thymine, the Rydberg states have not been studied as extensively as the valence states. CC2 predicts the first Rydberg state to be at 5.53 eV [124], while TDDFT locates it at 5.26 eV [131], and DFT/MRCI above 6 eV [133].

   An extended table of available experimental results for cytosine in given by Roos and coworkers [134]. The first experimental peak is at ca. 4.6 and the second at 5.2 eV [132, 135, 136]. Gas phase REMPI spectra locate the origin of $S_1$ at 4 eV [97].

### 11.3.1.3.   Adenine

Although adenine has many tautomers that have similar energies and can be present in gas phase experiments, here we will focus only on the tautomer present in natural DNA, the 9H-adenine, shown in Figure 11-2.

   The lowest four states are two $\pi\pi^*$ states and two $n\pi^*$ states, but their ordering is not satisfactorily determined despite the considerable work done both theoretically and experimentally. The different computational methods vary widely in the ordering and energies of the states they predict. The orbitals mainly contributing to these excited states are shown in Figure 11-4 and the vertical excitation energies are given in Table 11-4. The second $\pi\pi^*$ state has the largest oscillator strength and it is designated as $L_a$, using Platt's nomenclature [137]. This is the state that the dominant configuration is mainly an excitation from the HOMO orbital (as indicated in Figure 11-4) to the LUMO. The oscillator strength of the first $\pi\pi^*$ state is 5–10 times smaller, and the dominant configuration for this state involves a mixing of HOMO → LUMO +1 and HOMO −1 → LUMO. This is the $L_b$ state. The $L_a$ state is very difficult to describe computationally and its energy depends greatly on dynamical correlation. So if one compares CASSCF and CASPT2 values using the same active space and basis sets [138], at the CASSCF level its energy is 6.8 eV, and when dynamical correlation is included at the CASPT2 level it drops to 4.9 eV. By comparison the $L_b$ state with inclusion of dynamical correlation drops from 5.2 to 4.8 eV, only by 0.4 eV. The location of Rydberg states in adenine is a point of extended discussion, especially in relation to the photophysical properties of adenine. CASPT2 locates the first Rydberg states at 5.68 eV, while DFT/MRCI locates them at 5.31 eV [139] and CC2 at 5.53 eV [124]. The CC2 method predicted considerable mixing of valence and Rydberg character for higher excited states of adenine.

As was already discussed for uracil, the choice of the active space and basis sets in multireference methods is very important. This dependence can be clearly seen in the variation of CASPT2 results obtained for adenine, as seen in Table 11-4. Serrano-Andrés and coworkers used an active space with 12 electrons in 11 orbitals and a 6-31G(d,p) basis set [140]. Domcke and coworkers used an active space (12,10), of either $10\pi$, $8\pi+2n$ or $9\pi+1n$, and MP2 geometry [138]. Blancafort used an active space (16,12) of 10 $\pi$ and 2 $n$ orbitals [141]. The energies of the $\pi\pi^*$ states vary by 0.3 eV, while the energies of the $n\pi^*$ states change by as much as 1 eV, depending on the active space used in the CASPT2 calculations.

Experimentally the gas phase UV/V is absorption spectrum of adenine has a maximum at 4.92 eV, which is red shifted in aqueous solution to 4.77 eV [132]. This band contains at least two electronic transitions with maxima at about 4.6 and 5.0 eV, with the first transition being weaker than the second [142, 143]. A second band appears at higher energies, around 6 eV. In the gas phase R2PI spectra of jet-cooled adenine there are two origins at 4.40 and 4.48 eV, where the low energy one corresponds to a $n\pi^*$ transition, and the higher energy one corresponds to $\pi\pi^*$ [105].

### 11.3.1.4. Guanine

The ground state geometry of guanine is almost planar with the amino group being slightly pyramidalized. The tautomer presented here is 9H-guanine, shown in Figure 11-2. The excited states of guanine calculated using a variety of methods are shown in Table 11-5. The orbitals mainly involved are shown in Figure 11-4. As in adenine the first four singlet excited states are two $\pi\pi^*$ states and two $n\pi^*$ states. The first two states, $S_1$ and $S_2$, are predicted by all methods to be $\pi\pi^*$ states with similar oscillator strengths. This differs from adenine where the brightest state is the second $\pi\pi^*$ state. The dominant configurations in these states differ from adenine as well. In guanine the $S_1$ state is mostly HOMO $\rightarrow$ LUMO excitation ($L_a$) while $S_2$ is mostly HOMO $\rightarrow$ LUMO +1 ($L_b$) without much excitation from HOMO $-1$ [143]. So the ordering between $L_a$ and $L_b$ is reversed in guanine compared to adenine. The energy for $S_1$ is lower than the corresponding state in adenine and this trend is predicted by all methods applied. This should be expected since it is easier to predict qualitative trends between molecules that differ in their substituents than predict accurate absolute excitation energies [144]. The next states arise by excitation from the lone pairs, which here are one on nitrogen and one on oxygen. $S_3$ is an $n_O\pi^*$ state while the second $A''$ state is an $n_N\pi^*$ state. These states appear to be well separated from the $\pi\pi^*$ states located at higher energies. The first Rydberg state at the CASPT2 level is at 4.90 eV [145]. The CC2 method again predicts considerable mixing of valence and Rydberg character with the first mainly Rydberg state at 5.08 eV [124]. In general, in these calculations where Rydberg states were considered, the first one was found to be low in energy and possibly the $S_2$ or $S_3$ state. Involvement of Rydberg states may be important in the photophysical properties of guanine, as has been seen in adenine, although there have not been studies that explicitly include Rydberg states

when studying the photophysical properties of guanine. The relaxation studies will be discussed in more detail in the following section.

Experimentally guanine shows spectral features with similarities to those of adenine. There are two peaks at 4.51 and 4.95 eV with relatively low intensity and much stronger bands higher than 6 eV. The oscillator strengths for the first two peaks are similar also, 0.14 and 0.21, respectively [143, 146].

### 11.3.2.     Photophysics

All the nucleic acid bases absorb UV radiation, as seen in Tables 11-1, 11-2, 11-3, 11-4, and 11-5, making them vulnerable to the UV radiation of sunlight, since the energy of the photons absorbed could lead to photochemical reactions. As already mentioned above, the excited state lifetimes of the natural nucleobases, and their nucleotides, and nucleosides are very short, indicating that ultrafast radiationless decay to the ground state takes place [6]. The mechanism for nonradiative decay in all the nucleobases has been investigated with quantum mechanical methods. Below we summarize these studies for each base and make an effort to find common mechanisms if they exist.

#### 11.3.2.1.     Uracil

The electronically excited singlet states of uracil, and how they can lead to efficient radiationless decay to the ground state, were initially investigated using MRCI methods by Matsika [147] and later with other methods that agree with the MRCI results in the more general features [92, 126, 128, 148–150]. The discussion in this section describes MRCI results for free uracil in detail, along with the studies that have been



*Figure 11-5.* Energy level diagram of minima and conical intersections involved in the radiationless decay in uracil. Energies and structures taken from Ref. [147, 224]

published after that. Figure 11-5 summarizes the features in the PES that were found to be important in the photophysical behavior of uracil.

Vertical excitation energies for the first two excited singlet states, $S_1$ and $S_2$, are given in Table 11-1. The first excited state is a dark state involving excitation from the lone pair on oxygen to a $\pi^*$ orbital ($n_O\pi^*$) while the second is a bright $\pi\pi^*$ state. Initial absorption of UV radiation excites uracil to the $S_2$ state and subsequent radiationless decay necessarily involves both $S_1$ and $S_2$ excited states. The other singlet states are higher in energy and are not considered further, although they could be involved indirectly in the photophysical behavior of uracil, especially at higher energies. The vertical excitation energies depend strongly on the level of correlation, and the MRCI expansion has to be designed so that some type of $\sigma$ correlation is included. See original work for more details [147]. Vertical excitation energies at the MRCI level are 4.8 and 5.8 eV for the $S_1$ and $S_2$ states respectively.

Pathways that can lead to radiationless decay to the ground state after initial excitation to the bright $S_2$ state have been calculated using MRCI. These pathways are provided via conical intersections, as shown in Figure 11-5. The first conical intersection seam is between the first two excited states, $S_2$ and $S_1$, and is easily accessible from the Franck Condon region. The molecule is significantly distorted from planarity at this geometry as seen in Figure 11-5. The $C^5C^6$ bond is stretched to 1.48 Å loosing its double bond character and simultaneously the $C^4O^8$ bond is stretched from 1.20 to 1.25 Å. Minimum energy paths on the $S_1$ surface starting from this conical intersection can lead to two different directions. In one direction the $S_1$ minimum can be reached, while in another direction a conical intersection seam between $S_1$ and the ground state will be accessed. Access to this conical intersection will lead to ultrafast radiationless decay to the ground state, while access to the $S_1$ minimum will result in a long lifetime, since this is a dark state with low oscillator strength for radiative emission. The radiative lifetime calculated for this state is $2 \times 10^{-5}$ s. Both of these possibilities have been observed in experiments. He et al. [101, 102] first observed the dark state in the gas phase of methylated uracil molecules and in complexes of thymine with water molecules. Recently a dark state has also been observed in aqueous solutions of 1-cyclohexyluracil by Kohler and coworkers [151]. In that work the branching between the dark state and final decay to the ground state is observed to be 60% for radiationless decay.

The $S_1$ minimum and the minimum energy point on the conical intersection seam $S_1/S_0$ are almost isoenergetic with an energy ca. 4.1 eV. The lifetime of the dark state, if the $S_1$ minimum is reached, depends on the ability of a wavepacket on the $S_1$ minimum to reach the seam. Although these points are isoenergetic, they have very different conformations and lie in different regions of the PES, so it is important to investigate pathways connecting them. A transition state exists on the $S_1$ surface that connects the $S_1$ minimum with the $S_1/S_0$ conical intersection seam. This transition state creates a barrier of 0.65 eV, which makes radiationless relaxation from the $S_1$ minimum to the ground state difficult. The imaginary frequency is 2045 cm$^{-1}$. The motion involves mainly twist of the $C^5$—$C^6$ bond, especially the $C^6H^6$ fragment, which is the motion needed to reach the conical intersection. The wavefunction of the

$S_1$ state at that transition state involves mixing of the $n_O\pi^*$ and $\pi\pi^*$ configurations. The geometry also reflects this mixing. The $C^4O^8$ bond is 1.25 Å, a value between that of the $S_1$ minimum (1.35 Å) and the conical intersection (1.20 Å). The $C^4C^5$ bond is 1.44 Å, again between the values at the $S_1$ minimum, 1.36 Å, and at the conical intersection, 1.49 Å. Alternative radiationless decay pathways from the $S_1$ minimum have been investigated by searching for different conical intersections between $S_1$ and the ground state. A conical intersection seam at 5.5 eV was located but the energy is much higher than the conical intersection $S_1/S_0$ and the barrier, so this conical intersection is not expected to be important for radiationless decay from the $S_1$ surface. Overall these calculations suggest that nonadiabatic transition from $S_2$ to $S_1$ is very fast, but subsequent relaxation to the ground state depends on the branching between the pathway leading to the conical intersection $S_1/S_0$ and the pathway leading to the dark state minimum. Escape of population from this minimum is hindered by a barrier. This seems to be responsible for the nanosecond lifetimes observed in the gas phase [101, 102].

Pathways leading to the ground state have been investigated using different methods, CIS, CR-EOM-CCSD(T) [128, 148], TDDFT [92, 149] and CASPT2 [126, 150]. Specifically the $S_1/S_0$ conical intersection has also been found by others to be the leading path for deactivation to the ground state in uracil and thymine[92, 128, 148–150]. Zgierski et al. [128, 148] describe the $S_1/S_0$ conical intersection as an intersection between a biradical state and the ground state, but the structure is very similar to that found by the other methods.

Hudock et al. [126] used the ab initio molecular dynamics multiple spawning method to go beyond the static picture based on PES and include the time dependent dynamical behavior and predict time-resolved photoelectron spectroscopy results. According to these results the first ultrafast component of the photoelectron spectra of uracil corresponds to relaxation on the $S_2$ minimum rather than nonadiabatic transitions to the $S_1$ state. The authors suggest that the radiationless relaxation from $S_2$ to $S_1$ through the $S_2/S_1$ conical intersection is responsible for the second picosecond lifetime observed. Several trajectories were calculated starting from vertical excitation on the $S_2$ surface and it was observed that only a small fraction of them ended up on the $S_1$ surface in the first 500 fs. Furthermore, the electronic states of the uracil cation, formed upon photoionization, play an important role on the photoelectron spectra. This study points to the fact that time-dependent calculations mimicking the experimental observable may reveal aspects of excited state dynamics that are difficult to extract relying only on time-independent electronic structure potential energy surfaces. So the ultimate fate of excited states can only be determined through dynamical studies. Since these studies depend on the quality of the PES and require many quantum chemical calculations, they have been limited so far, but it is expected that a lot more contributions will be made in this area in the near future. For uracil the current study explores initial dynamics on the $S_2$ surface, but further work is needed to follow the dynamics after that and finally to the ground state.

### 11.3.2.2. Thymine

As discussed earlier, thymine is very similar to uracil in its excited states pattern. This is also true for its radiationless decay mechanism except from the fact that the excited state lifetime in thymine is somewhat longer than in uracil. Theoretically the mechanism for radiationless decay has been studied using CASPT2 electronic structure methods [150, 152].

The excites states in thymine, as in uracil, are an $S_1$ $n\pi^*$ state and an $S_2$ $\pi\pi^*$ bright state. So initial UV absorption will lead to the bright $S_2$ state. This state crosses with the $S_1$ state first. Perun et al. [152] found three different conical intersections between $S_1$ and $S_0$ present in thymine. All conical interections are characterized by strongly out-of-plane distorted geometries of the heterocyclic ring. Figure 11-6 shows the structure of thymine at two of these conical intersections. The energetically lowest one is similar to the one found in uracil and corresponds to a crossing between the ground state and the $\pi\pi^*$ state. The structure involves twisting of the $C^5C^6$ bond and the methyl group being almost perpendicular to the plane of the molecule, similarly to the structure found in uracil (see Figure 11-6a,b). This conical intersection is accessible without a barrier from the minimum of the $\pi\pi^*$ state and can provide a pathway for direct quenching of the population of the lowest $\pi\pi^*$ state of thymine. The energy at this point is 4.35 eV at the CASPT2 level while the energy at the $\pi\pi^*$ minimum is 4.48 eV. The other two conical intersections involve the $n\pi^*$



*Figure 11-6.* Structures of representative conical intersections $S_1/S_0$ in the pyrimidine bases, uracil, thymine, and cytosine. Uracil structures (**a,d**) are taken from Ref. [147, 210]. Thymine structures (**b,e**) are taken from Ref. [152]. Cytosine structures (**c,f**) are taken from Ref. [157]

state and $S_0$ and are much higher in energy, about 7.5 eV. Merchan et al. [150] have also calculated the $C^5C^6$ twisted conical intersection using CASPT2 methods located at 4.0 eV.

The ab initio molecular dynamics study by Hudock et al. discussed above for uracil included thymine as well [126]. Similarly to uracil, it was found that the first ultrafast component of the photoelectron spectra corresponds to relaxation on the $S_2$ minimum. Subsequently a barrier exists on the $S_2$ surface leading to the conical intersection between $S_2$ and $S_1$. The barrier involves out-of-plane motion of the methyl group attached to $C^5$ in thymine or out-of-plane motion of $H^5$ in uracil. Because of the difference of masses between these two molecules, kinematic factors will lead to a slower rate (longer lifetime) in thymine compared to uracil. Experimentally there are three components for the lifetimes of these systems, a subpicosecond, a picosecond and a nanosecond component. The picosecond component, which is suggested to correspond to the nonadiabatic $S_2/S_1$ transition, is 2.4 ps in uracil and 6.4 ps in thymine. This difference in the lifetimes could be explained by the barrier described above.

### 11.3.2.3.    *Cytosine*

Cytosine was the first nucleobase whose radiationless decay was studied with quantum mechanical methods. Nevertheless, its first excited states are not so clearly separated as in uracil and thymine, and this causes complications in the computational studies of the photophysics. So, many computational studies have been reported to elucidate the mechanisms for radiationless decay to the ground state but, not always with the same conclusions.

Overall the different calculations have reported that three different conical intersections between $S_1$ and $S_0$ exist in cytosine. One of them involves twisting of the $C^5C^6$ double bond and the states crossing are the ground state and the $\pi\pi^*$ state which, however, looks more like a diradical at this distorted geometry (ci1). The second conical intersection involves $N^3C^4$ torsion and the $S_1$ state has been called either an $n\pi^*$ or $\pi\pi^*$ state, although it is the same state and only the nomenclature differs (ci2). The third conical intersection involves the $n_O\pi^*$ excited state and has the CO bond elongated (ci3). All three conical intersections have somewhat similar energies and, depending on the method used, their relative importance varies.

The first study, by Ismail et al. [153], used the CASSCF method with a 6-31G* basis set and an active space of 14 electrons in 10 orbitals to locate conical intersections and pathways connecting them to the Franck Condon region. Two such conical intersections were identified in that work, the ci2 and ci3, as defined above. In that work the barrier leading to ci2 was calculated to be 10 kcal/mol, too high to make this conical intersection relevant. But the barrier leading to ci3 was found to be much smaller, 3.6 kcal/mol, and it was concluded that ci3 is involved in the dominant decay path . Reaching this intersection requires first a conical intersection between the $\pi\pi^*$ state, which is vertically the $S_1$ state, and the $n_O\pi^*$ state, which is vertically the $S_2$ state. Merchan and Serrano-Andrés followed up this study [140] using a method

that included dynamical correlation though, CASPT2, and predicted that an $S_1/S_2$ conical intersection is not necessary and the $S_1$ $\pi\pi^*$ state can directly evolve to a conical intersection with the ground state. The structure of cytosine at this conical intersection is very similar to that of ci3, but the character of $S_1$ is $\pi\pi^*$ instead of $n_O\pi^*$. In fact, it has been shown that in the vicinity of these two conical intersections all three states, $S_0$, $S_1$, and $S_2$, are close in energy and upon optimization a three-state conical intersection has been located [154].

Later, several other studies appeared where a different conical intersection, the twisted conical intersection (ci1), was located and was connected to the $S_1$ minimum with small barriers [133, 155–157]. MRCI calculations by Kistler and Matsika located all the above mentioned conical intersections between $S_1$ and $S_0$ and constructed detailed pathways that connect them to the Franck Condon region. Initial absorption of a UV photon excites the ground state system to $S_1$, the bright $\pi\pi^*$ state at 5.14 eV. The minimum energy path (MEP) on $S_1$ leads to a minimum at 4.31 eV by stretching the carbonyl by about 0.1 Å and folding the ring slightly in a butterfly fashion along the $N^1/C^4$ axis. From this minimum the base can continue to distort on the MEP in two different ways, both of which lead to conical intersections through barriers of about 0.15 eV. Figure 11-7 shows the pathways connecting the $S_1$



*Figure 11-7.* Minimum energy pathways on the $S_1$ surface of cytosine connecting the $S_1$ minimum to two $S_1/S_0$ conical intersections. The five singlet state energies at the MRCI level are shown. Energies are given in eV with respect to the ground state minimum energy. ci$IJ$ represents conical intersection between states $S_I$, $S_J$. (From Ref. [157])

minimum to the two conical intersections. Progression in the "sofa" direction of distortion (as is indicated in the figure), where $N^3$ puckers out of plane, leads to a conical intersection between $S_1$ and the closed-shell $S_0$ surface with energy of 4.27 eV. This is ci2 as defined above. Interestingly the MRCI calculations predict a very small barrier leading to ci2, although CASSCF had predicted a much higher barrier as discussed in the previous paragraph [153]. CASPT2 calculations reported recently also predict a small barrier of 0.2 eV [156]. The $S_1$ minimum is also connected to a second conical intersection that involves twisting of the $C^5C^6$ bond (ci1), as shown in Figure 11-7. A transition state along the pathway creates a low barrier of 0.15 eV. This conical intersection has an energy of 3.98 eV, which is about 0.3 eV lower than the global stationary minimum or the sofa ci2, making this the more energetically favored conical intersection [157]. The third conical intersection, ci3, is too high at the MRCI level, and also at the CASPT2 level as calculated by Blancafort [156]. The low barriers connecting the $S_1$ minimum to the conical intersections are consistent with REMPI experiments that show sharp lines only in a very narrow range of energies [97].

### 11.3.2.4.     Adenine

The photophysical properties of adenine have intrigued chemists from early on. Broo studied adenine and 2-aminopurine (2AP) in order to understand their differences in photophysical properties. Adenine like all natural nucleobases has very short excited state lifetimes and low quantum yields of fluorescence, while 2AP, which differs from adenine in the position of the amino group, has long lifetimes and strong fluorescence, making it a very useful fluorescent probe. In Broo's work it was observed that the first excited state is a $\pi\pi^*$ at vertical excitation but crosses with an $n\pi^*$ state which becomes the $S_1$ state adiabatically at the minimum. The large out-of-plane distortion on the $n\pi^*$ state opens up a deactivation channel in adenine compared to 2AP. In 2AP, on the other hand, the $S_1$ state always has a $\pi\pi^*$ character.

After the recent experimental studies [83, 84] measured more accurately the ultrashort excited state lifetimes, the photophysical properties of adenine became again highly interesting. In 2002 Domcke and coworkers [158] published a study that showed that a Rydberg state is present close to the valence excited states and it has a conical intersection with the ground state along the $N^9H$ stretching coordinate of the azine group. They argued that the predissociation of the $\pi\pi^*$ and $n\pi^*$ states by the Rydberg state (a $\pi\sigma^*$ state) and the conical intersection of the $\pi\sigma^*$ state with the ground state provide the mechanism for the ultrafast deactivation of the excited states of adenine. Similar mechanism has been observed in other aromatic biomolecules, i.e. pyrrole [159]. This mechanism however can only be operative in free bases since in nucleotides and nucleosides $N^9$ forms a bond with the sugar rather than the H. Discussions about the importance of this mechanism have not been settled and experimentalists try to provide evidence in support or against its importance

[108, 112, 114, 115]. Further theoretical calculations have been done including these conical intersections between the $\pi\sigma^*$ state and the ground state [160, 161].

This initial study in adenine was very different from studies on other bases where the main vibrational motions involved in the radiationless deactivation are in the ring. Several other studies however have appeared since then, where the mechanism in adenine is shown to proceed via conical intersections involving out-of-plane deformations of the six-membered ring [138, 139, 141, 160, 162–165]. Two conical intersections between $S_1$ and $S_0$ involving such torsions have been calculated by several groups [138, 139, 141, 163, 164]. The conical intersections involve twisting of the $C^2N^3$ and $N^1C^6$ bonds in the ring (denoted CI23 and CI16 respectively by Domcke and coworkers [138]). Figure 11-8 shows the structure of adenine at these two conical intersections. The lowest minimum of the $S_1$ state has $n\pi^*$ character. This minimum is connected to the two different conical intersections with the ground state with paths that have very small barriers, ca. 0.1 eV. Although the ordering of the states does not agree in all calculations, it is agreed in all of these studies that the brightest state is the second $\pi\pi^*$ state ($L_a$). If initial absorption of light excites mostly to that state it will cross with the lower energy states and eventually become the $S_1$ state. It can then cross the ground state $S_0$ with a conical intersection that involves twisting along $C^2N^3$ (CI23). This mechanism is consistent with R2PI experiments that show sharp vibronic peaks in a rather narrow wavelength range ($1100\,\text{cm}^{-1}$ from the origin) [105]. Figure 11-9 shows the pathways connecting the $S_1$ minimum to the two conical intersections as calculated by Domcke and coworkers [138]. The left panel of the figure also shows the higher bright $L_a$ state falling rapidly in energy and crossing with the ground state a CI23. Zgierski and coworkers proposed a biradical



a) $^1La/S_0$ (CI23)

c) $^1La/S_0$ (CI23)

b) $^1n\pi^*/S_0$ (CI16)

9H-Adenine    9H-Guanine

*Figure 11-8.* Structures of representative conical intersections $S_1/S_0$ in the purine bases, 9H-adenine and 9H-guanine. Adenine structures (**a,b**) are taken from Ref. [138]. Guanine structure (**c**) is taken from Ref. [171]

*Figure 11-9.* CASSCF potential-energy profiles of the ground-state $S_0$ (*circles*), the $^1n\pi^*$ state (*triangles*), the $^1L_b$ state (*squares*), and the $^1L_a$ state (*filled squares*) of the 9H-adenine along the linear interpolation reaction path from the equilibrium geometry of the $^1n\pi^*$ state to the CI32 (**a**) and CI16 (**b**) conical intersections. The diabatic correlation of the states is shown in (**a**). (From Ref. [138])

mechanism, as in the pyrimidine nucleobases, where the structure distorts as in the $\pi\pi^*/S_0$ intersection proposed by the other authors [166].

The conical intersections involving out-of-plane deformations are probably responsible for radiationless deactivation at lower UV energies whereas at higher energies it is proposed that the Rydberg state may be accessed opening the channel for hydrogen abstraction [160]. Additionally, a mechanism leading to opening of the 5-membered ring has been found, but the barrier needed to overcome to reach a conical intersection involved in this mechanism is high (ca. 1.5 eV) [160].

### 11.3.2.5.    Guanine

Radiationless decay in guanine has been studied less compared to the other nucleobases. CASPT2, DFT and CCSD methods have been employed to study the photophysical properties of guanine [145, 167–172]. The first excited state is a bright $\pi\pi^*$ $L_a$ state at 4.65 eV using CASPT2 theory. After initial excitation to this bright state it relaxes to a minimum at 3.97 eV. This minimum is connected to a conical intersection with the ground state at 3.7 eV. Thus the energy needed to reach the conical intersection after relaxing to the minimum is only 0.17 eV [145]. Pathways from the other higher energies have been calculated as well. Figure 11-10 taken from the work of Chen and Li [145] shows the different possible pathways calculated at the CASPT2 level. On the $n_N\pi^*$ surface a conical intersection between $n_N\pi^*$

*Figure 11-10.* The nonradiative deactivation paths for the lowest two $^1\pi\pi^*$ state and the $^1n_N\pi^*$ state of 9H-guanine. (From Ref. [145])

and $\pi\pi^*$ ($L_a$) is located, which suggests that the $n_N\pi^*$ could transform to the $\pi\pi^*$ excited state first and then follow the deactivation pathway described for that state. The Rydberg $\pi\sigma^*$ state was also considered, but it was found that dissociation of the NH bond of the six-membered ring is difficult because of a significant barrier.

Zgierski and coworkers proposed for guanine the same biradical mechanism that was proposed for all the other bases [172]. Furthermore, on the fly molecular dynamics using density functional theory have been used to study initial evolution along the $S_1$ surface of methylated guanine [167–169].

Guanine has six tautomers that are relatively close energetically and are probably present in gas phase experiments. This made the assignment of the spectra much more complicated, and problems in assignments arose necessitating theoretical work. The situation was additionally complicated because the excitation energies of the different tautomers have different energies, lifetimes and relaxation mechanisms. For example, the 9H tautomer is believe now to relax to the ground state so fast that it has probably not been seen at all in any spectra [109–111, 113]. The first excited state of the six tautomers and its radiationless decay has been studied theoretically by Chen and Li using CASPT2 methods [170], and by Marian using DFT/MRCI methods [171]. These calculations suggest that the deactivation of 9H-guanine is much more rapid than the other tautomers.

### 11.3.2.6. General Features of Conical Intersections in Nucleobases

An important question in the topic of radiationless decay in nucleobases is whether a common mechanism exists that operates in all of them. The studies above show that the mechanism in each nucleobase depends on the details of the excited states

present and their relative energies. Nevertheless, there are some similarities between the different systems. Specifically, there are always conical intersections between a $\pi\pi^*$ state and the ground state, and the structures of these conical intersections share some similarities.

Figures 11-6 and 11-8 show the structures of some representative $S_1/S_0$ conical intersections located in the nucleobases. Pyrimidine bases have common a $C^5C^6$ double bond. The first $\pi\pi^*$ excited state in all of them has a conical intersection with the ground state where this bond is twisted. The structures for uracil, thymine, and cytosine at the conical intersections are shown in Figure 11-6a,b,c. The dihedral angle $<XC^5C^6H$ (X=H,CH$_3$) is not the same in all bases but the main twisted structure is. The wavefunction shows biradical character at this point. Cytosine also has a double bond between $N^3C^4$. A second conical intersection with the ground state exists where this bond is elongated and twisted forming a 'sofa' like conformation as seen in Figure 11-6f. Uracil and thymine also have additional conical intersections arising from crossing of the $n_O\pi^*$ with the ground state (see Figure 11-6d,e).

Purine bases are bicyclic, but it is the six-membered ring that has been seen to participate mostly in conical intersection distortions. Some of the structures of the conical intersections observed are shown in Figure 11-8. All purines have a $C^4C^5$ double bond but it is shared between the rings, so it cannot be twisted easily to destabilize the ground state and bring $S_1$ and $S_0$ together in energy. However, the $C^2N^3$ double bond can be distorted after excitation to the $\pi\pi^*$ state. In both adenine and guanine a conical intersection between the $L_a$ state and the ground state involves this bond twisting so that $C^2$ comes out of the plane (see Figure 11-8a,c). Adenine has another possibility with the $N^1C^6$ bond which also deforms to form a conical intersection with the ground state (Figure 11-8b).

A question that becomes obvious at this point is what happens to the molecules that have similar structures to the natural bases but have different photophysical properties, i.e. they fluoresce. These molecules have similar main structure to the bases, similar ring systems and double bonds, and so, according to the previous discussion, similar conical intersections should be expected. If that is true, and conical intersections facilitate efficient radiationless decay, why do these molecules fluoresce instead of decaying nonadiabatically? That is a question that has occupied a number of scientists and some answers and insights are given in the following section.

### 11.3.3.    Fluorescent Analogs

Analogs of the DNA bases have been synthesized which display a significant fluorescence quantum yield upon excitation, unlike the natural DNA bases, while retaining much of the structural requirements necessary for base-pairing and formation of a helical stack with the natural DNA bases. Because their fluorescence is often strongly influenced by their immediate environment, with significant quenching occurring when the analog is stacked with surrounding bases, and because of their lower, and thus selective, excitation energies compared to the natural bases, these analogs have

been used as probes for studying conformational dynamics of DNA strands. This section will review the literature for these analogs, starting with the pyrimidine analogs, and then the purine analogs.

Some examples of fluorescent pyrimidine analogs which have been synthesized and used as probes are 5-methyl-2-pyrimidin-(1H)-one (5M2P) [173–176], pyrrolocytosine (p-Cyt) [177], and the $N^9$-H conjugate acid of N3,N4-ethenocytosine ($\epsilon$-CytH$^+$) [178, 179]. The structures with atomic numbering are shown in Figure 11-11. Kistler and Matsika studied the photophysics of 5M2P theoretically using MRCI and compared the results to those for cytosine, the DNA base 5M2P most closely resembles [157, 180]. It was shown that the stationary points on the $S_1$ $\pi\pi^*$ PES, as well as two $S_0$–$S_1$ conical intersections located on this surface, display almost identical ring distortion for the two bases, but energetic differences for these points predict the very different photophysical behaviors observed experimentally for cytosine and 5M2P. The $S_1$ PES for these two bases are shown in Figure 11-12 where they are superimposed for comparison. Cytosine excites initially 0.8 eV higher than 5M2P, with its $S_1$ population having a virtually unimpeded access to the two conical intersection channels to the ground state PES, supporting ultrafast radiationless decay of the excited state. The $S_1$ $\pi\pi^*$ PES for 5M2P, however, has a minimum which is energetically deep enough to make the two $S_0$–$S_1$ conical intersections located much less energetically accessible than they are in cytosine, thus binding vibrational states and promoting fluorescence decay of the excited state. The largest energetic difference for the $S_1$ PES of these two bases was at vertical excitation,



*Figure 11-11.* The structures of pyrimidine analogs described in the text, along with the general structure of $C^4$- and $C^5$-substituted 2P derivatives

*Figure 11-12.* S$_1$ pathways for cytosine and 5M2P from vertical excitations to the "sofa" and "twist" conical intersections. Cytosine paths are shown in *blue*, and 5M2P are shown in *green*. MRCI energies are given in eV. (From Ref. [144])

implying that cytosine is more vibrationally excited than 5M2P, and so has more initial force to distort from its FC geometry than 5M2P. The reasons for this excitation energy difference were further studied theoretically, by calculating the excited state energies for 2-pyrimidin-(1H)-one (2P) derivatives with electron-donating or electron-withdrawing groups at the C$^4$ or C$^5$ positions [144]. Using a simple frontier



*Figure 11-13.* HOMO and LUMO orbitals are shown for several 2P derivatives, where X means electron-donating group (here it is amino), and Z means electron-withdrawing group (here it is nitrosyl). Orbital levels are relative and qualitative. MRCI S$_1$ ($\pi\pi^*$) energies are shown at *bottom*, in eV. (From Ref. [144])

molecular orbital model, it was shown that $\pi$-donating groups at $C^4$, such as amino in the case of cytosine (4-amino-2-pyrimidin-(1H)-one) (Label 2, Figure 11-13), increased the $\pi\pi^*$ (HOMO-LUMO) energy compared to 2P (Label 3, Figure 11-13) by destabilizing the LUMO relative to the HOMO. This effect in cytosine was almost entirely removed by rotating the amino group so that the amino group was perpendicular to the ring (Label 2, Figure 11-13), thus eliminating the $\pi$-donating character of the group. Likewise, electron-withdrawing groups at $C^4$ lowered the $\pi\pi^*$ energies relative to 2P by destabilizing the HOMO relative to the LUMO, and therefore lowering the energy gap (Label 3, Figure 11-13). By assigning values of the Hammett substituent parameter $\sigma_p^+$ to the $C^4$ groups, a predictive linear trend was shown between $\sigma_p^+$ and the $\pi\pi^*$ energies. This trend was parallel with that seen for the experimentally determined excitation energies for the handful of 2P derivatives that have been measured. Figure 11-14 shows these correlations between $\sigma_p^+$ and the $\pi\pi^*$ energies. The dark $n\pi^*$ state energies also follow this trend for the $C^4$-substituted derivatives. $C^5$ substituents on 2P were shown to have either a small lowering effect on the $\pi\pi^*$ energies compared to 2P, due to destabilizing the HOMO relative to that of 2P (Label 16, Figure 11-13), or no effect. The methyl group in 5M2P has almost no energetic influence compared to 2P. It should be noted that in order to theoretically predict if these pyrimidine analogs fluoresce, a complete topological exploration of the $S_1$ PES for each analog, including stationary points and



*Figure 11-14.* Plot of $S_1$ ($\pi\pi^*$) energies in eV, with respect to $\sigma_P^+$ for several $C^4$-substituted 2P derivatives. *Squares* are MRCI results, and *hour-glasses* are from experimentally determined absorption maxima. (From Ref. [144])

conical intersections, must be carried out. With regards to the pyrimidine analogs, this sort of theoretical analysis has only been done for 5M2P [180].

Like 5M2P, the fluorescent bicyclic cytosine analog p-Cyt (Figure 11-11) also excites lower than Cyt, by over 1 eV. It has been studied theoretically by Thompson and coworkers using CIS and TDDFT methods to calculate excited state energies [181, 182]. Kistler and Matsika, using MRCI, included this base in their report on 2P substituent effects [144], since this base has both $C^4$ and $C^5$ substitution on the 2P ring system. Because of the low $\pi\pi^*$ energy calculated for p-Cyt compared to 2P, it was concluded that the five-membered ring was acting as a conjugation-extending group or an electron-donating group at $C^5$, both types of groups tending to lower $\pi\pi^*$ energy compared to 2P, rather than as an electron donating group from $N^7$ at $C^4$, which would raise the $\pi\pi^*$ energy, as it does in Cyt. Like 5M2P, p-Cyt has a lower excitation energy compared to Cyt, but without a complete analysis of the $S_1$ PES for p-Cyt outside the FC region, as has been done for 5M2P [180], a theoretical support of fluorescence for p-Cyt cannot be made. Currently such a study on p-Cyt has not been reported.

The cytosine analog $\epsilon$-Cyt (Figure 11-11) does not itself fluoresce as the free base. This base absorbs at about the same energy as cytosine (4.65 eV) [178]. In a low pH environment, however, the absorption is red-shifted by about 0.8 eV, and fluorescence is displayed. The species which fluoresces was identified as the $N^9$-H conjugate acid of $\epsilon$-Cyt, $\epsilon$-CytH$^+$ (Figure 11-11). Kistler and Matsika [144] showed theoretically using MRCI that protonation of $N^9$ results in a similar red shift of the $\epsilon$-Cyt $\pi\pi^*$ energy, and that this cytosine analog follows the same trend in absorption energies of other 2P derivatives.

Three fluorescent purine analogs will be described here: 2-aminopurine (2AP) [183] and 8-vinyladenine (8VAD), both of which mimic adenine, and N1,N6-ethenoadenine ($\epsilon$-AD). Their structures are shown in Figure 11-15. 2AP is probably the most widely used fluorescent DNA base analog, being utilized as a probe for DNA conformational dynamics [184, 185], due to the environmental specificity of



*Figure 11-15.* The structures of 2AP and 8VAD, $\epsilon$-AD, and $\epsilon$-ADH$^+$

its quantum yield. More recently, it has been used by Fiebig as a probe to study excited state delocalization in nucleotide oligomers [186], and its excited state electronic structure has been studied with double resonance spectroscopy by Marian and coworkers [187], and also in the presence of electric fields using Stark spectroscopy by Stanley and coworkers [188]. Like the fluorescent pyrimidine analogs compared to the natural pyrimidines they mimic, the first absorption band of 2AP is red shifted compared to that of 9H-adenine, by about $3700 \, cm^{-1}$ [96, 189]. The excited states base have been studied theoretically with a wide variety of theoretical methods, including MCSCF [190], MCQDPT [191], CASPT2 [192], DFT [183], and MRCI [188]. These studies all point to the bright absorption at vertical excitation being the first $\pi\pi^*$ state, with HOMO-LUMO character. Broo [193] used CIS and semi-empirical methods to show that the nonradiative decay channel present in excited adenine is not present in 2AP. This was further supported by Marian and coworkers [187], who used TDDFT to show that the conical intersection between the $S_1$ $\pi\pi^*$ state and the ground state, where 2AP is distorted out of plane at the amino-bearing ring carbon, is blocked by a barrier and is too high in energy to be easily accessed by the excited state population. A similar scenario was reported by Domke and coworkers [194], using CASSCF and CASPT2, showing similar distortion of 2AP at the gs/$\pi\pi^*$ conical intersection and a barrier blocking access to it on the $S_1$ PES. Serrano-Andres and coworkers used CASSCF to map out the PES and CASPT2 to refine the energies [164]. In their study, they presented an analysis of all the low-lying excited states, including the $n\pi^*$ states, and they contrasted their 2AP results with those for adenine, in which the excited state population had unimpeded access to a conical intersection between the $S_1$ $\pi\pi^*$ state and the ground state, supporting a nonradiative decay mechanism for excited adenine.

8VAD (Figure 11-15) is among the most recently utilized fluorescent nucleobase analogs [195–197]. Its absorption maximum is red-shifted compared to adenine by about 30 nm, at 291 nm, and it emits at 382 nm. Its base-pairing is the same as that of adenine. Currently only one theoretical study has been reported for 8VAD, by Kenfack and coworkers, using TDDFT and CIS to calculate the excited state properties and transitions [198]. In that study, the bright state was shown to be a $\pi\pi^*$ state wth HOMO-LUMO character.

$\epsilon$-AD (Figure 11-15) has been utilized as a fluorescent probe in DNA since the 1970s [199]. Unlike $\epsilon$-Cyt, $\epsilon$-AD is fluorescent as the free base, but not in its protonated form, $\epsilon$-ADH$^+$. Its absorption energy is about 30 nm red-shifted compared to that of adenine. Theoretical calculations to study its excited states have been done using CIS, TDDFT, and MCQDPT [179]. Both absorption and emission energies were studied in this report. Protonation of $\epsilon$-AD was shown to increase the $\pi\pi^*$ absorption energy compared to the free base by about 0.4 eV.

### 11.3.4. Three State Conical Intersections

Although the importance of two-state conical intersections in nonadiabatic processes has been established [10, 26], the occurrence and relevance of accidental three-state

conical intersections is only now emerging. Three-state degeneracies imposed by symmetry have been studied in the context of the Jahn-Teller problem for many decades [200], but only minor attention has been given to accidental three-state degeneracies in molecules [201]. Since most molecular systems in nature have low or no symmetry, these accidental intersections may have a great impact on the photophysics and photochemistry of molecular systems, as has been found in accidental two-state intersections [10, 20, 24, 26]. Three-state degeneracies can provide a more efficient relaxation pathway when more than one interstate transition is needed, they introduce more complicated geometric phase effects [202–204], and they can affect the dynamics and available pathways [205, 206]. The exact way these features affect the dynamical behavior of molecules is a question that has yet to be studied extensively. According to the noncrossing rule [13] three states can become degenerate in a subspace of dimension $N^{int} - 5$, where $N^{int}$ is the number of internal coordinates. For this reason they cannot be found in molecular systems with less than five internal coordinates. Even in systems with enough degrees of freedom the dimensionality makes their location difficult. Until recently the perception was that three-state degeneracies are extremely rare, mostly due to the inability to locate them. Recent developments of efficient algorithms, however, have made these calculations possible and have revealed that three-state conical intersections exist in many systems [154, 205–207, 207–209] and may play an important role on their spectroscopic or photophysical properties.

The involvement of three-state conical intersections in the photophysics and radiationless decay processes of the nucleobases has been investigated using MRCI and CASPT2 methods [154, 210]. Three-state conical intersections have been located for the pyrimidine base, uracil, and the purine base, adenine by Matsika using MRCI methods [210]. A diagram showing the energy of these three-state conical intersections relative to vertical excitations and other two-state conical intersections is shown in Figure 11-16. In uracil a three-state degeneracy between the $S_0$, $S_1$, and $S_2$ states has been located 6.2 eV above the ground state minimum energy. This energy is 0.4 eV higher than vertical excitation to $S_2$ and at least 1.3 eV higher than the two-state conical intersections found previously. In adenine two different three-state degeneracies between the $S_1$, $S_2$, and $S_3$ states have been located at energies close to the vertical excitation energies. The energetics of these three-state conical intersections suggest they can play a role in a radiationless decay pathway present in adenine. The existence of two different seams of three-state conical intersections indicates that these features are common and complicate the potential energy surfaces of adenine and possibly many other aromatic molecules. CASSCF combined with CASPT2 methods have been used to locate a three-state conical intersection in cytosine [154]. This crossing between $S_0/\pi\pi^*/n_o\pi^*$ states has an energy 0.6 eV above the $S_1$ minimum at the CASSCF level but 2.42 eV at the CASPT2 level. So CASPT2 predicts the three-state conical intersection is too high to be important to any photophysical properties in cytosine. It is expected that three-state conical intersections will be present in other nucleobases and similar organic molecules. It remains to be seen how their presence will affect dynamics.

(a) Uracil        (b) Adenine



*Figure 11-16.* Diagram of the energy levels at the two- and three-state conical intersections in uracil and adenine calculated at the MRCI level. ci$IJK$ represents conical intersection between states $S_I$, $S_J$, $S_K$. (From Ref. [210])

### 11.3.5. Solvent Effects on Excited States of DNA Bases

Most theoretical calculations of nucleobases have been done in vacuo because of the extra computational difficulties associated with modeling the surrounding environment. This is a good first step to understand the fundamental excited state properties, but in order to learn more about the bases in their natural environment one has to go beyond that. Most experimental spectroscopic work has been carried out in solution and these experiments need theoretical results in the same medium to be interpreted correctly. There is great interest in understanding the effect of the environment on the photophysical properties of nucleobases, and there have been a number of studies that have included solvent effects in the calculation of excited states and considered solvatochromic shifts. Studies on the effect of solvent on the photophysical pathways discussed above are a lot less common. The methods used range from the polarizable continuum model (PCM) to mixed quantum mechanics/molecular mechanics methods (QM/MM). What is found in general terms is that the $n\pi^*$ states are usually blue-shifted in polar solvents, such as water, while the $\pi\pi^*$ states have much smaller shifts which can be to the red or to the blue ends of the spectrum. The shift in $n\pi^*$ states is mainly electrostatic. The dipole moment of these states is usually smaller than that of the ground state since excitation involves moving charge from a localized lone pair to a delocalized orbital. When the molecule is surrounded by a polar solvent the general rule is that the larger the dipole moment of a state is,

the larger is the stabilization of that state. So, if the dipole moment of the ground state is larger than the excited state, the ground state will become more stable, and the final results will be an increased gap between the states, i.e. a blue shift. Using the same argument a red shift is expected when the dipole moment of the ground state is smaller than that of the excited state. Other effects, besides electrostatic, are also present which are more difficult to model. These are more important in nonpolar solvents.

### 11.3.5.1.    Solvatochromic Shifts

Tables 11-6, 11-7, and 11-8 show calculated solvatochromic shifts for the nucleobases. Solvation effects on uracil have been studied theoretically in the past using both explicit and implicit models [92, 94, 130, 149, 211–214] (see Table 11-6). Initial studies used clusters of uracil with a few water molecules. Marian et al. [130] calculated excited states of uracil and uracil-water clusters with two, four and six water molecules. Shukla and Lesczynski [122] studied uracil with three water molecules using CIS to calculate excitation energies. Improta et al. [213] used a cluster of four water molecules embedded into a PCM and TDDFT calculations to study the solvatochromic shifts on the absorption and emission of uracil and thymine. Zazza et al. [211] used the perturbed matrix method (PMM) in combination with TDDFT and CCSD to calculate the solvatochromic shifts. The shift for the $S_1$ state ranges between $(+0.21) - (+0.54)$ eV and the shift for the $S_2$ is calculated to be between $(-0.07) - (-0.19)$ eV. Thymine shows very similar solvatochromic shifts as seen in Table 11-6 [92].

*Table 11-6.* Vertical excitation energy shifts in eV for uracil and thymine in aqueous phase compared to gas phase, using various levels of theory. The ordering of the states is according to the gas phase energies

| Level of theory | $S_1(n\pi^*)$ | $S_2(\pi\pi^*)$ |
|---|---|---|
| **Uracil** | | |
| Scaled CIS $+3H_2O$ [122] | +0.20 | −0.07 |
| PMM/TDDFT [211] | +0.38 | −0.18 |
| PMM/TD-PBE0 [211] | +0.54 | −0.10 |
| PMM/CCSD [211] | +0.34 | −0.12 |
| SCRF/CCSD [211] | +0.21 | −0.07 |
| TDDFT-PCM [213] | +0.29 | −0.09 |
| TDDFT-PCM $+4H_2O$ [213] | +0.48 | −0.10 |
| INDO-CIS $+200H_2O$ [214] | +0.50 | −0.19 |
| **Thymine** | | |
| TDDFT-PCM [92] | +0.29 | −0.08 |
| TDDFT-PCM $+4H_2O$ [92] | +0.41 | −0.09 |

*Table 11-7.* Vertical excitation energy shifts in eV for cytosine in aqueous phase compared to gas phase using various levels of theory. The ordering of the states is according to the gas phase energies

| Level of theory | $S_1(\pi\pi^*)$ | $S_2(n\pi^*)$ | $S_3(n\pi^*)$ |
|---|---|---|---|
| TDDFT [215]+$3H_2O$ | +0.18 | +0.45 | +0.43 |
| Scaled CIS+$3H_2O$ [215] | +0.10 | +0.35 | +0.31 |
| EOM-CCSD/MM [263][a] | +0.25 | +0.57 | |
| CR-EOM-CCSD(T)/MM [263][a] | +0.25 | +0.54 | |
| CPCM-CASPT2 [216] | +0.2 | +0.6 | +0.8 |

[a] Cytosine in native DNA environment.

*Table 11-8.* Vertical excitation energy shifts (in eV) for adenine in aqueous phase compared to gas phase using various levels of theory. The ordering of the states is according to the gas phase energies

| Level of theory | $S_1$ | $S_2$ | $S_3$ |
|---|---|---|---|
| **Adenine** | | | |
| CIS-SCRF [217] | −0.03[a] | +0.01[a] | +0.40[b] |
| CIPSI-IEF-PCM [219] | +0.07[b] | −0.36[a] | −0.36[a] |
| TDDFT-PCM [219] | +0.14[b] | −0.06[a] | −0.06[a] |
| LRFE-MRMP [223] | −0.01[a] | +0.11[a] | +0.04[b] |
| **Guanine** | | | |
| CIS-SCRF [218] | +0.07[a] | +0.71[b] | −0.06[a] |
| CIPSI-IEF-PCM [220] | −0.35[a] | −0.70[a] | |
| CASPT2-SCRF [143] | −0.03[a] | +0.02[a] | |

[a] $\pi\pi^*$ state.
[b] $n\pi^*$ state.

Solvatochromic shifts for cytosine have also been calculated with a variety of methods (see Table 11-7). Shukla and Lesczynski [215] studied clusters of cytosine and three water molecules with CIS and TDDFT methods to obtain solvatochromic shifts. More sophisticated calculations have appeared recently. Valiev and Kowalski used a coupled cluster and classical molecular dynamics approach to calculate the solvatochromic shifts of the excited states of cytosine in the native DNA environment. Blancafort and coworkers [216] used a CASPT2 approach combined with the conductor version of the polarizable continuous (CPCM) model. All of these methods predict that the first three excited states are blue-shifted. $S_1$, which is a $\pi\pi^*$ state, is blue-shifted by 0.1–0.2 eV in water and 0.25 eV in native DNA. $S_2$ and $S_3$ are both $n\pi^*$ states and, as expected, the shift is bigger, 0.4–0.6 eV for $S_2$ and 0.3–0.8 eV for $S_3$. $S_2$ is predicted to be blue-shifted by 0.54 eV in native DNA.

Shukla and coworkers have studied the excited states of purine bases, adenine and guanine, in water using CIS with the self-consistent reaction field (SCRF) to model the water [217, 218]. Tomasi and coworkers have also studied the purine bases

using an integral equation formalism PCM (IEF-PCM) model for solvation. The multireference perturbation configuration interaction (CIPSI) and TDDFT methods were used for the calculation of the excited states [219, 220]. Explicit clusters of guanine with a few water molecules have also been used to study the first excited state and proton transfer [221, 222]. Perturbation theory methods have also been applied to these systems. The linear-response free energy (LRFE) method combined with MRMP has also been used for adenine [223]. CASPT2 combined with SCRF has been used for guanine [143]. The solvatochromic shifts obtained vary widely as can be seen in Table 11-8. The variations in the purines are bigger than those in the pyrimidine bases, and there is no agreement even on the sign of solvatochromic shift. Shifts in the $\pi\pi^*$ states of adenine vary between $+0.11$ to $-0.36$ eV, while the $n\pi^*$ states shift by $+0.04$ to $+0.40$ eV. In guanine the first $\pi\pi^*$ state has a shift predicted to be between ($+0.07$) to ($-0.35$) and on the second $\pi\pi^*$ state between ($+0.02$) and ($-0.70$). The $n\pi^*$ state has only been calculated with the CIS-SCRF method to be blue-shifted by $+0.71$ eV. In general, there seems to be little consensus, and it is difficult to estimate which of these values are more accurate.

### 11.3.5.2.    Solvent Effects on Photophysics

These studies discuss vertical and adiabatic excitation energies but the photophysical behavior requires calculations along the PES and at highly distorted geometries, which are more difficult to carry out in the presence of solvent. Some theoretical work has been done in this area, but it is quite limited.

The effect of solvation on uracil and thymine photophysics has been studied by Gustavvson and coworkers, who have studied uracil with four explicit water molecules and PCM to study distorted geometries [92, 93, 149]. The conical intersection connecting $S_1$ to the ground state that was found in the gas phase is also present in solution. The barrier connecting the $S_1$ minimum to the conical intersection is lower in solution, however, causing much shorter lifetimes. So the nanosecond lifetime which is observed in the gas phase is not observed in solution but a picosecond lifetime is observed.

The effect of a hydrogen bonded water molecule on the relaxation pathways of uracil has also been examined by Yoshikawa and Matsika [224]. These results are mainly related to the experimental work carried out by Kong and coworkers on the clusters of water with methylated uracils [101, 102]. The electronically excited singlet states of the four most stable complexes of uracil with one water molecule have been studied theoretically using MRCI methods. Four cyclic isomers of uracil forming two hydrogen bonds with the water molecule have been located with energies within 0.2 eV from the lowest energy isomer. It was examined how complexation of uracil with water affects previously reported radiationless decay pathways to the ground state after initial excitation to the bright $S_2$ state [224]. The features on the excited state potential energy surfaces found previously for free uracil exist also in the presence of water. The summary of the pathways found in free uracil is shown in Figure 11-5. Their energies, however, are shifted because of the hydrogen bonds

formed, and, since the shifts are not the same for all stationary points and conical intersections, an effect on the excited state dynamics and lifetimes is expected. The first step involving radiationless decay from $S_2$ to $S_1$ is not affected strongly by the water molecule. An experimentally observed effect is seen in the shorter lifetimes for the dark $S_1$ state in methylated uracil upon complexation with one or more water molecules. According to the calculations, the $S_1$ lifetime depends on the barrier that connects the dark state minimum to the conical intersection with the ground state. This barrier is lower for some of the isomers compared to free uracil leading to shorter $S_1$ lifetimes, in agreement with the experimental observations. Interconversion between the different isomers [225] is expected to occur faster than accessibility to the barrier, so the lowest barrier and fastest radiationless decay will be the kinetically favorable path, and the shortest lifetime will be observed.

The photophysical pathways in aqueous cytosine have been recently studied by Blancafort and Migani [216]. The effect of water on the two pathways shown in Figure 11-7 has been investigated. The effect of water on the lowest-energy path through the ethylenic-like conical intersection (see Figure 11-6c) is a lowering of the barrier. In contrast, the second path is destabilized by hydrogen bonding, although the bulk solvent effect reduces the destabilization. The barrier in water for this pathway is ca. 0.3 eV, which may still be accessible. Overall, both pathways are present in solution, with one of them having a smaller barrier than the other.

Tomasi and coworkers, in addition to calculating solvatochromic shifts, have examined the excited state PES of adenine and 2AP for finding mechanisms for radiationless relaxation in solution, using CIPSI and IEF-PCM [219]. They explored the proximity effect of interaction between the $n\pi^*$ and $\pi\pi^*$ states, and intramolecular twisting of the amino group. They concluded that the proximity effect seemed more plausible in explaining the differences in photophysical properties between adenine and 2-aminopurine in solution. Solvent effects on the conical intersections and relaxation pathways have also been studied for adenine in aqueous and acetonitrile solutions using the MRMP method [223]. The LRFE methods was used to locate the energy minima and conical intersections in solution. Based on these results the authors concluded that in solution the pathway involving the $L_a/S_0$ conical intersection will be operative, but the pathway involving the $n\pi^*/S_0$ conical intersection will be suppressed. The pathway involving the $\pi\sigma^*$ state will also have a higher barrier in solution. Further studies are needed in order to have a better understanding of the effect of solvent on the photophysical pathways of the nulceobases.

## 11.3.6. Beyond Monomers

In DNA or RNA the nucleobases are the building blocks and interact with each other via hydrogen bonding and $\pi$-stacking. These interactions change the energetics of the excited states and provide additional photophysical or photochemical channels. Understanding the mechanisms involved when the bases are stacked in a single strand or hydrogen bonded to form the double strand is considerably more

complicated and poses serious theoretical and experimental challenges. The polymers and their excited state dynamics are currently under investigation, with several questions in debate, among which the nature of excited states and the relative importance of base stacking and base pairing for the relaxation of excited states in DNA [119, 120, 226–230]. The nature of the excited states is being discussed in terms of excimers, which here are loosely bound excited-state dimers of bases, and excitons, which are excited states delocalized among more than one base. One of the most important differences between the monomers and polymers is the appearance of long-lived emissive states in polymers that are not present in monomers [3, 6].

Although the main topic of this review is the photophysics of single bases, a brief summary of some recent theoretical investigations for systems beyond the monomers is given here. Progress in this area is limited compared to what has been done computationally for the monomers but there is a great interest in this topic and intense future research is expected to follow.

It has been shown that photoexcitation of the guanine-cytosine (G-C) base pair leads to proton transfer [231]. Watson-Crick (WC) base pairs have excited state lifetimes much shorter than other non-WC base pairs indicating once again that the natural occurring WC base pairs are more photostable than other alternative configurations [115, 118, 232–235]. Much work has been done in the gas phase where many different base pair isomers exist. The ultrafast relaxation of the WC base pair has also been confirmed in solution using fluorescence up-conversion measurements [117].

Several ab initio calculations have investigated the influence of base-pairing on the excited state dynamics [234–243]. Electronic structure calculations, as well as dynamic simulations have been carried out for G-C in the gas phase or embedded in DNA [235, 237, 241–243]. These ab initio studies suggest that following photoexcitation to a $^1\pi\pi^*$ state a proton transfer reaction from guanine to cytosine occurs. The proton transfer involves mainly the azine hydrogen of guanine moving to the $N^3$ nitrogen of cytosine. This proton transfer is followed by an efficient radiationless decay of the excited state via a conical intersection seam to the ground state.

Several base pairs of adenine–thymine, including the WC pair, have also been studied [238]. It is found that a charge transfer state exists about 1.5 eV higher than the local $\pi\pi^*$ states. Proton transfer between the bases stabilizes a charge transfer state which then crosses with the ground state facilitating radiationless relaxation. This mechanism is not energetically favorable for non WC pairs.

Bases stacked rather than hydrogen bonded have also been studied with quantum chemical methods [182, 244–247]. The nature of excited states in these systems has been debated and theoretical calculations are called to decide on the degree of excited state localization or delocalization, as well as the presence and energy of charge transfer states. The experimentally observed hypochromism of DNA compared to its individual bases has been known for decades [248]. Accurate quantum chemical calculations are limited in these systems because of their increased size. Many of the reported studies have used TDDFT to calculate excited states of bases stacked with other bases [182, 244, 246, 247]. However, one has to be cautious when us-

ing TDDFT especially when the states of interest are charge transfer states, since it is well known that TDDFT cannot describe these states accurately. Hardman and Thompson used CIS for stacked complexes similar to those studied by TDDFT and did not find charge transfer states at low energies as was reported before [182, 244]. A CASPT2 study on a cytosine $\pi$-stacked dimer has also be done focusing on the formation of excimer on the excited state and its fluorescence [245]. Clearly higher level calculations are needed for stacked bases.

Model systems containing more than three stacked bases are even more challenging to study. Examination of excitonic states delocalized over several bases would require enormous computer resources with the usual ab initio methods. However, some researchers are studying such systems within the Frenkel exciton model [249]. Markovitsi and coworkers have contributed considerably in this theoretical area [227, 250–252], and recently they have also studied these systems experimentally [253–255]. In their theoretical treatment, a set of $N$ stacked bases are first assigned a set of $N$ energies, where each energy term is a linear combination of one isolated excited base energy and $N-1$ isolated ground state energies, which are calculated in the usual way for isolated bases, with CIS or CASSCF. An $N \times N$ Hamiltonian matrix is then constructed with these energies as diagonal elements. The off-diagonal elements are treated perturbatively as interactions of transition charge densities, which can be approximated at different levels as point dipoles or extended dipoles. The matrix is then diagonalized, with resulting eigenstates given as linear combinations of the eigenstates of the individually excited base basis set. This method, however, has limitations from the inherent complete neglect of orbital overlap between bases, which would be included in a full ab initio calculation of the supersystem, and as such it is unable to reproduce some experimentally observed phenomena, such as the hypochromism of base oligomers [250]. On the other hand, because of the relatively inexpensive computation involved with this method, this model lends itself well to the study of the excitation effects of dynamic fluctuations and solvation of the strands, by combining it with molecular dynamics, and studies can be carried out on current workstations.

Bittner and coworkers have recently offered a similar theoretical method for studying the excited states in oligomers within the Frenkel exciton model, where the off-diagonal terms are calculated using a transition density cube method [256], while diagonal elements were calculated using TDDFT. Results of this method compared with those that treat the off-diagonal elements using the ideal dipole approximation (IDA) showed that IDA strongly overestimates the coupling contributions. Bittner has also studied the dynamics of excitons charge-transfer in rigid oligomers of poly(dA)poly(dT) modeled as a lattice, with theoretical framework borrowed from lattice fermion models of quantum chromodynamics [257]. It was found that the charge-transfer states are low-lying dark states, followed by delocalized excitonic states with only weak mixing between the adenosine side and the thymidine side, underscoring the importance of the competition between charge-transfer and energy delocalization.

### 11.3.7. Photochemistry

Finally a few sentences are deserved for the vast area of DNA photochemistry. Thymine dimerization is the most common photochemical reaction with the quantum yield of formation in isolated DNA of all-thymine oligodeoxynucleotides 2–3% [3]. Furthermore, a recent study based on femtosecond time-resolved transient absorption spectroscopy showed that thymine dimers are formed in less than 1 ps when the strand has an appropriate conformation [258]. The low quantum yield of the reaction in regular DNA is suggested to be due to the infrequency of these appropriate reactive conformations.

The measured ultrafast nature of the reaction suggests rapid transitions through conical intersections and this suggestion has inspired some recent calculations which investigate the mechanism for the thymine photodimerization. Initial theoretical investigations used DFT and extracted information about the excited state based on the ground state PES [259, 260]. More recently the CASSCF method was used to study the involvement of conical intersections in the mechanism. Independently two groups reported the presence of conical intersections between the ground and excited state that can lead to the photoproduct [261, 262]. These studies represent the first steps into accurate quantum mechanical calculations of the photochemistry in DNA and are expected to grow rapidly.

### 11.4. CONCLUSIONS

Great advances have occurred during recent years in understanding the photophysical behavior of nucleobases based on femtosecond spectroscopy and high level ab initio calculations. The quantum chemical methods can provide insight into the photophysical and photochemical processes of chemical and biological systems. Current electronic structure methods can predict reasonably accurate excited states and PES of relatively small systems, such as the nucleobases. More efficient methods are needed however for expanding these studies to larger systems, especially in view of the potential applications in photobiology. Excited states of DNA monomers are now believed to be understood to a large degree, and efforts are moving now to the next challenge of being able to understand multimers.

### ACKNOWLEDGMENTS

### REFERENCES

1. Kraemer KH (1997) Sunlight and skin cancer: another link revealed. Proc Natl Acad Sci USA 94:11
2. Mukhtar H, Elmets CA (1996) Photocarcinogenesis: mechanisms, models and human health implications. Photochem Photobiol 63:355

3. Cadet J, Vigny P (1990) In: Morrison H (ed.) Bioorganic photochemistry. John Wiley, New York, pp. 1–272

4. Daniels M (1976) In: Wang SY (ed.) Photochemistry and Photobiology of Nucleic Acids, vol 1. Academic Press, New York, p 23

5. Ruzsicska BP, Lemaire DGE In: Horspool WM, Song P-S (eds.) CRC Handbook of Organic Photochemistry and Photobiology, CRC Press, Boca Raton, FL, p 1289

6. Crespo-Hernandez CE, Cohen B, Hare PM, Kohler B (2004) Ultrafast excited-state dynamics in nucleic acids. Chem Rev 104:1977

7. Born M, Oppenheimer R (1927) Zur Quantentheorie der Molekeln. Ann Phys 84:457

8. Landau LD (1932) Z. Sowjetunion U.R.S.S. 2:46

9. Zener C (1932) Non-adiabatic crossing of energy levels. Proc Roy Soc A137:696

10. Domcke W, Yarkony DR, Köppel H (2004) Conical intersections. World Scientific, Singapore,

11. Rice OK (1929) Phys Rev 33:748

12. Dantus M, Zewail A (2004) Introduction: femtochemistry. Chem Rev 104:1717

13. von Neumann J, Wigner EP (1929) On the behaviour of eigenvalues in adiabatic processes. Physik Z 30:467

14. Teller E (1937) The crossing of potential surfaces. J Phys Chem 41:109

15. Manaa MR, Yarkony DR (1993) On the intersection of two potential energy surfaces of the same symmetry. Systematic characterization using a lagrange multiplier constrained procedure. J Chem Phys 99:5251

16. Bearpark MJ, Robb MA, Schlegel HB (1994) A direct method for the location of the lowest energy point on a potential surface crossing. Chem Phys Lett 223:269

17. Yarkony DR (1998) Conical intersections: diabolical and often misunderstood. Acc Chem Res 31:511

18. Yarkony DR (1996) Diabolical conical intersections. Rev Mod Phys 68:985

19. Yarkony DR (1996) Current issues in nonadiabatic chemistry. J Phys Chem 100:18612

20. Yarkony DR (2001) Conical intersections: the new conventional wisdom. J Phys Chem A 105:6277

21. Bernardi F, Olivucci M, Robb MA (1990) Predicting forbidden and allowed cycloaddition reactions: potential surface topology and its rationalization. Acc Chem Res 23:405

22. Bernardi F, Olivucci M, Robb MA (1996) Potential energy surface crossings in organic photochemistry. Chem Soc Rev 25:321

23. Barckholtz TA, Miller TA (1998) Quantitative insights about molecules exhibiting jahn-teller and related effects. Int Rev Phys Chem 17:435

24. Robb MA, Garavelli M, Olivucci M, Bernardi F (2000) A computational strategy for organic photochemistry. In Lipkowitz KB Boyd DB (eds) Reviews in computational chemistry, vol. 15 Wiley-VCH, New York; pp 87–146

25. Matsika S (2007) Conical intersections in molecular systems. In: Lipkowitz KB and Cundari TR (eds.) Reviews in computational chemistry, vol. 23, Wiley-VCH, New Jersey, pp 83–124

26. Klessinger M, Michl J (1995) Excited states and photochemistry of organic molecules. VCH Publishers, Inc., New York

27. Michl J, Bonačić Koutecký V (1990) Electronic aspects of organic photochemistry. Wiley Interscience, New York

28. Garavelli M, Gelani P, Bernardi F, Robb MA, Olivucci M (1997) The $C_5H_6NH_2^+$ protonated shiff base: an ab initio minimal model for retinal photoisomerization. J Am Chem Soc 119:6891

29. Kobayashi T, Saito T, Ohtani H (2001) Real-time spectroscopy of transition states in bacteriorhodopsin during retinal isomerization. Nature 414:531

30. Ben-Nun M, Molnar F, Schulten K, Martinez TJ (2000) The role of intersection topography in bond selectivity of cis-trans photoisomerization. Proc Natl Acad Sci USA 97:9379

31. Warshel A, Chu ZT (2001) Nature of the surface crossing process in bacteriorhodopsin: computer simulations of the quantum dynamics of the primary photochemical event. J Phys Chem B 105:9857

32. Migani A, Sinicropi A, Ferr N, Cembran A, Garavelli M, Olivucci M (2004) Structure of the intersection space associated with Z/E photoisomerization of retinal in rhodopsin proteins. Faraday discuss 127:179

33. Toniolo A, Olsen S, Manohar L, Martinez TJ (2004) Conical intersection dynamics in solution: the chromophore of green fluorescent protein. Faraday Discuss 127:149

34. Martin ME, Negri F, Olivucci M (2004) Origin, nature, and fate of the fluorescent state of the green fluorescent protein chromophore at the CASPT2//CASSCF resolution. J Am Chem Soc 126:5452

35. Worth GA, Cederbaum LS (2001) Mediation of ultrafast transfer in biological systems by conical intersections. Chem Phys Lett 338:219

36. Atchity GJ, Xantheas SS, Ruedenberg K (1991) Potential energy surfaces near intersections. J Chem Phys 95:1862

37. Köppel H, Domcke W, Cederbaum LS (1984) Multimode molecular dynamics beyond the born-oppenheimer approximation. Adv Chem Phys 57:59

38. Domcke W, Stock G (1997) Theory of ultrafast nonadiabatic excited-state processes and their spectroscopic detection in real time. Adv Chem Phys 100:1–170

39. Yarkony DR (2001) Nuclear dynamics near conical intersections in the adiabatic representation. I. The effects of local topography on interstate transition. J Chem Phys 114:2601

40. Schlegel HB (1995) Geometry optimization on potential energy surfaces, in: Yarkony DR (ed) modern electronic structure theory part I World Scientifc, Singapore p 459–500, Advance Series in Physical Chemistry.

41. Bartlett RJ, Stanton JF (1994) Applications of post-hartree-fock methods: a tutorial. In: Lipkowitz KB Boyd DB (eds) Reviews in computational chemistry, vol. 5. Wiley-VCH, New York, pp 65–169

42. Shavitt I (1977) The method of configuration interaction. In: Schaefer III HF (ed) Methods of Electronic Structure Theory, vol. 4 of Modern Theoretical Chemistry; Plenum Press, New York; pp 189–275

43. Roos BO, Taylor PR (1980) A complete active space SCF method (CASSCF) using a density-matrix formulated super-CI approach. Chem Phys 48:157

44. Roos BO (1987) The complete active space self consistent field method and its applications in electronic structure calculations. Adv Chem Phys 69:399

45. Shepard R (1987) Geometrical energy derivative evaluation with MRCI wave functions. Int J Quantum Chem 31:33

46. Shepard R (1995) The analytic gradient method for configuration interaction wave functions. Yarkony DR (ed) In: Modern electronic structure theory part I, World Scientific, Singapore, p 345

47. Lischka H, Dallos M, Shepard R (2002) Analytic MRCI gradient for excited states: formalism and application to the $n\pi^*$ valence- and $n-(3s, 3p)$ rydberg states of formaldehyde. Mol Phys 100:1647

48. Lengsfield BH, Yarkony DR (1992) Nonadiabatic interactions between potential energy surfaces: theory and applications. In Baer M, Ng CY (eds) State-selected and state-to-state ion-molecule reaction dynamics: part 2 theory, Vol. 82 of Advances in Chemical Physics, John Wiley and Sons, New York, p 1–71.

49. Lischka H, Shepard R, Pitzer RM, Shavitt I, Dallos M, Müller Th, Szalay PG, Seth M, Kedziora GS, Yabushita S, Zhang Z (2001) High-level multireference methods in the quantum-chemistry program system COLUMBUS: analytic MR-CISD and MR-AQCC gradients and MR-AQCC-LRT for excited states, GUGA spin-orbit CI and parallel CI density. Phys Chem Chem Phys 3:664

50. Lischka H, Dallos M, Szalay PG, Yarkony DR, Shepard R (2004) Analytic evaluation of nonadiabatic coupling terms at the MR-CI level. I. Formalism. J Chem Phys 120:7322

51. Dallos M, Lischka H, Shepard R, Yarkony DR, Szalay PG (2004) Analytic evaluation of nonadiabatic coupling terms at the MR-CI level. II. Minima on the crossing seam: formaldehyde and the photodimerization of ethylene. J Chem Phys 120:7330

52. Werner H-J, Knowles PJ (1988) An efficient internally contracted multiconfiguration reference CI method. J Chem Phys 89:5803

53. Werner HJ, Knowles PJ, Lindh R, Schütz M, Celani P, Korona T, Manby FR, Rauhut G, Amos RD, Bernhardsson A, Berning A, Cooper DL, Deegan MJO, Dobbyn AJ, Eckert F, Hampel C, Hetzer G, Lloyd AW, McNicholas SJ, Meyer W, Mura ME, Nicklass A, Palmieri P, Pitzer R, Schumann U, Stoll H, Stone AJ, Tarroni R, Thorsteinsson T (2003) Molpro, version 2002.6, a package of ab initio programs

54. Andersson K, Malmqvist PA, Roos BO, Sadlej AJ, Wolinski K (1990) Second-order perturbation-theory with a CASSCF reference function. J Phys Chem 94:5483

55. Anderson K, Malmqvist PA, Roos BO (1992) Second-order perturbation-theory with a complete active space self-consistent field reference function. J Chem Phys 96:1218

56. Karlström G, Lindh R, Malmqvist P-Å, Roos B, Ryde U, Veryazov V, Widmark P-O, Cossi M, Schimmelpfennig B, Neogrady P, Seijo L (2003) Molcas: a program package for computational chemistry. Comput material sci 28:222

57. Roos BO, Andersson K, Fulscher MP, Malmqvist PA, Serrano-Andrès L, Pierloot K, Merchán M (1996) Multiconfigurational perturbation theory: applications in electronic spectroscopy. In: Prigogine I Rice SA (eds) New methods in computational quantum mechanics, Vol. 93 of Advances in Chemical Physics, Wiley, New York, p 219

58. Nakano H (1993) Quasidegenerate perturbation theory with multiconfigurational self-consistent-field reference functions. J Chem Phys 99:7983

59. Finley J, Malmqvist PA, Roos BO, Serrano-Andrès L (1998) The multi-state CASPT2 method. Chem Phys Lett 288:299

60. Serrano-Andrès L, Merchán M, Lindh R (2005) Computation of conical intersections by using perturbation techniques. J Chem Phys 122:104107

61. Hirao K (1992) Multireference möller-plesset method. Chem Phys Lett 190:374

62. Glaesemann KR, Gordon MS, Nakano H (1999) A study of feCO$^+$ with correlated wavefunctions. Phys Chem Chem Phys 1:967

63. Schmidt MW, Baldridge KK, Boatz JA, Elbert ST, Gordon MS, Jensen JH, Koseki S, Matsunaga N, Nguyen KA, Su S, Windus TL, Dupuis M, Montgomery JA (1993) Computation of conical intersections by using perturbation techniques. J Comput Chem 14:1347

64. Nakano H (1993) Quasidegenerate perturbation theory with multiconfigurational self-consistent-field reference functions. J Chem Phys 99:7983–7992

65. Nakano H (1993) MCSCF reference quasidegenerate perturbation theory with EpsteinNesbet partitioning. Chem Phys Lett 207:372–378

66. Dreuw A, Head-Gordon M (2005) Single-reference ab initio methods for the calculation of excited states of large molecules. Chem Rev 105:4009–4037

67. Head-Gordon M, Oumi M, Maurice D (1999) Quasidegenerate second-order perturbation corrections to single-excitation configuration interaction. Mol Phys 96:593

68. Laikov D, Matsika S (2007) Inclusion of second-order correlation effects for the ground and singly-excited states suitable for the study of conical intersections: the CIS(2) model. Chem Phys Lett 448:132–137

69. Koch H, Jensen HJA, Jorgensen P, Helgaker T (1990) Excitation-energies from the coupled cluster singles and doubles linear response function (CCSDLR) – applications to be, CH$^+$, CO, and H$_2$O. J Chem Phys 93:3345

70. Stanton JF, Bartlett RJ (1993) The equation of motion coupled-cluster method - a systematic biorthogonal approach to molecular-excitation energies, transition-probabilities, and excited-state properties. J Chem Phys 98:7029

71. Krylov AI (2001) Size-consistent wave functions for bond-breaking: the equation-of-motion spin-flip model. Chem Phys Lett 338:375

72. Krylov AI (2006) Spin-flip equation-of-motion coupled-cluster electronic structure method for a description of excited states, bond breaking, diradicals, and triradicals. Acc Chem Res 39:83–91

73. Nakatsuji H, Hirao K (1978) Cluster expansion of the wavefunction. symmetry-adapted-cluster expansion, its variational determination, and extension of open-shell orbital theory. J Chem Phys 68:2053

74. Christiansen O, Koch H, Jorgensen P (1995) The second-order approximate coupled cluster singles and doubles model CC2. Chem Phys Lett 243:409–418

75. Hättig C, Weigend F (2000) CC2 excitation energy calculations on large molecules using the resolution of the identity approximation. J Chem Phys 113:5154

76. Ahlrichs R, Bär M, Häser M, Horn H, Kölmel C (1989) Electronic structure calculations on workstation computers: The program system turbomole. Chem Phys Lett 162:165–169

77. Bickelhaupt FM, Baerends EJ (2000) Kohn-sham density functional theory: predicting and understanding chemistry. In: Lipkowitz KB Boyd DB (eds) Reviews in computational chemistry, vol. 15. Wiley-VCH, New York, pp 1–86

78. Runge E, Gross EKU (1984) Density-functional theory for time-dependent systems. Phys Rev Lett 52:997

79. Levine BG, Ko C, Quenneville J, Martinez TJ (2006) Conical intersections and double excitations in time-dependent density functional theory. Mol Phys 104:1039

80. Grimme S, Waletzke M (1999) A combination of KohnSham density functional theory and multi-reference configuration interaction methods. J Chem Phys 111:5645–5655

81. Daniels M, Hauswirth W (1971) Fluorescence of the purine and pyrimidine bases of the nucleic acids in neutral aqueous solution at 300 K. Science 171:675

82. Callis PR (1983) Electronic states and luminescence of nucleic acid systems. Ann Rev Phys Chem 34:329

83. Pecourt J-ML, Peon J, Kohler B (2000) Ultrafast internal conversion of electronically excited RNA and DNA nucleosides in water. J Am Chem Soc 122:9348

84. Pecourt J-ML, Peon J, Kohler B (2001) DNA excited-state dynamics: ultrafast internal conversion and vibrational cooling in a series of nucleosides. J Am Chem Soc 123:10370

85. Cohen B, Hare P, Kohler B (2003) Ultrafast excited-state dynamics of adenine and monomethylated adenines in solution: implications for the nonradiative decay mechanism. J Am Chem Soc 125:13594

86. Cohen B, Crespo-Hernandez CE, Kohler B (2004) Strickler-berg analysis of excited singlet state dynamis in DNA and RNA nucleosides. Faraday Discuss 127:137

87. Peon J, Zewail AH (2001) DNA/RNA nucleotides and nucleosides: direct measurement of excited-state lifetimes by femtosecond fluorescence up-conversion. Chem Phys Lett 348:255

88. Gustavsson T, Sharonov A, Markovitsi D (2002) Thymine, thymidine and thymidine 5'-monophosphate studied by femtosecond fluorescence upconversion spectroscopy. Chem Phys Lett 351:195

89. Gustavsson T, Sharonov A, Onidas D, Markovitsi D (2002) Adenine, deoxyadenosine and deoxyadenosine 5'-monophosphate studied by femtosecond fluorescence upconversion spectroscopy. Chem Phys Lett 356:49

90. Onidas D, Markovitsi D, Marguet S, Sharonov A, Gustavsson T (2002) Fluorescence properties of DNA nucleosides and nucleotides: a refined steady-state and femtosecond investigation. J Phys Chem B 106:11367

91. Sharonov A, Gustavsson T, Carré V, Renault E, Markovitsi D (2003) Cytosine excited state dynamics studied by femtosecond absorption and fluorescence spectroscopy. Chem Phys Lett 380:173–180

92. Gustavsson T, Banyasz A, Lazzarotto E, Markovitsi D, Scalmani G, Frisch MJ, Barone V, Improta R (2006) Singlet excited-state behavior of uracil and thymine in aqueous solution: a combined experimental and computational study of 11 uracil derivatives. J Am Chem Soc 128:607–619

93. Gustavsson T, Sarkar N, Lazzarotto E, Markovitsi D, Barone V, Improta R (2006) Solvent effect on the singlet excited-state dynamics of 5-fluorouracil in acetonitrile as compared with water. J Phys Chem B 110:12843–12847

94. Gustavsson T, Sarkar N, Lazzarotto E, Markovitsi D, Improta R (2006) Singlet excited-state dynamics of uracil and thymine derivatives: a femtosecond fluorescence upconversion study in acetonitrile. Chem Phys Lett 429:551–557

95. Brady BB, Peteanu L, Levy DH (1988) The electronic spectra of the pyrimidine bases uracil and thymine in a supersonic molecular beam. Chem Phys Lett 147:538

96. Nir E, Kleinermanns K, Grace L, de Vries MS (2001) On the photochemistry of purine nucleobases. J Phys Chem A 105:5106

97. Nir E, Müller M, Grace LI, de Vries MS (2002) REMPI spectroscopy of cytosine Chem Phys Lett 355:59

98. Nir E, Kleinermanns IHK, de Vries MS (2003) The nucleobase cytosine and the cytosine dimer investigated by double resonance laser spectroscopy and ab initio calculations. Phys Chem Chem Phys 5:4780

99. Piuzzi F, Mons M, Dimicoli I, Tardivel B, Zhao Q (2001) Ultraviolet spectroscopy and tautomerism of the DNA base quanine and its hydrate formed in a supersonic jet. Chem Phys 270:205

100. Lührs DC, Viallon J, Fischer I (2001) Excited state spectroscopy and dynamics of isolated adenine and 9-methyladenine. Phys Chem Chem Phys 3:1827

101. He Y, Wu C, Kong W (2003) Decay pathways of thymine and methyl-substituted uracil and thymine in the gas phase. J Phys Chem A 107:5145

102. He Y, Wu C, Kong W (2004) Photophysics of methyl-substituted uracils and thymines and their water complexes in the gas phase. J Phys Chem A 108:943

103. Kang H, Lee KT, Jung B, Ko YJ, Kim SK (2002) Intrinsic lifetimes of the excited state of DNA and RNA bases. J Am Chem Soc 124:12958

104. Kang H, Jung B, Kim SK (2003) Mechanism for ultrafast internal conversion of adenine, J Chem Phys 118:6717

105. Kim NJ, Jeong G, Kim YS, Sung J, Kim SK, Park YD (2000) Resonant two-photon ionization and laser induced fluorescence spectroscopy of jet-cooled adenine. J Chem Phys 113:10051

106. Ullrich S, Schultz T, Zgierski MZ, Stolow A (2004) Direct observation of electronic relaxation dynamics in adenine via time-resolved photoelectron spectroscopy. J Am Chem Soc 126:2262

107. Ullrich S, Schultz T, Zgierski MZ, Stolow A (2004) Electronic relaxation dynamics in DNA and RNA bases studied by time-resolved photoelectron spectroscopy. Phys Chem Chem Phys 6:2796

108. Satzger H, Townsend D, Zgierski MZ, Patchkovskii S, Ullrich S, Stolow A (2006) Primary processes underlying the photostability of isolated DNA bases: adenine. Proc Natl Acad Sci USA 103: 10196–10201

109. Nir E, Grace LI, Brauer B, de Vries MS (1999) REMPI spectroscopy of jet-cooled guanine. J Am Chem Soc 121:4896–4897

110. Crews B, Abo-Riziq A, Grace LI, Callahan M, Kabelác M, Hobza P, de Vries MS (2005) IR-UV double reonance spectroscopy of guanine-H$_2$O clusters. Phys Chem Chem Phys 7:3015–3020

111. Abo-Riziq A, Crews BO, Compagnon I, Oomens J, Meijer G, Helden GV, Kabelác M, Hobza P, de Vries MS (2005) The mid-IR spectra of 9-ehtyl guanine, guanosine, and 2-deoxyguanosine. Phys Chem Chem Phys 7:3015–3020

112. Canuel C, Mons M, Piuzzi F, Tardinel B, Dimicoli I, Elhanine M (2005) Excited states dynamics of DNA and RNA bases: characterization of a stepwise deactivation pathway in the gas phase. J Chem Phys 122:074316

113. Mons M, Piuzzi F, Dimicoli I, Gorb L, Lesczynski J (2006) Near-UV resonant two-photon ionization spectroscopy of gas phase guanine: evidence for the observation of three rare tautomers. J Phys Chem A 110:10921–10924

114. Ritze H-H, Lippert H, Samoylova E, Smith VR, Hertel IV, Radloff W, Schultz T (2005) Relevance of $\pi\sigma^*$ states in the photoinduced processes of adenine, adenine dimer, and adeninewater complexes. J Chem Phys 122:224320

115. Samoylova E, Lippert H, Ullrich S, Hertel IV, Radloff W, Schultz T, (2005) Dynamics of photoinduced processes in adenine and thymine base pairs. J Am Chem Soc 127:1782–1786

116. Hare PM, Crespo-Hernandez CE, Kohler B (2007) Internal conversion to the electronic ground state occurs via two distinct pathways for the pyrimidine bases in aqueous solution. Proc Natl Acad Sci USA 104:435–440

117. Schwalb NK, Temps F (2007) Ultrafast electronic relaxation in guanosine is promoted by hydrogen bonding with cytidine. J Am Chem Soc 129:9272

118. Abo-Riziq A, Grace L, Nir E, Kabelac M, Hobza P, de Vries MS (2005) Photochemical selectivity in guanine-cytosine base-pair structures. Proc Natl Acad Sci USA 102:20–23

119. Crespo-Hernandez CE, Cohen B, Kohler B (2005) Base stacking controls excited-state dynamics in A-T DNA. Nature 436:1141–1144

120. Markovitsi D, Talbot F, Gustavsson T, Onidas D, Lazzarotto E, Marguet S (2006) Molecular spectroscopy: complexity of excited state dynamics in DNA. Nature 442: E7

121. Broo A, Holmén A (1997) Calculations and characterization of the electronic spectra of DNA bases based on ab initio MP2 geometries of different tautomeric forms. J Phys Chem A 101:3589

122. Shukla MK, Leszczynski J (2002) Phototautomerism in uracil: a quantum chemical investigation. J Phys Chem A 106:8642

123. Shukla MK, Mishra PC (1999) A gas phase ab initio excited state geometry optimization study of thymine, cytosine and uracil. Chem Phys 240:319

124. Fleig T, Knecht S, Hättig C (2007) Quantum-chemical investigation of the structures and electronic spectra of the nucleic acid bases at the coupled cluster CC2 level. J Phys Chem A 111:5482–5491

125. Lorentzon J, Fülscher MP, Roos BO (1995) Theoretical study of the electronic spectrum of uracil and thymine. J Am Chem Soc 117:9265

126. Hudock HR, Levine BG, Thompson AL, Satzger H, Townsend D, Gador N, Ullrich S, Stolow A, Martinez TJ (2007) Ab initio molecular dynamics and time-resolved photoelectron spectroscopy of electronically excited uracil and thymine. J Phys Chem A 111:8500–8508

127. Petke JD, Maggiora GM, Christoffersen RE (1992) Ab-initio configuration interaction and random phase approximation caclulations of the excited singlet and triplet states of uracil and cytosine. J Phys Chem 96:6992

128. Zgierski MZ, Patchkovskii S, Fujiwara T, Lim EC (2005) On the origin of the ultrafast internal conversion of electronically excited pyrimidine bases. J Phys Chem A 109:9384–9387

129. Neiss C, Saalfrank P, Parac M, Grimme S (2003) Quantum chemical calculation of excited states of flavin-related molecules. J Phys Chem A 107:140

130. Marian CM, Schneidaer F, Kleinschmidt M, Tatchen J (2002) Electronic excitation and singlet-triplet coupling in uracil tatutomers and uracil-water complexes. Eur Phys J D 20:357

131. Shukla MK, Leszczynski J (2004) TDDFT investigation on nucleic acid bases: comparison with experiments and standard approach. J Comput Chem 25:768–778

132. Clark LB, Peschel GG, Tinoco I Jr (1965) Vapor spectra and heats of vaporization of some purine and pyrimidine bases. J Phys Chem 69:3615

133. Tomić K, Jörg T, Marian CM (2005) Quantum chemical investigation of the electronic spectra of the keto, enol, and keto-imine tautomers of cytosine. J Phys Chem A 109:8410–8418

134. Fülscher MP, Roos BO (1995) Theoretical study of the electronic spectrum of cytosine. J Am Chem Soc 117:2089

135. Clark LB, Tinoco I Jr (1965) Correlations in the ultraviolet spectra of the purine and pyrimidine bases. J Am Chem Soc 87:11

136. Voelter W, Records R, Bunnenberg E, Djerassi C (1968) Magnetic circular dichroism studies. vi. investigation of some purines, pyrimidines, and nucleosides. J Am Chem Soc 90:6163–6170

137. Platt JR (1949) Classification of spectra of cata-condensed hydrocarbons. J Chem Phys 17:484

138. Perun S, Sobolewski AL, Domcke W (2005) Ab initio studies on the radiationless decay mechanisms of the lowest excited singlet states of 9H-adenine. J Am Chem Soc 127:6257–6265

139. Marian CM (2005) A new pathway for the rapid decay of electronically excited adenine. J Chem Phys 122:104314

140. Merchán M, Serrano-Andrés L (2003) Ultrafast internal conversion of excited cytosine via the lowest $\pi\pi^*$ electronic singlet state. J Am Chem Soc 125:8108

141. Blancafort L (2006) Excited-state potential energy surface for the photophysics of adenine. J Am Chem Soc 128:210–219

142. Clark LB (1990) Electronic spectrum of the adenine chromophore. J Phys Chem 94:2873–2879

143. Fülscher MP, Serrano-Andres L, Roos BO (1997) A theoretical study of the electronic spectra of adenine and guanine. J Am Chem Soc 119:6168

144. Kistler KA, Matsika S (2007) Cytosine in context: a theoretical study of substituent effects on the excitation energies of 2-pyrimidinone derivatives. J Phys Chem A 111:8708–8716

145. Chen H, Li SH (2006) Ab initio study on deactivation pathways of excited 9H-guanine. J Chem Phys 124:154315

146. Clark LB (1977) Electronic spectra of crystalline 9-ethylguanine and guanine hydrochloride. J Am Chem Soc 99:3934

147. Matsika S (2004) Radiationless decay of excited states of uracil through conical intersections. J Phys Chem A 108:7584

148. Zgierski MZ, Fujiwara T, Kofron WG, Lim EC (2007) Highly effective quanching of the ultrafast radiationless decay of photoexcited pyrimidine bases by covalent modification: photophysics of 5,6-trimethylenecytosine and 5,6-trimethyleneuracil. Phys Chem Chem Phys 9:3206–3209

149. Santoro F, Barone V, Gustavsson T, Improta R (2006) Solvent effect on the singlet excited-state lifetimes of nucleic acid bases: a computational study of 5-fluorouracil and uracil in acetonitrile and water. J Am Chem Soc 128:16312–16322

150. Merchán M, Gonzalez-Luque R, Climent T, Serrano-Andrés L, Rodriguez E, Reguero M, Pelaez D (2006) Unified model for the ultrafast decay of pyrimidine nucleobases. J Phys Chem B 110:26471–26476

151. Hare PM, Crespo-Hernandez CE, Kohler B (2006) Solvent-dependent photophysics of 1-cyclohexyluracil: ultrafast branching in the initial bright state leads nonradiatively to the electronic ground state and a long-lived $^1n\pi^*$ state. J Phys Chem B 110:18641

152. Perun S, Sobolewski AL, Domcke W (2006) Conical intersections in thymine. J Phys Chem A 110:13238

153. Ismail N, Blancafort L, Olivucci M, Kohler B, Robb MA (2002) Ultrafast decay of electronically excited singlet cytosine via $\pi, \pi^*$ to $n, \pi^*$ state switch. J Am Chem Soc 124:6818

154. Blancafort L, Robb MA (2004) Key role of a threefold state crossing in the ultrafast decay of electronically excited cytosine. J Phys Chem A 108:10609

155. Zgierski MZ, Patchkovskii S, Lim EC (2005) Ab initio study of a biradical radiationless decay channel of the lowest excited electronic state of cytosine and its derivatives. J Chem Phys 123:081101

156. Blancafort L (2007) Energetics of cytosine singlet excited-state decay paths – A difficult case for CASSCF and CASPT2. Photochem Photobiol 83:603–610

157. Kistler KA, Matsika S (2007) Radiationless decay mechanism of cytosine: an ab initio study with comparisons to the fluorescent analogue 5-methyl-2-pyrimidinone. J Phys Chem A 111:2650–2661

158. Sobolewski AL, Domcke W (2002) On the mechanism of nonradiative decay of DNA bases: ab initio and TDDFT results for the excited states of 9H-adenine. Eur Phys J D 20:369

159. Sobolewski AL, Domcke W, Dedonder-Lardeux C, Jouvet C (2002) Excited-state hydrogen detachment and hydrogen transfer driven by repulsive (1)pi sigma* states: a new paradigm for nonradiative decay in aromatic biomolecules. J Phys Chem Chem Phys 4:1093–1100

160. Perun S, Sobolewski AL, Domcke W (2005) Photostability of 9h-adenine: mechanisms of the radiationless deactivation of the lowest excited singlet states. Chem Phys 313:107–112

161. Chung WC, Lan Z, Ohtsuki Y, Shimakura N, Domcke W, Fujimura Y (2007) Conical intersections involving the dissociative $^1\pi\sigma^*$ state in 9H-adenine: a quantum chemical ab initio study. Phys Chem Chem Phys 9:2075–2084

162. Nielsen SB, Solling TI (2005) Are conical intersections responsible for the ultrafast processes of adenine, protonated adenine, and the corresponding nucleosides?. Chem Phys Chem 6:1276

163. Serrano-Andrés L, Merchán M, Borin AC (2006) A three-state model for the photophysics of adenine. Chem Eur J 12:6559–6571

164. Serrano-Andrés L, Merchán M, Borin AC (2006) Adenine and 2-aminopurine: paradigms of modern theoretical photochemistry. Proc Natl Acad Sci USA 103:8691–8696

165. Chen H, Li SH (2005) Theoretical study toward understanding ultrafast internal conversion of excited 9H-adenine. J Phys Chem A 109:8443–8446

166. Zgierski MZ, Patchkovskii S, Lim EC (2007) Biradical radiationless decay channel in adenine and its derivatives. Can J Chem 85:124

167. Langer H, Doltsinis NL (2003) Selected photostabilisation of guanine by methylation. Phys Chem Chem Phys 5:4516–4518

168. Langer H, Doltsinis NL (2003) Excited state tautomerism of the DNA base guanine: a restricted open-shell kohn-sham study. J Chem Phys 118:5400–5407

169. Langer H, Doltsinis NL (2004) Nonradiative decay of photoexcited methylated guanine. Phys Chem Chem Phys 6:2742–2748

170. Chen H, Li SH (2006) Theoretical study on the excitation energies of six tautomers of guanine: evidence for the assignment of the rare tautomers. J Phys Chem A 110:12360–12362

171. Marian CM (2007) The guanine tautomer puzzle: quantum chemical investigation of ground and excited states. J Phys Chem A 111:1545–1553

172. Zgierski MZ, Patchkovskii S, Fujiwarab T, Lim EC (2007) The role of out-of-plane deformations in subpicosecond internal conversion of photoexcited purine bases: absence of the ultrafast decay channel in propanodeoxyguanosine. Chem Phys Lett 440:145–149

173. Laland SG, Serck-Hanssen G (1964) Synthesis of pyrimidin-2-one deoxyribosides and their ability to support the growth of the deoxyriboside-requiring organism lactobacillus acidophilus r26. Biochem J 90:76–81

174. Rappaport HP (1988) The 6-thioguanine/5-methyl-2-pyrimidinone base pair. Nucl Acids Res 16(15):7253–7267

175. Rappaport HP (1989) Artificial DNA base pair analogues

176. Wu P, Norland TM, Gildea B, McLaughlin LW (1990) Base stacking and unstacking as determined from a DNA decamer containing a fluorescent base. Biochem 29:6508–6514

177. Berry DA, Jung K-Y, Wise DS, Sercel AD, Pearson WH, Mackie H, Randolph JB, Somers RL (2004) Pyrrolo-dc and pyrrolo-c: fluorescent analogs of cytidine and 2'-deoxycytidine for the study of oligonucleotides. Tet Lett 45:2457–2461

178. Barrio JR, Sattsangi PD, Gruber GA, Dammann LG, Leonard NJ (1976) Species responsible for 3,n4-ethenocytidine. J Am Chem Soc 98(23):7408–7414

179. Major DT, Fisher B (2003) Theoretical study of the ph-dependent photophysics of n1, n6-ethenoadenine and n3,n4-ethenocytosine. J Phys Chem A 107:8923–8931

180. Kistler KA, Matsika S (2007) The fluorescence mechanism of 5-methyl-2-pyrimidinone: an ab initio study of a fluorescent pyrimidine analog. Photoch Photob 83:611–624

181. Thompson KC, Miyake N (2005) Properties of a new fluorescent cytosine analogue, pyrrolocytosine. J Phys Chem B 109:6012–6019

182. Hardman SJO, Thompson KC (2006) Influence of base stacking and hydrogen bonding on the fluorescence of 2-aminopurine and pyrrolocytosine in nucleic acids. Biochem 45: 9145–9155

183. Jean JM, Hall KB (2002) 2-aminopurine electronic structure and fluorescence properties in DNA. Biochem 41:13152–13161

184. Sowers LC, Fazakerley GV, Eritja R, Kaplan BE (1986) Base pairing and mutagenesis: observation of a protonated base pair between 2-aminopurine and cytosine in an oligonucleotide by proton NMR. Proc Natl Acad Sci USA 83:5434–5438

185. Rachofsky EL, Osman R, Ross JBA (2001) Probing structure and dynamics of DNA with 2-aminopurine: effects of local environment on fluorescence. Biochem 40:946–956

186. Wang Q, Raytchev M, Fiebig T (2007) Ultrafast energy delocalization and electron transfer dynamics in 2-aminopurine-containing trinucleotides. Photochem Photobiol 83:637–641

187. Seefeld KA, Plützer C, Löwenich D, H"aber T, Linder R, Kleinermanns K, Tatchen J, Marian CM (2005) Tautomers and electronic states of jet-cooled 2-aminopurine investigated by double resonance spectroscopy and theory. Phys Chem Chem Phys 7:3021–3026

188. Kodali G, Kistler KA, Matsika S, Stanley RJ (2007) 2-aminopurine excited state electronic structure measured by stark spectroscopy. J Phys Chem B 111:10615–10625

189. Häupl T, Windolph C, Jochum T, Brede O, Hermann R (1997) Picosecond fluorescence of nucleic acid bases. Chem Phys Lett 280:520–524

190. Rachofsky EL, Ross JBA, Krauss M, Osman R (1998) Spectroscopy of 2-aminopurine: an MCSCF study. Acta Phys Polon A 94:735–748

191. Rachofsky EL, Ross JBA, Krauss M, Osman R (2001) CASSCF investigation of electronic excited states of 2-aminopurine. J Phys Chem A 105:190–197

192. Borin AC, Serrano-Andrés L, Ludwig V, Coutinho K, Canuto S (2006) Theoretical electronic spectra of 2-aminopurine in vapor and in water. Int J Quant Chem 106: 2564–2577

193. Broo A (1998) A theoretical investigation of the physical reason for the very diffferent luminescence properties of the two isomers adenine and 2-aminopurine. J Phys Chem A 102:526

194. Perun S, Sobolewski AL, Domke W (2006) Ab initio studies of the photophysics of 2-aminopurine. Mol Phys 104:1113–1121

195. Mamos P, Aerschot AAV, Weyns NJ, Hardewijn PA (1992) Straightforward c-8 alkylation of adenosine analogues with tetraalkyltin reagents. Tet Lett 33:2413–2416

196. Lang P, Gerez C, Tritsch D, Fontecave M, Biellmann J-F, Burger A (2003) Synthesis of 8-vinyladenosine 50-di- and 50-triphosphate: evaluation of the diphosphate compound on ribonucleotide reductase, Tetrahedron 59:7315–7322

197. Gaied NB, Glasser N, Ramalanjaona N, Beltz H, Wolff P, Marquet R, Burger A, Mly Y (2005) 8-vinyl-deoxyadenosine, an alternative fluorescent nucleoside analog to 2'-deoxyribosyl-2-aminopurine with improved properties. Nucl Acids Res 33:1031–1039

198. Kenfack CA, Burger A, Mély Y (2006) Excited-state propertes and transitions of fluorescent 8-vinyl adenosine in DNA. J Phys Chem A 110:26327–26336

199. Spencer RD, Weber G, Tolman GL, Barrio JR, Leonard NJ (1974) Species responsible for the fluorescence of 1:N6-ethenoadenosine. Eur J Biochem 45:425–429

200. Bersuker IB (1984) The jahn-teller effect and vibronic interactions in modern chemistry. Plenum Press, New York

201. Katriel J, Davidson ER (1980) The non-crossing rule: triply degenerate ground-state geometries of $CH_4^+$. Chem Phys Lett 76:259

202. Keating SP, Mead CA (1985) Conical intersections in a system of four identical nuclei. J Chem Phys 82:5102

203. Han S, Yarkony DR (2003) Conical intersections of three states. Energies, derivative couplings, and the geometric phase effect in the neighborhood of degeneracy subspaces. Application to the allyl radical. J Chem Phys 119:11562

204. Han S, Yarkony DR (2003) Nonadiabatic processes involving three electronic states. I. Branch cuts and linked pairs of conical intersections. J Chem Phys 119:5058

205. Coe JD, Martinez TJ (2005) Competitive decay at two- and three-state conical intersections in excited-state intramolecular proton transfer. J Am Chem Soc 127:4560

206. Coe JD, Martinez TJ (2006) Ab initio molecular dynamics of excited-state intramolecular proton transfer around a three-state conical intersection in malonaldehyde. J Phys Chem A 110:618–630

207. Matsika S, Yarkony DR (2002) Accidental conical intersections of three states of the same symmetry. I. Location and relevance. J Chem Phys 117:6907

208. Matsika S, Yarkony DR (2003) Beyond two-state conical intersections. Three-state conical intersections in low symmetry molecules: The allyl radical. J Chem Soc 125:10672

209. Matsika S, Yarkony DR (2003) Conical intersections of three electronic states affect the ground state of radical species with little Or no symmetry: pyrazolyl. J Am Chem Soc 125:12428

210. Matsika S (2005) Three-state conical intersections in nucleic acid bases. J Phys Chem A 109:7538

211. Zazza C, Amadei A, Sanna N, Grandi A, Chillemi G, Nola AD, D'Abramo M, Aschi M (2006) Theoretical modeling of the valence UV spectra of 1,2,3-triazine and uracil in solution. Phys Chem Chem Phys 8:1385–1393

212. Zhang RB, Zeegers-Huyskens T, Ceulemans A, Nguyen MT (2005) Interaction of triplet uracil and thymine with water. Chem Phys 316:35

213. Improta R, Barone V (2004) Absorption and fluorescence spectra of uracil in the gas phase and in aqueous solution: a TDDFT quantum mechanical study. J Am Chem Soc 126:14320

214. Ludwig V, Coutinho K, Canuto S (2007) A Monte Carlo-quantum mechanics study of the lowest $n - \pi^*$ and $\pi - \pi^*$ states of uracil in water. Phys Chem Chem Phys 9:4907–4912

215. Shukla MK, Leszczynski J (2002) Interaction of water molecules with cytosine tautomers: An excited-state quantum chemical investigation. J Phys Chem A 106:11338

216. Blancafort L, Migani A (2007) Water effect on the excited-state decay paths of singlet excited cytosine. J Photochem Photobiol A 190:283–289

217. Mishra SK, Shukla MK, Mishra PC (2000) Electronic spectra of adenine and 2-aminopurine: an ab initio study of energy level diagrams of different tautomers in gas phase and aqueous solution. Spectrochim Acta A 56:1355

218. Shukla MK, Mishra SK, Kumar A, Mishra PC (2000) An ab initio study of excited states of guanine in the gas phase and aqueous media: Electronic transitions and mechanism of spectral oscillations. J Comput Chem 21:826

219. Mennucci B, Toniolo A, Tomasi J (2001) Theoretical study of the photophysics of adenine in solution: tautomerism, deactivation mechanisms, and comparison with the 2-aminopurine fluorescent isomer. J Phys Chem A 105:4749

220. Mennucci B, Toniolo A, Tomasi J (2001) Theoretical study of guanine from gas phase to aqueous solution: role of tautomerism and its implications in absorption and emission spectra. J Phys Chem A 105:7126

221. Shukla MK, Leszczynski J (2005) Effect of hydration on the lowest singlet $\pi\pi^*$ excited-state geometry of guanine: a theoretical study. J Phys Chem B 109:17333–17339

222. Shukla MK, Leszczynski J (2005) Excited state proton transfer in guanine in the gas phase and in water solution: a theoretical study. J Phys Chem A 109:7775–7780

223. Yamazaki S, Kato S (2007) Solvent effect on conical intersections in excited-state 9H-adenine: radiationless decay mechanism in polar solvent. J Am Chem Soc 129:2901–2909

224. Yoshikawa A, Matsika S (2008) Excited electronic states and photophysics of uracil-water complexes. Chem Phys 347:393–404

225. Mourik TV, Price SL, Clary DC (1999) Ab initio calculations on uracil-water. J Phys Chem A 103:1611

226. Crespo-Hernandez CE, Cohen B, Kohler B (2006) Molecular spectroscopy: complexity of excited-state dynamics in DNA - Reply. Nature 441:E8

227. Emanuele E, Zakrzewska K, Markovitsi D, Lavery R, Millie P (2005) Exciton states of dynamic dna double helices: alternating dcdg sequences. J Phys Chem B 109:16109–16118

228. Bittner ER (2007) Frenkel exciton model of ultrafast excited state dynamics in AT DNA double helices. J Photochem Photobiol A 190:328–334

229. Buchvarov I, Raytchev QWQM, Trifonov A, Fiebig T (2007) Electronic energy delocalization and dissipation in single- and double-stranded DNA. Proc Natl Acad Sci USA 104:4794–4797

230. Kwok W-M, Ma C, Phillips DL (2006) Femtosecond time- and wavelength-resolved fluorescence and absorption spectroscopic study of the excited states of adenosine and an adenine oligomer. J Am Chem Soc 128:11894–11905

231. Nir E, Kleinermanns K, de Vries MS (2000) On the photochemistry of purine nucleobases. Nature 408:949

232. Weinkauf R, Schermann JP, de Vries MS, Kleinermanns K (2002) Molecular physics of building blocks of life under isolated or defined conditions. Eur Phys J D 20:309

233. Plützer C, Hünig I, Kleinermanns K, Nir E, de Vries MS (2003) On the photochemistry of purine nucleobases. Chem Phys Chem 4:838–842

234. Schultz T, Samoylova E, Radloff W, Hertel IV, Sobolewski AL, Domcke W (2004) Efficient deactivation of a model base pair via excited-state hydrogen transfer. Science 306:1765

235. Sobolewski AL, Domcke W, Hättig C (2005) Tautomeric selectivity of the excited-state lifetime of guanine/cytosine base pairs: The role of electron-driven proton-transfer processes. Proc Natl Acad Sci USA 102:17903–17906

236. Sobolewski AL, Domcke W (2003) Ab initio study of the excited-state coupled electron-proton-transfer process in the 2-aminopyridine dimer. Chem Phys 294:2763

237. Sobolewski AL, Domcke W (2004) Ab initio studies on the photophysics of the guanine-cytosine base pair. Phys Chem Chem Phys 6:2763

238. Sobolewski AL, Domcke W (2006) Role of electron-driven proton-transfer processes in the excited-state deactivation adenine–thymine base pair. J Phys Chem A 110: 9031–9038

239. Sobolewski AL, Domcke W (2007) Computational studies of the photophysics of hydrogen-bonded molecular systems. J Phys Chem A 111:11725–11735

240. Shukla MK, Leszczynski J (2002) A theoretical investigation of excited-state properties of the adenine–uracil base pair. J Phys Chem A 106:1011–1018

241. Groenhof G, Schafer LV, Boggio-Pasqua M, Goette M, Grubmuller H, Robb MA (2007) Ultrafast deactivation of an excited cytosine-guanine base pair in DNA. J Am Chem Soc 129:6812–6819

242. Markwick PRL, Doltsinis NL (2007) Ultrafast repair of irradiated DNA: nonadiabatic ab initio simulations of the guanine-cytosine photocycle. J Chem Phys 126:175102

243. Markwick PRL, Doltsinis NL, Schlitter J (2007) Probing irradiation induced DNA damage mechanisms using excited state Car-Parrinello molecular dynamics. J Chem Phys 126:045104

244. Jean JM, Hall KB (2001) 2-aminopurine fluorescence quenching and lifetimes: role of base stacking. Proc Natl Acad Sci USA 98:37–41

245. Olaso-González G, Roca-Sanjuán D, Serrano-Andrés L, Merchán M (2006) Toward the understanding of DNA fluorescence: the singlet excimer of cytosine. J Chem Phys 125:231102

246. Santoro F, Barone V, Improta R (2007) Influence of base stacking on excited-state behavior of polyadenine in water, based on time-dependent density functional calculations. Proc Natl Acad Sci USA 104:9931–9936

247. Varsano D, Di Felice R, Marques MAL, Rubio A (2006) A TDDFT study of the excited states of DNA bases and their assemblies. J Phys Chem B 110:7129–7138

248. Tinoco I (1960) Hypochromism in polynucleotides. J Am Chem Soc 82:4785–4790

249. Frenkel J (1931) On the transformation of light into heat in solids. II. Phys Rev 37:1276–1294

250. Bouvier B, Gustavsson T, Markovitsi D, Millié P (2002) Dipolar coupling between electronic transitions of the DNA bases and its relevance to exciton states in double helices. Chem Phys 275:75–92

251. Bouvier B, Dognon J-P, Lavery R, Markovitsi D, Millié P, Onidas D, Zakrzewska K (2003) Influence of conformational dynamics on the exciton states of DNA oligomers. J Phys Chem B 107:13512–13522

252. Emanuele E, Markovitsi D, Milli P, Zakrzewska K (2005) UV spectra and excitation delocalization in DNA: influence of the spectral width. Chem Phys Chem 6:1387–1392

253. Onidas D, Gustavsson T, Lazzarotto E, Markovitsi D (2007) Fluorescence of the DNA double helices (dAdT)n·(dAdT)n studied by femtosecond spectroscopy. Phys Chem Chem Phys 9:5143–5148

254. Onidas D, Gustavsson T, Lazzarotto E, Markovitsi D (2007) Fluorescence of the DNA double helix (dA)20·(dT)20 studied by femtosecond spectroscopy effect of the duplex size on the properties of the excited states. J Phys Chem B 111:9644–9650

255. Miannay F-A, Bányász Á, Gustavsson T, Markovitsi D (2007) Ultrafast excited-state deactivation and energy transfer in guanine-cytosine DNA double helices. J Am Chem Soc 129:14574–14575

256. Bittner ER (2006) Lattice theory of ultrafast excitonic and charge-transfer dynamics in DNA. J Chem Phys 125:094909

257. Czader A, Bittner ER (2008) Calculations of the exciton coupling elements between DNA bases using the transition density cube method. J Chem Phys 128:035101

258. Schreier WJ, Schrader TE, Koller FO, Gilch P, Crespo-Hernandez CE, Swaminathan VN, Carell T, Zinth W, Kohler B (2007) Thymine dimerization in DNA is an ultrafast photoreaction. Science 315:625–629

259. Durbeej B, Eriksson LA (2002) Reaction mechanism of thymine dimer formation in DNA induced by UV light. Photochem Photobiol A 152:95–101

260. Zhang RB, Eriksson LA (2006) A triplet mechanism for the formation of cyclobutane pyrimidine dimers in UV-irradiated DNA. J Phys Chem B 110:7556–7562

261. Boggio-Pasqua M, Groenhof G, Schäfer LV, Brubmüller H, Robb MA (2007) Ultrafast deactivation channel for thymine dimerization. J Am Chem Soc 129:10996–10997
262. Blancafort L, Migani A (2007) Modeling thymine photodimerizations in DNA: mechanism and correlation diagrams. J Am Chem Soc 129:14540–14541
263. Valiev M, Kowalski K (2006) Hybrid coupled cluster and molecular dynamics approach: Application to the excitation spectrum of cytosine in the native DNA environment. J Chem Phys 125:211101
264. Salter LM, Chaban GM (2002) Theoretical study of gas phase tautomerization reaction for the ground and first excited electronic states of adenine. J Phys Chem A 106:4251

CHAPTER 12

# AB INITIO QUANTUM MECHANICAL/MOLECULAR MECHANICAL STUDIES OF HISTONE MODIFYING ENZYMES

YINGKAI ZHANG

*Department of Chemistry, New York University, New York, NY 10003, USA,*
*e-mail: yingkai.zhang@nyu.edu*

**Abstract:**    Histone proteins that form the nucleosome core are subject to a variety of post-translational transformations. These histone modifications make up the histone code which extends the information in the genetic code and is emerging as an essential mechanism to regulate gene expression. In spite of a current flurry of significant advances in experimental studies, there has been little theoretical understanding regarding how enzymes generate or remove these modifications. Very recently, we have made excellent progresses in investigating two such important histone-modifying enzyme families: zinc-dependent histone deacetylases (HDACs) and histone lysine methyltransferases (HKMTs). Our studies on a histone-deacetylase-like protein HDLP suggested a novel catalytic mechanism. The simulations on HKMT SET7/9 have characterized the histone lysine methylation reaction and elucidated the origin of enzyme catalysis. Our computational approaches centered on the pseudobond ab initio quantum mechanical/molecular mechanical (QM/MM) method, which allows for accurate modeling of the chemistry at the reaction active site while properly including the effects of the protein environment

**Keywords:**    Ab initio QM/MM method, Pseudobond approach, Enzyme catalysis, Reaction mechanism, Molecular dynamics simulation

## 12.1.    INTRODUCTION

In the nuclei of all eukaryotic cells, DNA is tightly wrapped around an octamer of histone proteins and is compacted into a dense structure known as chromatin. In order to access the genetic information which is required in numerous essential cellular processes including DNA replication, gene expression and DNA repair, chromatin needs to be partially unwound. One important mechanism to regulate chromatin structure and thus to control the access of the genomic DNA is through histone modifications [1–6]. The histone octamer is composed of two copies of H2A, H2B, H3 and H4 core histone proteins. Their tails, that protrude out of the surface of the

chromatin assemblies, are subject to a variety of reversible covalent modifications including acetylation, methylation and phosphorylation. A specific pattern of histone tail transformation creates a distinguishing histone code to control the access of the genomic DNA thereby leading to a particular downstream event [1, 2, 7–9]. Failure of appropriate histone modifications can lead to aberrant gene regulation and is implicated in human diseases, notably cancer [10, 11].

There are two key questions that are central to the histone code concept [1, 2]: (1) How the code is written, i.e., how enzymes add or remove the modifications at the specific target sites in the histones. (2) How the code is translated, i.e., how the histone tails displaying appropriate modification patterns promote the recognition by effector proteins. Very recently, we have made excellent progresses on the first question by investigating histone deacetylases and histone lysine methyltransferases. To understand the inner workings of histone modifying enzymes is not only of fundamental importance, but also an essential starting point for the development of novel anti-tumor agents and new disease therapies. In spite of significant advances in experimental studies, many questions regarding mechanisms and structure-functional relationship remain open. Due to the short lifetime of enzymatic transition states, it is often extremely difficult for experimental methods alone to directly characterize them. Thus, theoretical help is much needed.

## 12.2.    PSEUDOBOND AB INITIO QM/MM APPROACH

To simulate enzyme reactions, high level QM methods are required to describe chemical bond forming and breaking. Meanwhile, the catalytic power of enzymes is not only determined by its active site, but also be controlled by its heterogeneous protein environment. Thus the enzyme environment needs to be properly described. In light of these issues, the field of development and application of combined quantum mechanical/molecular mechanical (QM/MM) approaches [12–25] is expanding rapidly, and the QM/MM methods have emerged as the method of choice in simulating chemical reactions in complex systems. In the QM/MM framework a small chemically active region is treated by a quantum mechanical method, while the remainder of the system containing a large number of atoms is described by a molecular mechanical force field. Such a combined QM and MM approach can take advantage of the applicability and accuracy of the QM methods for chemical reactions and of the computational efficiency of the MM calculations, and thus extend the realm of quantum mechanical calculations to large systems.

In our simulations of histone modifying enzymes, the computational approaches centered on the pseudobond ab initio quantum mechanical/molecular mechanical (QM/MM) approach. This approach consists of three major components [20, 26–29]: a pseudobond method for the treatment of the QM/MM boundary across covalent bonds, an efficient iterative optimization procedure which allows for the use of the ab initio QM/MM method to determine the reaction paths with a realistic enzyme environment, and a free energy perturbation method to take account

*Figure 12-1.* Illustration of the difference between the pseudobond approach and the conventional link atom approach in the treatment of the QM/MM boundary problem

of protein dynamics. In comparison with the conventional link-atom approach, a key advantage of the the pseudobond approach for the treatment of the QM/MM covalent boundary problem is that it does not introduce additional degrees of freedom into the system, as illustrated in Figure 12-1. Instead, the boundary carbon atom of the environment part, is replaced by a special atom, $C_{ps}$. This $C_{ps}$ atom has the following properties: it is a seven-valence-electron atom, which means that it only has one free valence. Therefore, the $C_{ps}$ and the atoms in the active form a well-defined closed shell system, which can be described by the quantum mechanical method. This $C_{ps}$ atom also has an effective core potential. By designing the effective core potential of the $C_{ps}$ atom, this carbon—carbon pseudobond is made to mimic the original carbon—carbon bond with the similar bond length and strength, and also similar effects on the rest of the active part. The pseudobond ab initio QM/MM approach [20, 26–29] has been employed in the study of several enzymes, including enolase [30], acetylcholinesterase [31, 32], 4-oxalocrotonate tautomerase [33], and kinase [34, 35]. These studies demonstrate that the method is powerful in providing detailed insights into enzyme reactions. Some theoretical predictions [30, 33] were subsequently confirmed by experimental studies [36–38]. In spite of its success, the pseudobond method is still much in need of further development. In the original paper [26], only $C(sp^3)$—$C(sp^3)$ pseudobond has been successfully developed, which limits its applicability. In the case of proteins, it can only be used to cut the protein side chains. In order to cut the protein backbones, we need the $C(sp^3)$—$N(sp^3)$ as well as $C(sp^3)$—Carbonyl carbon pseudobond. For the cut of DNA and RNA bases, we also need the $C(sp^3)$—$N(sp^3)$ pseudobonds. So the challenge is whether we can develop better and more pseudobonds.

In order to improve the accuracy and applicability of the pseudobond approach, very recently we have developed a new formulation to construct this seven-valence-electron boundary atom [29]. *The key difference is that the seven-valence-electron boundary atom has its own basis set instead of that of fluorine.* Here a STO-2G

basis set has been employed for the seven-valence-electron boundary atom, which has four parameters. For its effective core potential, we use an angular momentum independent formula which has only two parameters. In this way, we introduce a minimum number of parameters. By parameterizing both the basis set and the effective core potential, we have been able to not only significantly improve the $C_{ps}(sp^3)$—$C(sp^3)$ pseudobond, but also develop accurate $C_{ps}(sp^3)$—$C(sp^2$, carbonyl) and $C_{ps}(sp^3)$—$N(sp^3)$ pseudobonds for the cutting of protein backbones and nucleic acid bases for the first time. It should be noted that all rest atoms in QM sub-system are described with the 6-31G* basis set. The developed pseudobonds are independent of the molecular mechanical force field. Although the parameterization is performed with density functional calculations using hybrid B3LYP exchange-correlation functional, it is found that the same set of parameters is also applicable to Hartree-Fock and MP2 methods, as well as DFT calculations with other exchange-correlation functionals. Tests on a series of molecules yield very good structural, electronic and energetic results in comparison with the corresponding full ab initio quantum mechanical calculations. The standard deviations (SD) between the pseudobond QM calculations and the corresponding standard B3LYP(6-31G*) calculations for bond lengths, angles and atomic Mulliken charge are 0.013 Å, 0.8 degree and 0.03 for eight molecules. We have also tested the energy differences for these four pairs of molecules, which are the deprotonation energies for molecules $CH_3CH_2OH$, $CH_3CH_2NH_3^+$, $CH_3CH_2SH$, $CH_3CH_2COOH$ respectively. This is a quite stringent test since the pseudobond is only one or two bonds away from the reaction bond X—H, X=O,N,S,O . For this test, the mean absolute error is only $1.9\,kcal\,mol^{-1}$. The performance for the other two pseudobonds are also very similar [29]. Meanwhile, it should be noted that these developed pseudobond parameters are semi-empirical in nature because they are parameterized against molecular properties with a limited set of molecules. It would be ideal to design the boundary atom to mimic atomic properties of the corresponding atom instead of molecular properties, thus it could be more fundamental and transferable.

## 12.3. AN UNEXPECTED REACTION MECHANISM FOR HISTONE DEACETYLATION [39]

Histone deacetylases (HDACs) function in opposition to histone acetylases by catalyzing the cleavage of acetyl groups from acetyl-lysine residues in histone N-terminal tails [6, 40, 41]. The deacetylation of histones often leads to closed chromatin structure and is generally associated with transcriptional repression and gene silencing. Based on the sequence analysis, 18 distinguished human HDACs have been classified into four different classes [42]. Class I and II HDACs are zinc dependent hydrolases, which are among the most promising anti-cancer targets. The aberrant recruitment of class I and II enzymes has been associated with a variety of diseases, and their inhibitors have shown to be able to inhibit the growth and survival of tumor cells in model systems and clinical trials [43–46]. The detailed knowledge

of their catalytic mechanism for class I and II HDACs is of great interest and high importance because it may open the way for the design of novel HDAC inhibitors with enhanced potency and specificity.

Before our work [39], only one catalytic mechanism for zinc dependent HDACs has been proposed in the literature, which was originated from the crystallographic study of HDLP [47], a histone-deacetylase-like protein that is widely used as a model for class-I HDACs. In the enzyme active site, the catalytic metal zinc is penta-coordinated by two asp residues, one histidine residues as well as the inhibitor [47]. Based on their crystal structures, Finnin *et al.* [47] postulated a catalytic mechanism for HDACs in which the first reaction step is analogous to the "hydroxide mechanism" for zinc proteases: zinc-bound water is a nucleophile and $Zn^{2+}$ is five-fold coordinated during the reaction process. However, recent experimental studies by Kapustin *et al.* suggested that the transition state of HDACs may not be analogous to zinc-proteases [48], which cast some doubts on this mechanism.

To characterize the reaction mechanism for histone deacetylation, we have carried out B3LYP QM/MM calculations on the deacetylation reaction catalyzed by the HDLP [39]. Our results do not support the previous mechanistic hypothesis, but suggested an unexpected reaction mechanism for histone deacetylation. In this new mechanism, both histidine residues in the second coordination shell are singly coordinated in the reactant state. Although the water molecule is weakly coordinated to the zinc atom in the reactant, the water immediately departs from the metal center when the reaction proceeds. During the rest of the reaction process, the Zinc atom clearly has a tetrahedral coordination. Therefore, the key catalytic role of the zinc atom is to activate the carbonyl group of the amide towards nucleophilic attack. This tetrahedral zinc coordination during the reaction process can explain the experimental fact that the phosphorous based SAHA ligands are not strong inhibitors to HDAC [48]. Besides the reaction mechanism, we can also identify important residues in catalysis by determining the individual residue energetic contribution to the stabilization of this transition state. It was found that Tyr297, Asp116 and Asp173 play important roles to stabilize the transition state, which are consistent with experimental mutational studies [47]

## 12.4. ENZYME MECHANISM AND CATALYSIS OF HISTONE LYSINE METHYLATION [49, 50]

In addition to histone deacetylation, histone lysine methylation can also lead to gene silencing which is not blocked by the HDAC inhibitors [6, 51]. Several lines of evidence have suggested a connection between cancer and histone lysine methyltransferases (HKMTs) [52]. HKMTs catalyze the transfer of methyl group(s) from the cofactor *S*-adenosyl-methionine (AdoMet) to some specific lysine residues in the N-terminal histone tails [53, 54]. With one exception of Dot1 [55], all known HKMTs contain the SET domain which represents a novel structural fold [53, 56]. Among SET-domain HKMTs, SET7/9 is one of the best characterized experimentally. It is a

mono-methyltransferase which only catalyzes the transfer of one methyl group to the unmodified histone lysine residue H3—K4 [57, 58]. In its active site, the cofactor AdoMet and the substrate peptide bind to opposite faces of the SET-domain and are connected by a narrow channel which has a hydrophobic inner wall [58, 59]. The target lysine residue is inserted into this narrow channel to access the methyl moiety of AdoMet.

To investigate the methyl transfer reaction catalyzed by SET7/9, we have carried out multiple ab initio quantum mechanical/molecular mechanical free energy (QM/MM-FE) calculations and molecular dynamics simulations. The QM sub-system consists of AdoMet and the side-chain of histone lysine residue H3—K4. With the $C_\alpha$—$C_\beta$ bond of H3—K4 treated as a pseudobond [26], the resulting QM sub-system has 66 atoms. All other 9497 atoms are described classically. In order to take account of enzyme dynamics, 11 different snapshots of the enzyme-substrate complex from a molecular dynamics trajectory have been used as initial structures for QM/MM studies [32]. Enzyme reaction paths were determined by B3LYP/(6-31G*) QM/MM calculations with an iterative minimization procedure and the reaction coordinate driving method [27]. For determined reaction paths, single point MP2 QM/MM calculations with both 6-31G* and 6-31+G* basis sets have been carried out. Free energies along the reaction paths were determined by free energy perturbation calculations and the harmonic approximation [27, 28]. The calculated average free energy reaction barrier is $20.4\pm1.4$ kcal mol$^{-1}$ and $20.7\pm1.4$ kcal mol$^{-1}$ for MP2(6-31+G*)/MM and MP2(6-31G*)/MM calculations respectively, which is in excellent agreement with the activation barrier of 20.9 kcal mol$^{-1}$ estimated from the experimental value of $k_{cat}$ [60].

QM/MM calculations from different starting structures yield a consistent mechanistic picture for this methyl-transfer reaction: a typical in-line $S_N2$ nucleophilic substitution reaction with a mainly dissociative transition state. In the transition state, the CH$_3$ plane is almost at the middle of the attacking N$_\zeta$ atom and the leaving S$_\delta$ atom. The average distances of $S_\delta \cdots C_\epsilon$ and $C_\epsilon \cdots N_\zeta$ bonds are $2.32\pm0.02$ Å and $2.30\pm0.02$ Å respectively, which are much longer than their covalent bond distances. Thus this methyl-transfer transition state can be considered to be loose. The loose character of the transition state is further confirmed by the bond order analysis [61]. At the transition state, the bond orders of $S_\delta \cdots C_\epsilon$ and $C_\epsilon \cdots N_\zeta$ bonds are $0.58\pm0.01$ and $0.30\pm0.01$, respectively. Since the fractional associativity of a mechanism is defined by the bond order of the attacking group [62], the methyl-transfer reaction catalyzed by SET7/9 can be quantified as 30% associative and 70% dissociative based on the calculated Wiberg bond order of $0.30\pm0.01$ for $C_\epsilon \cdots N_\zeta$ at the transition state.

With the characterized mechanism, the next key question is the origin of its catalytic power. A prerequisite for this investigation is to reliably compute free energy barriers for both enzyme and solution reactions. By employing on-the-fly Born-Oppenheimer molecular dynamics simulations with the ab initio QM/MM approach and the umbrella sampling method, we have determined free energy profiles for the methyl-transfer reaction catalyzed by the histone lysine methyltransferase SET7/9

and its corresponding uncatalyzed reaction in aqueous solution, respectively [50]. At each time step, the forces on atoms in both QM and MM sub-systems as well as the total energy are calculated with a pseudobond ab initio QM/MM method on the fly, and Newton equations of motion are integrated. Our calculated activation free energy barrier for the methyl transfer reaction catalyzed by SET7/9 is 22.5 kcal/mol, which agrees with the experimental value of 20.9 kcal/mol very well. The difference in potential of mean force between a corresponding pre-reaction state and the transition state for the solution reaction is computed to be 30.9 kcal/mol. Thus the enzyme SET7/9 lowers the barrier for the methyl-transfer reaction step by 8.4 kcal/mol compared with the uncatalyzed reaction, which corresponds to a rate enhancement of about one million fold.

In order to elucidate the origin of this enormous enzyme catalytic power, first we analyzed geometries of pre-reaction states and transition states. It was found that the mechanism difference between the enzyme reaction and the solution reaction is very small and is not likely to be a key source of enzyme catalysis. Then we examined the electrostatic field and hydrogen bond network in the reaction center. The calculated results indicate that the enzyme SET7/9 provides a pre-organized electrostatic environment to facilitate the methyl-transfer reaction, while for the reaction in solution, the hydrogen-bond network near the reaction center undergoes a a significant change and there is a strong shift in electrostatic field from the pre-reaction state to the transition state. Thus our results indicate that a combination of the electrostatic pre-organization in enzyme and the hydrogen bond network reorganization in solution is an essential contributor to the enormous catalytic power of the histone lysin methyltransferase SET7/9.

## 12.5.    CONCLUSIONS

In this chapter, I have summarized our most recent progresses in simulating histone modifying enzymes and the further development of the pseudobond approach. We not only suggested a novel catalytic mechanism of histone deacetylation based on the DFT QM/MM studies of a histone-deacetylase-like protein HDLP, but also characterized the histone lysine methylation reaction mechanism and elucidated the origin of enzyme catalysis by simulating HKMT SET7/9 with ab initio QM/MM approaches. To provide such detailed insights into these important enzymes is not only of fundamental importance, but also an essential starting point for the development of novel anti-tumor agents and new disease therapies. Meanwhile, these studies further demonstrate that ab initio QM/MM methods have become increasingly powerful in complementing experimental methods to elucidate the chemistry in complex systems, including biocatalysis. Since histone modification is emerging as an essential mechanism in regulating chromatin structure and gene expression, it is expected that much more theoretical and computational studies will be carried out to understand its inner workings.

## ACKNOWLEDGMENTS

## REFERENCES

1. Strahl BD, Allis CD (2000) The language of covalent histone modifications. Nature 403:41–45
2. Jenuwein T, Allis CD (2001) Translating the histone code. Science 293:1074–1080
3. Iizuka M, Smith MM (2003) Functional consequences of histone modifications. Curr Opin Genet Dev 13:154–160
4. Khorasanizadeh S (2004) The nucleosome: from genomic organization to genomic regulation. Cell 116:259–272
5. Khan AU, Krishnamurthy S (2005) Histone modifications as key regulators of transcription. Front Biosci 10:866–872
6. Biel M, Wascholowski V, Giannis A (2005) Epigenetics – an epicenter of gene regulation: histones and histone-modifying enzymes. Angew Chem-Int Edit 44:3186–3216
7. Turner BM (2002) Cellular memory and the histone code. Cell 111:285–291
8. Fischle W, Wang YM, Allis CD (2003) Binary switches and modification cassettes in histone biology and beyond. Nature 425:475–479
9. Margueron R, Trojer P, Reinberg D (2005) The key to development: interpreting the histone code? Curr Opin Genet Dev 15:163–176
10. Hake SB, Xiao A, Allis CD (2004) Linking the epigenetic 'language' of covalent histone modifications to cancer. Br J Cancer 90:761–769
11. Santos-rosa H, Caldas C (2005) Chromatin modifier enzymes, the histone code and cancer. Eur J Cancer 41:2381–2402
12. Warshel A, Levitt M (1976) Theoretic studies of enzymic reactions: dielectric electrostatic and steric stabilization if the carbonium ion in the reaction of lysozyme. J Mol Bio 103:227
13. Singh UC, Kollman PA (1986) A combined ab initio quantum mechanical and molecular mechanical method for carrying out simulations on complex molecular systems: applications to the ch_3cl + cl – exchange reaction and gas phase protonation of polyethers. J Comp Chem 7:718–730
14. Field MJ, Bash PA, Karplus M (1990) A combined quantum mechanical and molecular mechanical potential for molecular dynamics simulations. J Comp Chem 11:700–733
15. Gao J, Truhlar DG (2002) Quantum mechanical methods for enzyme kinetics. Annu Rev Phys Chem 53:467–505
16. Gao JL, Ma SH, Major DT, Nam K, Pu JZ, Truhlar DG (2006) Mechanisms and free energies of enzymatic reactions. Chem Rev 106:3188–3209
17. Warshel A, Sharma PK, Kato M, Xiang Y, Liu HB, Olsson MHM (2006) Electrostatic basis for enzyme catalysis. Chem Rev 106:3210–3235
18. Friesner RA, Guallar V (2005) Ab initio quantum chemical and mixed quantum mechanics/molecular mechanics (QM/MM) methods for studying enzymatic catalysis. Annu Rev Phys Chem 56:389–427
19. Mulholland AJ (2005) Modelling enzyme reaction mechanisms, specificity and catalysis. Drug Discov Today 10:1393–1402

20. Zhang Y (2006) Pseudobond ab initio QM/MM approach and its applications to enzyme reactions. Theor Chem Acc 116:43–50

21. Riccardi D, Schaefer P, Yang Y, Yu HB, Ghosh N, Prat-resina X, Konig P, Li GH, Xu DG, Guo H, Elstner M, Cui Q (2006) Development of effective quantum mechanical/molecular mechanical (QM/MM) methods for complex biological processes. J Phys Chem B 110:6458–6469

22. Bruice TC (2006) Computational approaches: reaction trajectories, structures, and atomic motions: enzyme reactions and proficiency. Chem Rev 106:3119–3139

23. Hammes-schiffer S (2004) Quantum-classical simulation methods for hydrogen transfer in enzymes: a case study of dihydrofolate reductase. Curr Opin Struct Biol 14:192–201

24. Garcia-viloca M, Gao J, Karplus M, Truhlar DG (2004) How enzymes work: analysis by modern rate theory and computer simulations. Science 303:186–195

25. Senn HM, Thiel W (2007) QM/MM methods for biological systems. Top Curr Chem 268:173–290

26. Zhang Y, Lee TS, Yang W (1999) A pseudobond approach to combining quantum mechanical and molecular mechanical methods. J Chem Phys 110:46–54

27. Zhang Y, Liu H, Yang W (2000) Free energy calculation on enzyme reactions with an efficient iterative procedure to determine minimum energy paths on a combined ab initio QM/MM potential energy surface. J Chem Phys 112:3483–3492

28. Zhang Y, Liu H, Yang W (2002) Ab initio QM/MM and free energy calculations of enzyme reactions. In: Schlick T., Gan H. H., (ed) Methods for Macromolecular Modeling. Springer-Verlag; Berlin, pp 332–354

29. Zhang Y (2005) Improved pseudobonds for combined ab initio quantum mechanical/molecular mechanical methods. J Chem Phys 122:024114

30. Liu H, Zhang Y, Yang W (2000) How is the active site of enolase organized to catalyze two different reaction steps? J Am Chem Soc 122:6560–6570

31. Zhang Y, Kua J, McCammon JA (2002) Role of the catalytic triad and oxyanion hole in acetylcholinesterase catalysis: an ab initio QM/MM study. J Am Chem Soc 124:10572–10577

32. Zhang Y, Kua J, McCammon JA (2003) Influence of structural fluctuation on enzyme reaction energy barriers in combined quantum mechanical/molecular mechanical studies. J Phys Chem B 107: 4459–4463

33. Cisneros GA, Liu H, Zhang Y, Yang W (2003) Ab initio QM/MM study shows there is no general acid in the reaction catalyzed by 4-oxalocrotonate tautornerase. J Am Chem Soc 125:10384–10393

34. Cheng Y, Zhang Y, McCammon JA (2005) How does the camp-dependent protein kinase catalyze the phosphorylation reaction: an ab initio QM/MM study. J Am Chem Soc 127:1553–1562

35. Cheng Y, Zhang Y, McCammon JA (2006) How does activation loop phosphorylation modulate catalytic activity in the camp-dependent protein kinase: a theoretical study. Protein Sci 15: 672–683

36. Poyner RR, Larsen TM, Wong SW, Reed GH (2002) Functional and structural changes due to a serine to alanine mutation in the active-site flap of enolase. Arch Biochem Biophys 401:155–163

37. Cisneros GA, Wang M, Silinski P, Fitzgerald MC, Yang WT (2004) The protein backbone makes important contributions to 4-oxalocrotonate tautomerase enzyme catalysis: understanding from theory and experiment. Biochemistry 43:6885–6892

38. Metanis N, Brik A, Dawson PE, Keinan E (2004) Electrostatic interactions dominate the catalytic contribution of arg39 in 4-oxalocrotonate tautomerase. J Am Chem Soc 126:12726–12727

39. Corminboeuf C, Hu P, Tuckerman ME, Zhang Y (2006) Unexpected catalytic mechanism for histone deacetylase suggested by a density functional theory QM/MM study. J Am Chem Soc 128: 4530–4531

40. De ruijter AJM, Vangennip AH, Caron HN, Kemp S, Vankuilenburg ABP (2003) Histone deacetylases (hdacs): characterization of the classical hdac family. Biochem J 370:737–749

41. Holbert MA, Marmorstein R (2005) Structure and activity of enzymes that remove histone modifi-
    cations. Curr Opin Struct Biol 15:673–680
42. Gregoretti IV, Lee YM, Goodson HV (2004) Molecular evolution of the histone deacetylase family:
    functional implications of phylogenetic analysis. J Mol Biol 338:17–31
43. Acharya MR, Sparreboom A, Venitz J, Figg WD (2005) Rational development of histone deacetylase
    inhibitors as anticancer agents: a review. Mol Pharmacol 68:917–932
44. Drummond DC, Noble CO, Kirpotin DB, Guo Z, Scott GK, Benz CC (2005) Clinical development
    of histone deacetylase inhibitors as anticancer agents. Annu Rev Pharmacol Toxicol 45:495–528
45. Kelly WK, Marks PA (2005) Drug insight: histone deacetylase inhibitors - development of the new
    targeted anticancer agent suberoylanilide hydroxamic acid. Nat Clin Pract Oncol 2:150–157
46. Marks PA, Rifkind RA, Richon VM, Breslow R, Miller T, Kelly WK (2001) Histone deacetylases
    and cancer: causes and therapies. Nat Rev Cancer 1:194–202
47. Finnin MS, Donigian JR, Cohen A, Richon VM, Rifkind RA, Marks PA, Breslow R, Pavletich NP
    (1999) Structures of a histone deacetylase homologue bound to the tsa and saha inhibitors. Nature
    401:188–193
48. Kapustin GV, Fejer G, Gronlund JL, Mccafferty DG, Seto E, Etzkorn FA (2003) Phosphorus-based
    saha analogues as histone deacetylase inhibitors. Org Lett 5:3053–3056
49. Hu P, Zhang Y (2006) Catalytic mechanism and product specificity of the histone lysine methyl-
    transferase set7/9: An ab initio QM/MM-FE study with multiple initial structures. J Am Chem Soc
    128:1272–1278
50. Wang S, Hu P, Zhang Y (2007) Ab initio quantum mechanical/molecular mechanical molecular
    dynamics simulation of enzyme catalysis: the case of histone lysine methyltransferase set7/9. J Phys
    Chem B ASAP
51. Martin C, Zhang Y (2005) The diverse functions of histone lysine methylation. Nat Rev Mol Cell
    Biol 6:838–849
52. Schneider R, Bannister AJ, Kouzarides T (2002) Unsafe sets: histone lysine methyltransferases and
    cancer. Trends Biochem Sci 27:396–402
53. Xiao B, Wilson JR, Gamblin SJ (2003) Set domains and histone methylation. Curr Opin Struct Biol
    13:699–705
54. Cheng X, Collins RE, Zhang X (2005) Structural and sequence motifs of protein (histone) methyla-
    tion enzymes. Annu Rev Biophys Biomolec Struct 34:267–294
55. Min J, Feng Q, Li Z, Zhang Y, Xu R (2003) Structure of the catalytic domain of human dot1l, a
    non-set domain nucleosomal histone methyltransferase. Cell 112:711–723
56. Yeates TO (2002) Structures of set domain proteins: protein lysine methyltransferases make their
    mark. Cell 111:5–7
57. Wilson JR, Jing C, Walker PA, Martin SR, Howell SA, Blackburn GM, Gamblin SJ, Xiao B (2002)
    Crystal structure and functional analysis of the histone methyltransferase set7/9. Cell 111:105–115
58. Xiao B, Jing C, Wilson JR, Walker PA, Vasisht N, Kelly G, Howell S, Taylor IA, Blackburn GM,
    Gamblin SJ (2003) Structure and catalytic mechanism of the human histone methyltransferase set7/9.
    Nature 421:652–656
59. Kwon T, Chang JH, Kwak E, Lee CW, Joachimiak A, Kim YC, Lee JW, Cho Y (2003) Mechanism
    of histone lysine methyl transfer revealed by the structure of set7/9-adomet. EMBO J., 22:292–303
60. Trievel RC, Beach BM, Dirk LMA, Houtz RL, Hurley JH (2002) Structure and catalytic mechanism
    of a set domain protein methyltransferase. Cell 111:91–103
61. Takusagawa F, Fujioka M, Spies A, Schowen RL (1998) S-adenosylmethionine (adomet)-dependent
    methyltransferases. In: Sinnott M., (ed), Comprehensive biological catalysis: a mechanistic reference.
    Academic Press, San Diego, pp 1–30
62. Mildvan AS (1997) Mechanisms of signaling and related enzymes. Proteins 29:401–416

CHAPTER 13

# INTERPRETING THE OBSERVED SUBSTRATE SELECTIVITY AND THE PRODUCT REGIOSELECTIVITY IN ORF2-CATALYZED PRENYLATION FROM X-RAY STRUCTURES

GUANGLEI CUI[1], XUE LI[1], NING YU[2], AND KENNETH M. MERZ[1]

[1]*Department of Chemistry and the Quantum Theory Project, University of Florida, 2328 New Physics Building, P.O. Box 118435, Gainesville, FL 32611-8435, USA, e-mail: merz@qtp.ufl.edu*
[2]*Simulations Plus, Inc. 42505 10th Street West, Lancaster, CA 93534, USA*

**Abstract:** The combined QM/MM based X-ray crystallography technique is described. Its relevant strengths and weaknesses relative to traditional refinement protocols are discussed. The method is illustrated by refining Orf2 protein–ligand complexes and comparing the QM/MM based method to CNS derived results. It is shown that in this instance the QM/MM based approach give superior results to traditional MM based refinements methods as implemented in CNS

**Keywords:** X-ray crystallography, X-ray refinement, Protein–ligand complexes, Orf2

## 13.1. INTRODUCTION

X-ray crystallography is perhaps one of the most successful and broadly used approaches employed in molecular biology where the general interest is to understand how life operates at the molecular level [1, 2]. It often provides the most concrete evidence by which experimental findings from other sources can be validated, solidified, and more relevant experiment designs can be proposed. The role of X-ray crystallography in computational biochemistry and biophysics is even more critical because the research in these areas heavily relies on the availability of accurate three-dimensional structural information. In recent years, homology modeling and ab initio protein structure prediction have made tremendous progress in terms of the qualities of secondary structure assignment, overall three dimensional fold recognition, and in some cases the details in residue side chain packing [3–7]. However, X-ray quality structures are rare [8, 9], which limits their usefulness in the subsequent investigations, particularly in enzyme reaction mechanism studies where high-resolution crystal structures are the first choice.

*Scheme 13-1.* The X-ray structure determination workflow (Scheme provided by Dagmar Ringe)

Unlike small molecule crystallography, obtaining atomic resolution crystal structures of biomacromolecules (Scheme 13-1) are significantly more difficult, despite numerous advances in both experimental approaches and refinement algorithms, which include the use of synchrotron radiation sources, cryo-cooling techniques, simulated annealing (SA) [10] and the maximum likelihood formalism (MLF) [11]. The number of crystal structures determined with resolutions better than 1.0 or 1.5 Å in the Protein Database Bank is less than 1 or 10%, respectively. The majority of the crystal structures have a resolution in the range between 2.0 and 3.0 Å. With empirical parameters of bonding (molecular mechanics force field), heavy atoms of protein residues can be solved at this resolution with reasonable confidence because there are only a limited number of amino acids and their connectivity is known given the primary structure. However, the accuracy of the representation of small molecule ligands, substrates, and cofactors, will likely suffer, especially when the chemical composition (both structurally and compositionally) is not well characterized by classical force fields. More difficulties arise when protein associated molecules have large thermal motions or undergoes frequent conformational transitions within the protein binding site. Success in computational investigations of enzymatic reaction mechanism may be shadowed if these potential issues in structure quality are not addressed.

Our group has been actively developing new crystal structure refinement approaches [12, 13] that introduce quantum mechanical (QM) treatments to replace the standard molecular mechanical (MM) potentials used today. A QM description is advantageous compared to a traditional MM one because (1) the use of QM in refinement eliminates the reliance on empirical force field parameters, particularly for regions of interest that often involve structures substantially different from those found in the gas phase (e.g., covalent complexes, systems with unusually close contacts, etc.); (2) because the electron density for hydrogen atoms is close to noise level in most density maps, even for structures at 1 Å or better resolutions, oftentimes the accurate energetics information may be very useful in resolving protonation or tautomeric states of key regions in the structure; (3) with higher-level QM methods, the model quality can be systematically improved if more CPU time can be afforded. Previously, Ryde et al. and we have demonstrated our QM or QM/MM X-ray refinement approach as a valuable tool in resolving subtle structure features [14, 15]

and deriving residue protonation states [16–18]. Herein we present our recent work in interpreting the relaxed substrate specificity and product regioselectivity observed in Orf2-catalyzed prenylation using the QM/MM X-ray refinement approach. Our intention is not just to present a case study, which is interesting in its own right, but to stress the profound insight that one can obtain by combining state-of-the-art free energy calculations and QM/MM refinement of X-ray diffraction data, which leads to the emergence of a dynamic picture of the reaction mechanism, not immediately apparent in the PDB structure.

In the next section, we first give a brief overview of the QM (QM/MM) X-ray refinement methodology and implementation, which has been reported previously and then follow this with the case study.

## 13.2.    STRUCTURE REFINEMENT USING QM

X-ray crystallographic experiments measure the intensity of the diffraction peaks resulting from the X-rays scattered by electron clouds, which is related to the thermal average of electron density distributions $<\rho(\mathbf{r})>$ in the crystal by a Fourier transform:

$$F(\mathbf{S}) = \int_V \langle \rho(\mathbf{r}) \rangle \, exp \, [2\pi i \mathbf{S} \cdot \mathbf{r}] d\mathbf{r} \qquad (13\text{-}1)$$

$$I(\mathbf{S}) \propto F(\mathbf{S}) \qquad (13\text{-}2)$$

where $\mathbf{S}$ is the scattering vector $(h,k,l)$, $F(\mathbf{S})$ is the structure factor of a diffracted X-ray, and $I(\mathbf{S})$ is the intensity of a diffraction peak. It follows from Eq. (13-1) that structure factors are complex quantities that contain both amplitudes $|F(\mathbf{S})|$, which can be derived from the X-ray signals according to Eq. (13-2), and the phase angles $\alpha(\mathbf{S})$. If the latter were accurately known, the electron density distribution within the unit cell of the crystal would be constructed by an inverse Fourier summation:

$$\langle \rho(\mathbf{r}) \rangle = \frac{1}{V_{cell}} \sum_{\mathbf{S}} |F_o(\mathbf{S})| \, exp \, [-2\pi i \mathbf{S} \cdot \mathbf{r} + ia(\mathbf{S})] \qquad (13\text{-}3)$$

Since the phase angles cannot be measured in X-ray experiments, structure solution usually involves an iterative process, in which starting from a rough estimate of the phases, the structure suggested by the electron density map obtained from Eq. (13-3) and the phase computed by Eq. (13-1) are gradually refined, until the computed structure factor amplitudes from Eq. (13-1) converge to the ones observed experimentally.

At a certain stage in the refinement, the electron density map is interpreted using a model representation of the charge density distribution to extract the atomic coordinates. A commonly used scattering formalism is the independent-atom model (IAM), in which the total charge density in the crystal is approximated by the superposition

of isolated atomic densities and each atom scatters the X-ray independently and contributes to the model structure factor according to the following equation,

$$F_c(\mathbf{S}) = \sum_{j=1}^{n} n_j f_j(\mathbf{S}) \, exp \, [2\pi i \mathbf{S} \cdot \mathbf{R}_j] \cdot exp \, [2\pi^2 \mathbf{S}^t \mathbf{U}_j \mathbf{S}] \qquad (13\text{-}4)$$

The terms involving the subscript $j$ represents the contribution of atom $j$ to the computed structure factor, where $n_j$ is the occupancy, $f_j$ is the atomic scattering factor, and $\mathbf{R}_j$ is the coordinate of atom $i$. In Eq. (13-4) the thermal effects are treated as anisotropic harmonic vibrational motion and $\mathbf{U}_j = <\mathbf{u}_j \mathbf{u}_j^t>$ is the mean-square atomic displacement tensor; when the thermal motion is treated as isotropic, Eq. (13-4) reduces to:

$$F_c \mathbf{S} = \sum_{j=1}^{n} n_j f_j(\mathbf{S}) exp \, [2\pi i \mathbf{S} \cdot \mathbf{R}_j] \cdot exp[-B_j \mathbf{S}^2] \qquad (13\text{-}5)$$

where $B_j$ is the familiar isotropic temperature factor of atom $j$. Thus Eqs. (13-4) or (13-5) is used in place of Eq. (13-1) and the atomic coordinates and temperature factors are adjusted as parameters to fit the calculated structure factors to the observed signals. Throughout the refinement process, the residual index, or $R$-factor, is computed and monitored as the measure of convergence,

$$R = \frac{\Sigma w \|F_o(hkl)| - |F_c(hkl)\|}{\Sigma w |F_o(hkl)|} \qquad (13\text{-}6)$$

where $w$ is the weighing factor associated with each reflection, $|F_o|$ is the experimentally observed structure factor and $|F_c|$ is the one calculated using either Eqs. (13-1), (13-4) or (13-5).

Crystal structure refinement at medium resolution (2.0 Å or higher) usually follows an energetically restrained refinement (EREF) formalism (Eq. 13-7, [19]) due to the insufficient number of diffractions as opposed to the number of atoms in the system, which can be conveniently combined with simulation techniques such as molecular dynamics (MD) and simulated annealing (SA). Empirical restraints are introduced as harmonic penalty functions for bond stretching, bending, and rotation in which the equilibrium values (force field parameters) are taken from statistical analyses on small molecule crystal structure database. Non-bond interactions such as van der Waals and electrostatic may be included, but it was found that the classical non-bond force field parameters often resulted in abnormally behaved MD simulations [20]. This part of the target function is termed $E_{chem}$. One such popular set of parameters for proteins was proposed by Engh and Huber [21] based on statistical analysis of the chemical moieties of proteins and polynucleotides from the Cambridge Structural Database (CSD) and is used in many crystal structure refinement packages. The second part of the target function comes from the model structure

diffraction calculation ($E_{X\text{-}ray}$) and is combined with $E_{chem}$ along with a weighting factor ($w$) to introduce $E_{X\text{-}ray}$ as an effective potential and bias refined structures toward X-ray observations. It is worthwhile to point out that if an arbitrarily large weight is assigned to $E_{X\text{-}ray}$ parameter over-fitting will be bound to happen, and the seminal work by Brunger [22] suggested setting aside a subset of X-ray data prior to refinement and monitoring the $R$ factor for this set (called $R_{\text{free}}$) as a cross-validated indicator of the progress toward convergence.

$$E = E_{chem} + wE_{X\text{-}ray} \tag{13-7}$$

One straightforward choice for the X-ray target function is the least square residual that represent the discrepancy between the observed and model-predicted structure factors:

$$E_{X\text{-}ray} = E^{LSQ} = \sum_{hkl} w(|F_o(hkl)| - k|F_c(hkl)|)^2 \tag{13-8}$$

However, as pointed out by Read and others [23–27], use of residual as the target function is only justified for models that are very close to the true structure, which is often not the case in macromolecule refinements. An improved target function can be derived using the maximum likelihood formalism (MLF),

$$E_{X\text{-}ray} = E^{MLF} = \sum_{hkl} \left(\frac{1}{\sigma_{ML}^2}\right) (|F_o(hkl)| - \langle |F_c(hkl)| \rangle)^2 \tag{13-9}$$

which is more suitable for the general case of incomplete models and models that contains initial biases. In the commonly used refinement force fields such as the one used in the CNS program [28], $E_{chem}$ is a sum of terms describing various types of interactions, i.e. [29]:

$$
\begin{aligned}
E_{chem} = E_{MM} = &\sum_{bonds} k_b(b - b_0)^2 \\
&+ \sum_{angles} k_\theta(\theta - \theta_0)^2 \\
&+ \sum_{dihedrals} k_\phi(n\phi + d) \\
&+ \sum_{chiral,planar} k_\omega(\omega - \omega_0)^2 \\
&+ \sum_{i<j} (ar_{ij}^{-12} - br_{ij}^{-6} + cr_{ij}^{-1})
\end{aligned}
\tag{13-10}
$$

where the attractive van der Waals and electrostatic terms are often omitted.

$$E = E_{QM/MM} + wE_{X-ray} \qquad (13\text{-}11)$$

In our QM flavored X-ray refinement, $E_{chem}$ is simply replaced with a quantum mechanic potential (Eq. 13-11). Compared to an MM energy function that uses fixed atomic charges to model electrostatic interactions, QM has the intrinsic advantage that it can represent charge fluctuations and dynamic polarization. In addition, a QM description is superior to an MM one when the regions of interest involve structures that differ substantially from those found in the gas phase (e.g., covalent complexes, systems with unusually close contacts, etc), where QM can represent these interactions more reliably than MM. To circumvent the great computational cost associated with electronic structure calculations of proteins, linear-scaling semiempirical QM treated was employed in our first QM X-ray refinement implementation. We have demonstrated that the QM energy restraints were capable of maintaining reasonable stereochemistry to the extent that the resultant $R$ and free $R$ values are at least comparable to those of the classical approach. The refinement is even more tractable when part of the system is treated with classical MM potentials with full electrostatic and van der Waals contributions, for example, the commonly used AMBER molecular mechanic force field descriptions. By adopting this QM/MM X-ray refinement strategy (Eq. 13-12, [30]), computation time is better spent on regions where subtle structure features are not resolved to satisfaction by the traditional refinement approach. These regions of interest may include substrates, cofactors, critical solvent molecules, and so on.

$$E_{QM/MM} = \langle \Phi_{QM} | H_{QM} + H_{QM/MM} | \Phi_{QM} \rangle + E_{MM} \qquad (13\text{-}12)$$

where $H_{QM/MM}$ is described by Eq. (13-13):

$$H_{QM/MM} = -\sum_i \sum_M \frac{Q_M}{r_{iM}} + \sum_\alpha \sum_M \frac{Z_\alpha Q_M}{R_{\alpha M}} \qquad (13\text{-}13)$$
$$+ \sum_\alpha \sum_M \epsilon_{\alpha M} \left[ \left( \frac{R_{\sigma M}^*}{R_{\alpha M}} \right)^{12} - \left( \frac{R_{\alpha M}^*}{R_{\alpha M}} \right)^{6} \right]$$

In the next section, we will demonstrate how the QM/MM refinement strategy was combined with advanced sampling techniques to resolve the ambiguous electron densities of a substrate complexed with Orf2, an aromatic prenyltransferase, and successfully clarify the relaxed substrate specificity and product regioselectivity.

## 13.3.   ORF2 AND ISOPRENOID BIOSYNTHESIS

Isoprenoids or terpenoids are a large class of naturally occurring organic compounds with tremendous chemical and structural diversity. They are organic materials produced in the HMG-CoA (3-hydroxy-3-methyl-glutaryl-CoA) reductase pathway

[31, 32] or the MEP/DOXP (2-C-methyl-D-erythritol 4-phosphate/deoxy-xylulose phosphate) pathway [33] having molecular structures containing carbon backbones made up of isoprene units. Isoprenoids are ubiquitous across a wide range of organisms, such as eubacteria, archaea, algae, plants, animals and fungi. Plant isoprenoids have been long appreciated for their "aromatic" qualities, such as those found in citral, menthol, camphor, etc., and have been commonly used in traditional herbal remedies [34]. It was recently discovered that certain isoprenoids and their derivatives exhibit desirable pharmaceutical characteristics, such as anti-microbial, anti-oxidant, anti-inflammatory, anti-viral, and anti-cancer effects [35–44]. More importantly, natural isoprenoids have low cellular toxicity and good membrane permeability, which make them ideal drug template compounds.

Several thousand natural terpenoids have been isolated and classified based on their origin, the number of isoprene units that they contain, and the roles that they play in their parent organisms. The biosynthesis and the chemical diversity of terpenoids have always been of great interests to biochemists, medicinal and organic chemists. A number of terpenoid synthase families have been categorized and extensively studied, including monoterpene, diterpene, triterpene, and sesquiterpene cyclases. The advances in the research of terpenoid biosynthesis have been summarized in a recent review article by Christianson [45], in which the structural foundation of various terpenoid cyclases was described in great detail. A common trait of these enzymes is the essential magnesium cluster and the aspartate-rich magnesium binding motifs, which implies certain similarity in their reaction mechanisms. In addition to terpenoid synthase, prenyltransferase contributes to terpenoid biosynthesis as well. One such example is Orf2 [46], a newly identified 300-residue prenyltransferase from *Streptomyces*. Although the natural substrate is unknown, Orf2 accepts a great variety of aromatic compounds, such as dihydroxynaphthalenes, flaviolin, and 4-hydroxyphenylpyruvate, and attaches geranyl ($C_{10}$) or farnesyl groups ($C_{15}$) to them (Figure 13-1). The geranylated products can then be further combined, modified, cyclized, and transferred by downstream enzymes. This certainly explains, at least in part, the apparent chemical diversity in plant and animal terpenoids. It also opens the opportunity of exploiting the relaxed substrate specificity of Orf2, possibly through protein engineering, to simplify the synthesis of terpenoids in drug discovery [47, 48]. In addition to the unusual substrate specificity, Orf2 also displays interesting regioselectivity in the prenylated products. Those of 1,6-dihydroxynaphthalene



*Figure 13-1*. The chemical structure of a prenyl group. N can be 0 (dimethylallyl), 1 (geranyl), 2 (farnesyl), 3 (geranylgeranyl), etc

*Figure 13-2.* The PT-barrel of Orf2 (shown in *red* ribbons) and the bound substrates, GSPP and 1,6-DHN (shown in licorice colored by atom types) from 1ZB6

(1,6-DHN) have been characterized by both mass spectroscopy and [1]H nuclear magnetic resonance analyses as *trans*-2-geranyl 1,6-DHN and *trans*-5-geranyl 1,6-DHN with a yield ratio of 10:1.

The 3-dimensional crystal structures of Orf2 and the Orf2 substrate complexes provided the initial clues to the observed variability in the observed reactants and products. Unlike two other members, farnesyltransferase (FTase) and geranylgeranyltransferase (GGTase), of the prenyltransferase family whose structures have been determined [49–51], Orf2 adopts a novel α/β barrel fold (termed PT-barrel by the authors, Figure 13-2) that results in a spacious binding site for geranyldiphosphate (GPP) and aromatic substrates. A magnesium ion, required for enzyme activity, is located near the diphosphate group of GPP and bound to the carboxylate of Asp62, the α-phosphate of GPP, and four water molecules. The commonly found trinuclear magnesium cluster and magnesium cluster binding motifs in terpenoid synthase were not observed in any of the crystal structures determined, which are summarized in Table 13-1. Although Orf2 is an Asp/Glu rich protein, most of the aspartate and glutamate residues are distributed on the outer α-barrel, and only two aspartate residues reside near where the magnesium and diphosphate are bound, Asp62 and Asp110. In 1ZB6, the carboxylate of Asp110 is 3.1 Å from a magnesium-bound water molecule, providing additional support for the metal binding possibly through hydrogen bond interactions. Geranyl *S*-thiolodiphosphate (GSPP), a GPP analogue, was used to allow the determination of the ternary complexes. The GSPP substitution and the binding of the aromatic substrates have little impact on the side chain conformations of the residues in the binding site. When superimposed, the all-atom root mean square deviation (RMSD) of any two Orf2 structures (residue 10–300) is no greater than 0.5 Å.

*Table 13-1.* The available crystal structures of Orf2 complexes

| PDB ID | Mean *B* value | Resolution (Å) | Bound substrate (s) |
|--------|----------------|----------------|---------------------|
| 1ZCW | 29.40 | 2.25 | geranyldiphosphate |
| 1ZDY | 16.70 | 1.44 | *N*-(tris(hydroxymethyl)methyl)-3-aminopropanesulfonic acid |
| 1ZB6 | 45.60 | 1.95 | 1,6-dihydroxynaphthalene and geranyl *S*-thiolodiphosphate |
| 1ZDW | 29.00 | 2.02 | Flaviolin and geranyl *S*-thiolodiphosphate |

We noted that 1,6-DHN and flaviolin have rather different conformations in the deposited structures (Figure 13-3). In 1ZB6, 1,6-DHN is oriented so that the C2 and C5 carbon are 4.1 and 6.7 Å from the C1 carbon of GSPP respectively, which can only explain the formation of the major product. In 1ZDW, the closest atom of flaviolin is about 6 Å away from the C1 carbon of GSPP. How the minor prenylated 1,6-DHN and prenylated flaviolin products are formed cannot be directly understood from the two crystal conformations. One possibility is that the aromatic substrates are not tightly bound and can have significant fluctuations in their active site orientation(s). The B-factors of 1,6-DHN and flaviolin were inspected as well as their electron density maps ($2F_o$–$F_c$) downloaded from the Uppsala Electron Density Server [52]. Unlike the residues that are located in the PT-barrel, both 1,6-DHN and



*Figure 13-3.* The crystal binding orientations (*left*) of 1,6-DHN and flaviolin observed in the deposited 1ZB6 and 1ZDW structures, and the 2D sketches of the prenylated 1,6-DHN (*right*). The C2 and C5 carbon atoms of 1,6-DHN and the C1 carbon atom of GSPP are shown as orange spheres

*Figure 13-4.* The crystal structure of 1ZB6 colored by *B*-factor values. The color scale (from *blue* to *white* to *red*) used has a range from 23.0 to 67.0 $Å^2$. 1,6-DHN is shown with thick cylinders with an average *B*-factor of 65.68 $Å^2$

flaviolin have much higher B-factors, on average 65.68 and 48.63 $Å^2$ respectively, similar to those of the surface residues (Figure 13-4). The electron density contour plots (Figure 13-5) at a density value of $1.5\sigma$ reveals the lack of clearly defined density for a unique 1,6-DHN conformation, which is a strong indication of substantial fluctuation in substrate binding. At the same contour level, however, flaviolin can be clearly placed with little ambiguity. As we examined the lower density contour levels in 1ZB6 (Figure 13-5), a rather large volume of space is observed and it is difficult to explain the shape of this volume by just using a single 1,6-DHN orientation.

This unusual substrate binding observed in Orf2 deserves further investigation to answer the questions raised above. To address the ambiguous electron density, multi-conformer refinement with different occupancies is typically used in MM augmented X-ray refinement packages. However, placing the different conformers in the current case can be problematic because of the poor quality of the local density. Ideally, simulated annealing MD should in principle render all plausible substrate conformations given an accurate energy function is used to describe the interactions in binding. We find this difficult to pursue as the transitions between different conformers may be hindered by large free energy barriers and the use of QM potential in X-ray refinement improves on accuracy, but also limits the amount of sampling

*Figure 13-5.* Electron density contour plots ($2F_o–F_c$) for 1,6-DHN (*left*) and flaviolin (*right*). The contour values are 1.5σ (high), 1.0σ (medium), and 0.5σ (low) from the *top* to the *bottom*

that one can afford due to a dramatic increase in computational cost comparing to the MM approach. Hence, a balanced solution is to conduct the sampling with a well parameterized classical potential, and then suggested candidates of alternative conformers are subject to QM/MM X-ray refinement. To achieve these, we explored the Orf2 binding site with a multi-nanosecond MD simulation and quantitatively characterized the binding of 1,6-DHN by computing the potential of mean force of binding as a function of two independent variables, the distances between the reaction centers. Two alternative binding conformations were identified and subsequently validated through our QM/MM X-ray refinement approach [53]. The technical details of the MD simulation and potential of mean force calculations used in this study are standard and ample examples can be found in the literature. Therefore, we will only focus on how the QM/MM X-ray refinement was carried out with our in-house integration of AMBER [54], DIVCON, and the Crystallography and NMR System (CNS) [28].

## 13.4.    QM/MM X-RAY REFINEMENT DETAILS

Unlike the native integration of DIVCON with the *sander* module in AMBER, the integration of CNS was done at the file system level. For every step of hybrid minimization, CNS, after receiving atomic coordinates from *sander*, calculates and outputs the forces as the gradient of $E_{X-ray}$ in Eq. (13-11), which is then read by *sander*. *Sander* combines the forces from $E_{X-ray}$ and $E_{QM/MM}$ potential and updates the coordinates accordingly. This process proceeds until the $R$ and free $R$ factors converge, which usually takes a few hundreds of steps.

In all the refinements carried out in this study, the complete set of 22,923 reflections between the resolution limits of 1.95 and 29.74 Å obtained from the Protein Data Bank were considered, 1,154 ($\sim$5%) of which were used for cross-validation. The alternative 1,6-DHN conformers identified from the free energy calculation were introduced into 1ZB6 by a simple superposition of the protein backbone after the original 1,6-DHN conformer and certain crystal water molecules that would cause steric clashes were removed. Both the CNS and QM/MM X-ray refinements were adapted from the minimization protocol provided by the CNS program suite. A two-stage refinement protocol was established through experimentation, in which the structure of any binding conformation was first optimized with the CNS program suite using the built-in MM energy function and the widely accepted Engh and Huber parameter set [21]. Subsequently, the results of the previous CNS calculations were used as the initial structures and were minimized with the AMBER MM energy function to relax the hydrogen atoms while keeping the heavy atoms fixed. In the following QM/MM X-ray refinement calculations, the residues chosen to be treated quantum mechanically include $Mg^{2+}$, GSPP, 1,6-DHN, and all the protein side chains and solvent molecules that are in direct contact through either hydrogen bonding, metal coordination, or van der Waals interactions. This created a neutral QM region of 666 atoms. Divide-and-Conquer [55–59], a linear-scaling semi-empirical technique, was employed to efficiently compute the QM energy at the PM3 [60, 61] level of theory. In both the CNS and QM/MM refinements, three different X-ray weighting factors (wa), 0.01, 0.2 and 1.0, were selected, allowing us to evaluate the effect of weighting factors on the quality of the final refinements.

## 13.5.    RESULTS AND DISCUSSION

In the deposited crystal structure 1ZB6, a single conformation was assigned to 1,6-DHN, which accounts for the formation of the major prenylated product, even though the deposited electron density, in our estimation, suggests that other binding conformations are also possible. The distribution of all possible binding conformations is determined by the interplay of protein–substrate, protein–solvent and substrate–solvent interactions. To better understand the Orf2 substrate selectivity and explain the observed regioselectivity in the prenylated product, we first characterized the binding landscape, and then evaluated the probability of this distribution.

The approach that we use is based on the principles of statistical mechanics and the techniques of numerical simulations.

### 13.5.1. The MD Simulation of Orf2 Ternary Complex with GPP and 1,6-DHN

To characterize and retrieve the distribution of binding of 1,6-DHN, the time evolution of the ternary complex of Orf2 with GPP and 1,6-DHN was calculated for a total duration of 16 ns starting with the model complex structure that we built from crystal structures 1ZB6 and 1ZCW. Comparing to our simulation results to 1ZCW, the backbone conformation of the protein was well preserved during the MD simulation with an average RMSD of 1.3 Å for the last 10 ns (Figure 13-6). A snapshot from the MD simulation is shown in Figure 13-7. The magnesium ion maintained a stable octahedral configuration, coordinated by four solvent molecules located in the same plane, the α-phosphate of GPP, and the side chain of Asp62. The diphosphate group of GPP strongly interacted with the side chains of several residues, including Lys119, Lys169, Arg228, Tyr216, and Lys284. Lys169 and Arg228, which are situated near the edge of the β barrel and were solvent exposed in the deposited structure. During the MD simulation, they were attracted by the negative charge on the diphosphate and formed salt bridges to further stabilize the bound GPP. The geranyl group of GPP remained mostly extended, but unlike what was observed in a previous MD simulation of the FTase ternary complex where the rotation of backbone torsion angles in farnesyldiphosphate (FPP) was seriously hindered, the backbone of GPP was more flexible.

In contrast to the relatively stable fold, the 16 ns MD simulation sampled a variety of different orientations for the bound 1,6-DHN, which we characterized using the distances D1 and D2, defined previously. The ranges of D1 and D2 sampled during the 16 ns MD simulation are plotted in Figure 13-8, which can be grouped into three clusters, an indication of three different binding states. We term these clusters/states



*Figure 13-6.* The root mean square deviation (Å) of the backbone of Orf2 over time (ps)

*Figure 13-7.* A snapshot of the MD simulation of the Orf2 ternary complex looking down the PT-barrel. The protein backbone is shown as ribbons (*red*) superimposed on top of the crystal conformation (*gray*). GPP and the magnesium ion are shown as licorice as well as 1,6-DHN, proximal protein side chains (within 3 Å) and 4 Mg-coordinating solvent molecules

as S1, S2 and S3 in the following discussion. The crystal binding conformation (shorter D1 and longer D2) is located in cluster S1. During the MD simulation that we started in cluster S1, the substrate first explored similar conformations that are within the same cluster and then migrated into cluster S2 and S3. Snapshots from the MD simulation that highlight the transitions among these binding states are provided in the supplementary material. Conformations with shorter D2 and longer D1 can be found in cluster S3, suggesting that this binding state may lead to the minor prenylated product. Cluster S2 is significantly less populated than cluster S1 and S3 and is possibly an intermediate state connecting S1 and S3.

### 13.5.2.    The Relative Free Energy Surface of 1,6-DHN Binding

To have a quantitative appreciation of these three binding states and the transitions between them, we computed the 2-dimensional potential of mean force as a function of D1 and D2 with the umbrella sampling technique followed by a WHAM analysis.

*Figure 13-8.* The distribution of D1 and D2 (in Å) sampled during the 16 ns MD simulation of the Orf2 ternary complex

Calculating free energies of conformational changes is generally difficult in biological systems due to the fact that the reaction coordinates of conformational changes are usually complicated and are not known a priori. In the current case, we were fortunate to have observed a trajectory from which the distribution of the substrate was retrieved and adequately characterized by two independent variables D1 and D2. However, the need to use two variables to discriminate different conformational states of the substrate definitely increases computational costs. The initial sampling consisted of 55 umbrella MD simulations carried out on a regular two-dimensional grid with a grid spacing of 0.7 Å and 0.8 Å in D1 and D2 respectively (Figure 13-9). However, we found that the overlap of the D1 and D2 distributions of adjacent simulations was not sufficient (data not shown) to yield a smooth PMF surface with low statistical uncertainty. One simple solution is to increase the sampling resolution by using a finer grid [62]. We first doubled the amount of sampling by placing additional umbrella points at the centers of the grid cells in Figure 13-9. This greatly reduced the statistical error in the areas of low free energies. To reduce the errors in the barrier regions and regions near the boundaries of the binding site, particularly in the region that may lead to the minor prenylated product, we further doubled the amount of sampling by placing umbrella points at the midpoints of the edges of the grid cells. Altogether over 200 umbrella MD simulations were carried out for the final PMF calculation and over 100 ns of sampling were used in the WHAM analysis.

The computed PMF surface is shown in Figure 13-10 as a 20-level contour plot. The statistical errors are less than 0.5 kcal/mol in most of the regions except the

*Figure 13-9.* The centers of sampling of the initial 55 umbrella MD simulations



*Figure 13-10.* The potential of mean force (kcal/mol) as a function of D1 and D2 (Å) shown as a 20-level contour plot and colored by the free energy values

bottom portion of cluster S3 where D2 is around 3.5 Å or less. The previously proposed binding states S1, S2, and S3 are clearly outlined in the contour plot by well-defined free energy minima. S1, the most stable binding state, is approximately 2.3 kcal/mol more favorable than S2 and S3. S1 and S2 are directly connected by a small energy barrier, which suggests S2 should be readily accessible from S1,

consistent with the frequent transitions between them observed in the 16 ns MD simulation. It may be possible to use a single variable to describe both S1 and S2 as the path connecting them is almost one dimensional. S1 and S3 are separated by a much larger and higher barrier, and the probability of direct barrier crossing (the shortest path) is very low. Hence, the low energy transition pathways between S1 and S3 must proceed through S2.

The MD snapshots corresponding to the S2 and S3 binding states were randomly picked from two umbrella points whose centers were the closest to the lowest points in the two free energy basins, and were superimposed with the crystal structure 1ZB6 and the low electron density contour (Figure 13-11). In both snapshots the naphthalene ring fits nicely into the volume enclosed by the electron density contour. Both S2 and S3 are stabilized by forming hydrogen bonds with the side chains of Ser214 and Tyr288. In the crystal structure 1ZDW, the conformation of flaviolin is similar to that of S2 (Figure 13-11) and is also hydrogen bonded to the side chains of Ser214 and Tyr288. The free energy basin where S3 is located is rather wide and flat in the direction of D2. Conformations with even shorter D2 [near (7.0, 3.5)] are only 0.4 kcal/mol higher in free energy compared to S3 and can be approached by crossing a small barrier (approximately 2.0 kcal/mol relative to S3) from S3, which may lead to the formation of the minor prenylated 1,6-DHN product. Multiple barrier crossing events captured during the 16 ns MD simulation revealed that the transitions S1↔S2 and S2↔S3 involve primarily the rotation of 1,6-DHN around the normal of the naphthalene plane accompanied by the translocation of the center of



*Figure 13-11.* The snapshots of the S2 (*green*) and S3 (*cyan*) binding states predicted from the calculated potential of mean force superimposed with 1ZB6 (*grey*) and 1ZDW (*purple*). The electron density contour of 1ZB6 at the level of 0.5σ is shown as a yellow mesh

the molecular mass. MD snapshots that highlight these transitions are provided in the supplementary material.

### 13.5.3.    The CNS and QM/MM X-Ray Refinement Calculations

To establish the optimal representation of the alternative binding conformations and critically assess the credibility of our potential of mean force study, the two selected MD snapshots were subjected to experimental validation against the observed X-ray diffraction data using a refinement protocol supplemented with a hybrid QM/MM energy function. The reported $R$ and $R_{free}$ factors for 1ZB6 are 0.233 and 0.263 respectively. The calculated $R$ and $R_{free}$ factors are compiled in Table 13-2 for the three X-ray weighting factors. Overall the two-stage QM/MM refinement protocol produced structures with comparable agreement with the X-ray diffraction data at high X-ray weighting factors. Due to the nature of the maximum likelihood refinement (MLR) algorithm, this agreement was weakened progressively when smaller weighting factors were used. However, we noted a remarkable difference between the two refinement methods in that the refinements using a more physical and realistic QM/MM energy function display less dependence on the input from the X-ray diffraction measurements. The three different conformations of the ternary complex are of similar quality, even though S1 appears to have the best agreement (lowest $R$

*Table 13-2.* The CNS and QM/MM X-ray refinement of S1, S2 and S3 conformers

| Conformers | Refinement protocol | X-ray weights | $R$ | $R_{free}$ | Distance (Å) | |
|---|---|---|---|---|---|---|
| | | | | | D1 | D2 |
| S1 | QM/MM | 0.01 | 0.2540 | 0.2674 | 3.96 | 7.09 |
| | | 0.2 | 0.2419 | 0.2629 | 3.97 | 7.12 |
| | | 1.0 | 0.2290 | 0.2628 | 4.01 | 7.21 |
| | CNS | 0.01 | 0.3735 | 0.4015 | 5.03 | 8.17 |
| | | 0.2 | 0.2606 | 0.3004 | 4.53 | 7.71 |
| | | 1.0 | 0.2307 | 0.2754 | 4.10 | 7.17 |
| S2 | QM/MM | 0.01 | 0.2604 | 0.2894 | 6.89 | 9.82 |
| | | 0.2 | 0.2432 | 0.2798 | 6.78 | 9.73 |
| | | 1.0 | 0.2285 | 0.2734 | 6.80 | 9.79 |
| | CNS | 0.01 | 0.3690 | 0.4021 | 8.15 | 10.48 |
| | | 0.2 | 0.2617 | 0.3015 | 7.53 | 10.28 |
| | | 1.0 | 0.2320 | 0.2763 | 7.10 | 10.11 |
| S3 | QM/MM | 0.01 | 0.2496 | 0.2795 | 5.91 | 4.04 |
| | | 0.2 | 0.2414 | 0.2749 | 5.81 | 3.96 |
| | | 1.0 | 0.2283 | 0.2699 | 5.87 | 3.97 |
| | CNS | 0.01 | 0.3709 | 0.4018 | 7.36 | 5.27 |
| | | 0.2 | 0.2642 | 0.3057 | 6.95 | 4.56 |
| | | 1.0 | 0.2315 | 0.2777 | 6.42 | 4.20 |

*Figure 13-12.* The snapshots of the S2 (*green*) and S3 (*cyan*) binding states from the PMF calculation (*top*) and the results from the QM/MM X-ray optimization (*bottom*) superimposed on *top* of the crystal structure 1ZB6 (*grey*) and the electron density contour at 0.5σ

and $R_{\text{free}}$ values) followed by S3 and S2. Therefore, the alternative conformations identified in the PMF calculation are both acceptable models for the Orf2 ternary complex. The distances D1 and D2 in the CNS and QM/MM refined conformations were measured and are given in Table 13-2. In general, all three refined conformations were shifted slightly away from the locations of the free energy minima on the PMF surface (Figure 13-10). This is not surprising because the potential of mean force calculation was carried out on a solvated Orf2 ternary complex with the real substrates at 300K, in which only the AMBER molecular mechanical force field was used to describe the interactions in the system, while the refined conformations correspond to local minima on the combined potential energy landscape that are shaped by a physical energy term and restraints from the X-ray diffraction measurements.

The difference in D1 and D2 between the two refinement methods is significant at smaller X-ray weighting factors, but gradually diminishes as larger weights are used. Although the detailed procedure and setup of the original CNS refinement of 1ZB6 was not reported, we estimated that an X-ray weighting factor of 1.5 is necessary to achieve the optimized $R$ and $R_{\text{free}}$ values. The fact that our QM/MM combined refinement method is less dependent on the diffraction data is potentially an important advantage over traditional methods based on simple energetic descriptions, which may be particularly well suited in applications such as evaluating possible side chain protonation states and hydrogen bonding possibilities at critical locations in system of interest. The final conformations of the QM/MM refined S2 and S3 binding states are shown in Figure 13-12.

## 13.6.    CONCLUSIONS

Unlike most enzymes that are known for being efficient, highly specific and highly selective biological catalysts, Orf2 can be considered as a "reaction chamber" for many small aromatic substrates, displaying not only reduced substrate selectivity but also interesting regioselectivity of the prenylated products. The notion of a "reaction chamber" gives rise to the opportunity to engineer the active-site of Orf2 to adjust the functionality of terpenoids for potentially novel therapeutic applications. The underlying 3-dimensional framework was revealed in a set of recently determined crystal structures; however, in this study the mechanism of substrate selectivity and product regioselectivity in Orf2 prenylation was further elucidated. We thoroughly explored the binding site of Orf2 and quantitatively evaluated the relative free energies of several binding states of 1,6-DHN in terms of a 2-dimensional potential of mean force. The deposited substrate conformation corresponds to the most stable binding state, but the others are certainly accessible at 300K. Given the moderate free energy barriers that separate them, it would not be unreasonable to assume that the substrate binds in a single conformation and rapidly equilibrates S1 and S3 (possibly via S2). Consequently, the observed substrate electron density is the thermodynamic average of the contributions of all these binding states. The binding site of Orf2 is narrow in one dimension, which is consistent with the observed substrate preference, but spacious in the others to admit aromatic molecules of different sizes, possibly up to three or four fused rings. The conformations of the other two binding states that we identified and validated against experimental X-ray diffraction data not only explained the electron density profile at low contour values and the high B-factors of the substrate, but also offered a plausible interpretation of the regioselectivity in the prenylated products. The product yields are dependent upon the thermodynamics of binding, the incubation time, and the relative reactivities. The thermodynamics of substrate binding is determined by the interplay of substrate–protein, protein–solvent and substrate–solvent interactions. For 1,6-DHN, a 10:1 product ratio (S1-product vs. S3-product) was observed under the conditions that the concentration of the substrate was ~160 fold in excess of the enzyme and both were incubated at 298K

for 4 h. We estimate that the concentration ratio of the two competing conformations is close to 100:1 from a 2.7 kcal/mol free energy difference between S1 and S3, assuming a constant concentration of the ternary complex. Based on the Curtin-Hammett principle, we speculate that the prenylation of S3 is a faster reaction with a lower reaction barrier height relative to that of S1. Ser214 and Tyr288 are likely the moderators of substrate binding, the mutations of which may affect the selectivity and regioselectivity, therefore are good candidates for site-directed mutagenesis experiments. Flaviolin is bound similarly to the S2 state of 1,6-DHN, also hydrogen bonded to the side chains of Ser214 and Tyr288. However, the electron density of 1ZDW is better defined around flaviolin, indicating a single preferred conformation. This is possibly due to the fact that an extra hydrogen bond can be formed between the carbonyl group of flaviolin and Tyr288, which favors this orientation and reduces conformational fluctuation.

The credibility of the results from computational enzyme mechanism studies rely on high quality structural information as starting points. Extra caution must be taken in thoroughly scrutinizing the structure details of the models offered by the initial refinement effort before any further computational effort is spent. The pivot of focus is best laid on the alignment of participating residues, substrates, or cofactors as these residues are likely to have unusual geometries and correspondingly poorly treated by the inadequate classical force field approach during X-ray refinement process. The sanity status of crystal structures is also important to other structure-based research, such as structure-based drug design and enzyme engineering where visual examination and extrapolation are frequently used. We have demonstrated, previously and now, that QM-based hybrid refinement approach can provide the crucial structure subtleties that may not be available from the conventional X-ray refinement, particularly when the model quality is undermined by the X-ray diffraction measurements. The strength of our QM flavored X-ray refinement approach at its current stage is not to offer significant improvement on $R$ and free $R$ factors of the refined model, which are often used as a primary quality indicator. $R$ and free $R$ factors reflect the overall consistency of a model with X-ray diffraction, contributed by all atoms in the system, and are unlikely to be affected to a great extent by local optimization. In the example presented here, we observed slightly lowered values for $R$ and free $R$ factors in all three different conformers contributed by the 666 QM treated atoms out of 2571. Perhaps more importantly, the distances, D1 in conformer S1 and D2 in conformer S3, are noticeably shorter in the QM/MM refined models, which is more intuitive and appealing from the perspective of enzymatic catalysis.

Incorporating quantum mechanical potential into the final refinement stage of crystal structure determination substantially improves the accuracy of energetics evaluation and thus offers an approach to probe subtle structural features unattainable solely from diffraction data due to experimental uncertainty. This can be seen from the response of $R$ and $R_{free}$ factors upon changes in X-ray weights. The deterioration in $R$ and $R_{free}$ due to the use of smaller weights is less pronounced in QM/MM refinement calculations (Table 13-2), which signifies a great opportunity to explore further. By using smaller weights, the X-ray bias is reduced, which increases

the chance to differentiate the proper protonation or tautomeric states by using an adequate quantum mechanical energy function.

Despite all the benefits that a QM-based X-ray refinement approach promises, we should be aware of the shortcomings that accompany the use of quantum mechanical potentials, including the increased computational cost, thereby, elongated refinement cycles, the limited sampling efficiency, and the computational instabilities (aka failed SCF convergence) due to distorted structural elements in the initial geometry. We expect that the conventional MM-based refinement approach will likely continue to be the first choice of crystallographers in the near future due to its robustness and speed especially in the early stages of refinement. However, QM-based refinement approaches, while complementary to current ones, should provide significant structural improvements in the latter stages of the refinement process, when non-standard co-factors are present, and when subtle but important structure features are poorly resolved by conventional methods.

## ACKNOWLEDGMENTS

## REFERENCES

1. Banaszak LJ, (2000) Foundation of structural biology, Academic Press, San Diego
2. Lesk AM, (2001) Introduction to protein architecture: the structural biology of proteins,Oxford University Press, Oxford
3. Marti-Renom MA et al (2000) Comparative protein structure modeling of genes and genomes. Annu Rev Biophys Biomol Struct 29:291–325
4. Bonneau R, Baker D (2001) Ab initio protein structure prediction: progress and prospects. Annu Rev Biophys Biomol Struct 30:173–189
5. Baker D, Sali A (2001) Protein structure prediction and structural genomics. Science 294(5540):93–96
6. Schueler-Furman O et al (2005) Progress on modeling of protein structures and interactions. Science 310(5748):638–642
7. Skolnick J (2006) In quest of an empirical potential for protein structure prediction. Curr Opin Struct Biol 16(2):166–171
8. Kuhlman B et al (2003) Design of a novel globular protein fold with atomic-level accuracy. Science 302(5649):1364–1368
9. Bradley P, Misura KMS, Baker D (2005) Toward high-resolution de novo structure prediction for small proteins. Science 309(5742):1868–1871
10. Brunger AT (1988) Crystallographic refinement by simulated annealing. Application to a 2.8 Å resolution structure of aspartate aminotransferase. J Mol Biol 203(3):803–816
11. Adams PD et al (1997) Cross-validated maximum likelihood enhances crystallographic simulated annealing refinement. Proc Natl Acad Sci USA 94(10):5018–5023
12. Yu N, Yennawar H, Merz KM Jr (2005) Refinement of protein crystal structures using energy restraints derived from linear-scaling quantum mechanics. Acta Crystallogr D Biol Crystallogr 61:322–332

13. Yu N, Merz KM Jr (2004) Theoretical study of the electron density distributions of glycyl-L-threonine dihydrate. Molecular Physics 102(23–24):2545–2557

14. Ryde U, Nilsson K (2003) Quantum chemistry can locally improve protein crystal structures. J Am Chem Soc 125(47):14232–14233

15. Yu N et al (2006) Critical assessment of quantum mechanics based energy restraints in protein crystal structure refinement. Protein Sci 15(12):2773–2784

16. Nilsson K et al (2004) The protonation status of compound II in myoglobin, studied by a combination of experimental data and quantum chemical calculations: quantum refinement. Biophys J 87(5):3437–3447

17. Nilsson K, Ryde U (2004) Protonation status of metal-bound ligands can be determined by quantum refinement. J Inorg Biochem 98(9):1539–1546

18. Yu N et al (2006) Assigning the protonation states of the key aspartates in beta-secretase using QM/MM X-ray structure refinement. J Chem Theory Comput 2(4):1057–1069

19. Jack A, Levitt M (1978) Refinement of large structures by simultaneous minimization of energy and *R* factor. Acta Crystallogr A 34:931–935

20. Brunger AT, Adams PD (2002) Molecular dynamics applied to X-ray structure refinement. Acc Chem Res 35(6):404–412

21. Engh R, Huber R (1991) Accurate bond and angle parameters for X-ray protein-structure refinement. Acta Crystallogr A 47:392–400

22. Brunger AT (1992) Free *R*-value – a novel statistical quantity for assessing the accuracy of crystal-structures. Nature 355(6359):472–475

23. Read RJ (1986) Improved fourier coefficients for maps using phases from partial structures with errors. Acta Crystallogr A 42:140–149

24. Read RJ (1990) Structure-factor probabilities for related structures. Acta Crystallogr A 46:900–912

25. Pannu NS, Read RJ (1996) Improved structure refinement through maximum likelihood. Acta Crystallogr A 52:659–668

26. Adams PD et al (1997) Cross-validated maximum likelihood enhances crystallographic simulated annealing refinement. Proc Natl Acad Sci USA 94(10):5018–5023

27. Brunger AT, Adams PD (2002) Molecular dynamics applied to X-ray structure refinement. Accounts of Chemical Research 35(6):404–412

28. Brunger AT, Adams PD, Clore GM, Delano WL, Gros P, Grosse-Kunstleve RW, Jiang JS, Kuszewski J, Nilges M, Pannu NS, Read RJ, Rice LM, Simonson T, Warren GL (1998) Acta Crystallogr D 54:905–921

29. Brunger AT, Karplus M, Petsko GA (1989) Crystallographic refinement by simulated annealing – application to crambin. Acta Crystallogr A 45:50–61

30. Warshel A, Levitt M, (1976) Theoretical studies of enzymic reactions: dielectric, electrostatic and steric stabilization of the carbonium ion in the reaction of lysozyme. J Mol Biol 103:227–249

31. Spurgeon SL, Porter JW (1981) In biosynthesis of isoprenoid compounds. In: Porter JW, Spurgeon SL (eds) John Wiley and Sons, New York, p 1

32. Qureshi N, Spurgeon SL (1981) In biosynthesis of isoprenoid compounds. In: Porter JW, Spurgeon SL (eds) John Wiley and Sons, New York, p 47

33. Rohmer M (1999) The discovery of a mevalonate-independent pathway for isoprenoid biosynthesis in bacteria, algae and higher plants. Nat Prod Rep 16(5):565–574

34. Piironen V et al (2000) Plant sterols: biosynthesis, biological function and their importance to human nutrition. J Sci Food Agric 80(7):939–966

35. Santos FA, Rao VSN (1998) Inflammatory edema induced by 1,8-cineole in the hindpaw of rats: a model for screening antiallergic and anti-inflammatory compounds. Phytomedicine 5(2):115–119

36. Blanco-Colio LM et al (2003) Anti-inflammatory and immunomodulatory effects of statins. Kidney Int 63(1):12–23
37. Grosser N et al (2004) The antioxidant defense protein heme oxygenase 1 is a novel target for statins in endothelial cells. Free Radic Biol Med 37(12):2064–2071
38. Chowdhury SA et al (2005) Tumor-specificity and apoptosis-inducing activity of stilbenes and flavonoids. Anticancer Res 25(3B):2055–2063
39. Jahangir T et al (2005) Alleviation of free radical mediated oxidative and genotoxic effects of cadmium by farnesol in Swiss albino mice. Redox Rep 10(6):303–310
40. Soria-Mercado IE et al (2005) Antibiotic terpenoid chloro-dihydroquinones from a new marine actinomycete. J Nat Prod 68(6):904–910
41. Zhou YD et al (2005) Terpenoid tetrahydroisoquinoline alkaloids emetine, klugine, and isocephaeline inhibit the activation of hypoxia-inducible factor-1 in breast tumor cells. J Nat Prod 68(6):947–950
42. Boucher K et al (2006) HMG-coa reductase inhibitors induce apoptosis in pericytes. Microvasc Res 71(2):91–102
43. Hwang DR et al (2006) Synthesis and anti-viral activity of a series of sesquiterpene lactones and analogues in the subgenomic HCV replicon system. Bioorg Med Chem 14(1):83–91
44. Jahangir T et al (2006) Farnesol prevents Fe-NTA-mediated renal oxidative stress and early tumour promotion markers in rats. Hum Exp Toxicol 25(5):235–242
45. Christianson DW (2006) Structural biology and chemistry of the terpenoid cyclases. Chem Rev 106(8):3412–3442
46. Kuzuyama T, Noel JP, Richard SB (2005) Structural basis for the promiscuous biosynthetic prenylation of aromatic natural products. Nature 435(7044):983–987
47. Botta B et al (2005) Novel prenyltransferase enzymes as a tool for flavonoid prenylation. Trends Pharmacol Sci 26(12):606–608
48. Koehl P Relaxed specificity in aromatic prenyltransferases. Nat Chem Biol 1(2):71–72
49. Taylor JS et al (2003) Structure of mammalian protein geranylgeranyltransferase type-I. EMBO J 22(22):5963–5974
50. Zhang H, Seabra MC, Deisenhofer J (2000) Crystal structure of Rab geranylgeranyltransferase at 2.0 angstrom resolution. Structure Fold Des 8(3):241–251
51. Park HW et al (1997) Crystal structure of protein farnesyltransferase at 2.25 angstrom resolution. Science 275(5307):1800–1804
52. Kleywegt G et al (2004) The uppsala electron-density server. Acta Crystallogr D Biol Crystallogr 60:2240–2249
53. Yu N, Hayik SA, Wang B, Liao N, Reynolds CH, Merz KM Jr (2006) Assigning the protonation states of the key aspartates in beta-secretase using QM/MM X-ray structure refinement. J Chem Theory Comput (Web release, June 7)
54. Case DA, Darden TA, Cheatham TE III, Simmerling CL, Wang J, Duke RE, Luo R, Merz KM, Wang B, Pearlman DA, Crowley M, Brozell S, Tsui V, Gohlke H, Mongan J, Hornak V, Cui G, Beroza P, Schafmeister C, Caldwell JW, Ross WS, Kollman PA (2004) AMBER 8
55. Yang W (1991) Direct calculation of electron-density in density-functional theory. Phys Rev Lett 66(11):1438–1441
56. Yang W, Lee T (1995) A density-matrix divide-and-conquer approach for electronic-structure calculations of large molecules. J Chem Phys 103(13):5674–5678
57. Dixon SL, Merz KM Jr (1996) Semiempirical molecular orbital calculations with linear system size scaling. J Chem Phys 104(17):6643–6649

58. Lee T, York D, Yang W (1996) Linear-scaling semiempirical quantum calculations for macro-molecules. J Chem Phys 105(7):2744–2750
59. Dixon SL, Merz KM Jr (1997) Fast, accurate semiempirical molecular orbital calculations for macromolecules. J Chem Phys 107(3):879–893
60. Stewart J (1989) Optimization of parameters for semiempirical methods. 1: method. J Comput Chem 10(2):209–220
61. Stewart J (1989) Optimization of parameters for semiempirical methods. 2: applications. J Comput Chem 10(2):221–264
62. Roux B (1995) The calculation of the potential of mean force using computer simulations. Comput Phys Comm 91:275–282

# CHAPTER 14

# UNRAVELING THE MECHANISMS OF RIBOZYME CATALYSIS WITH MULTISCALE SIMULATIONS

TAI-SUNG LEE[1,2], GEORGE M. GIAMBAŞU[1,2], ADAM MOSER[2], KWANGHO NAM[3], CARLOS SILVA-LOPEZ[2], FRANCESCA GUERRA[2], OLALLA NIETO-FAZA[2], TIMOTHY J. GIESE[1], JIALI GAO[2], AND DARRIN M. YORK[1,2]

[1]*Biomedical Informatics and Computational Biology, University of Minnesota, Minneapolis, MN 55455, USA*

[2]*Department of Chemistry, University of Minnesota, 207 Pleasant St. SE, Minneapolis, MN 55455, USA, e-mail: york@umn.edu (D.M. York)*

[3]*Department of Chemistry and Chemical Biology, Harvard University, Cambridge, MA 02138, USA*

**Abstract:** Description of a multiscale simulation strategy we have developed to attack problems of RNA catalysis is presented. Ribozyme systems give special challenges not present in typical protein systems, and consequently demand new methods. The main methodological components are herein summarized, including the assembly of the **QCRNA** database, parameterization of the AM1/d-PhoT Hamiltonian, and development of new semiempirical functional forms for improved charge-dependent response properties, methods for coupling many-body exchange, correlation and dispersion into the QM/MM interaction, and generalized methods for linear-scaling electrostatics, solvation and solvent boundary potentials. Results for a series of case studies ranging from noncatalytic reaction models that compare the effect of new DFT functionals, and on catalytic RNA systems including the hairpin, hammerhead and L1 ligase ribozymes are discussed

**Keywords:** Ribozyme catalysis, multiscale simulation, linear-scaling method, QM/MM, DFT functional

## 14.1. INTRODUCTION

Over the last several decades, the original notion that the only function of RNA molecules was as messenger intermediates in the pathway from the genetic code to protein synthesis has undergone a revolution. The role of RNA in cellular function is now known to be considerably more diverse, ranging from regulation of gene expression and signalling pathways to catalyze important biochemical reactions, including protein synthesis itself [1–7]. These discoveries have transformed our view of RNA

as a simple messenger to one more profoundly central in the evolution of life forms, our understanding and appreciation of which is still in its infancy. Ultimately, the elucidation of the mechanisms of RNA catalysis will yield a wealth of new insights that will extend our understanding of biological processes and facilitate the design of new RNA-based technologies [8–10].

Molecular simulations of RNA catalysis, in principle, offer a means of accessing the most intimate mechanistic details that may aid in the interpretation of experiments and provide predictive insight into design [11]. In order to study reactions catalyzed by biological macromolecules such as RNA, simulations are performed using so-called multiscale models. Here, by "multiscale model", we mean the integration of a hierarchy of models that work together to provide a computationally tractable representation of a complex biochemical reaction in a realistic environment. As a specific example, for enzyme systems, one typically treats the reactive chemical events with a sufficiently accurate quantum mechanical model, the microscopic solvent fluctuations and changes in molecular conformation using molecular mechanical force field model, and the macroscopic dielectric relaxation using a continuum solvation model. The simplest multiscale model to study enzyme reactions would be a combined quantum mechanical/molecular mechanical (QM/MM) potential [12–18].

Simulations of RNA enzymes, or *ribozymes*, however, are laden with challenges not apparent for most protein enzymes. RNA molecules are highly negatively charged, and exhibit strong, and often specific interactions with solvent [11, 19–21]. This requires special attention to the microscopic in silico model that requires consideration of a very large number of solvent molecules and counter and co-ions to be included. Electrostatic interactions need to be treated rigorously without cut-off, and long simulation times are typically needed to insure that the ion environment is properly equilibrated [22–24]. These issues are further complicated by the fact that RNA molecules bind divalent metal ions which play an important role in folding, and in many instances, also contribute actively to the catalytic chemical steps. The highly-charged nature of RNA and its interaction with divalent metal ions and other solvent components makes inclusion of explicit electronic polarization in the molecular models much more important that in typical protein enzyme systems. The chemistry involved in reactions of prototype ribozymes such as cleavage transesterification involves large changes in local charge state and hybridization around phosphorus, exacerbating the need to design QM/MM methods that can reliably model hypervalent states of phosphorus. There is a need to design new models that circumvent the need for "atom-type" parameters to be assigned to the QM system in order to compute QM/MM interactions, as the "atom-type" can change as a reaction proceeds. Finally, there is a growing precedent that many ribozyme reactions may involve large changes in conformation and metal ion binding along the reaction coordinate, creating the need to develop extremely fast semiempirical quantum models that can be practically applied in conjunction with long-time simulations to adequately sample relevant configurations and create multidimensional free energy surfaces along multiple reaction coordinates.

## 14.2. COMPUTATIONAL APPROACH

This section describes the main methodological advances that will be used in subsequent selected applications, including (1) Development of fast semiempirical methods for multiscale quantum simulations, (2) Directions for development of next-generation QM/MM models, and (3) Linear-scaling electrostatic and generalized solvent boundary methods.

### 14.2.1. Development of Fast Semiempirical Methods for Multiscale Quantum Simulations

The development of new-generation Hamiltonians is greatly benefited by the use and refinement of existing models. Conventional semiempirical Hamiltonians such as MNDO [25, 26], AM1 [27] and PM3 [28] are based on a minimal *s* and *p* orbital basis, and were parametrized to reproduce gas phase geometries and heats of formation for molecules in their ground states. Relatively little consideration was given to models that could accurately describe chemical reactions, especially in the case of phosphate hydrolysis where *d* orbitals are required for an accurate representation of hypervalent phosphorus transition states and intermediates. The *d*-orbital extension of the MNDO method [29, 30] greatly improved the description of hypervalent phosphorus species, but suffered from an extremely poor description of hydrogen bonding as did the original MNDO model. The empirical core-core functions of AM1 and PM3 allowed these models to reasonably predict hydrogen bonding, and if appropriately reparametrized, were demonstrated to be quite accurate for nucleic acid base pairing [31]. Consequently, it seemed that a functional form that merged the core-core functions of AM1 and PM3 with the *d*-orbital formulation of MNDO/d would offer a significantly improved base model that could be easily implemented into any *d*-orbital semiempirical program and thus be readily available to a wide scientific community. However, a prerequisite to performing such a reparameterization is to assemble a consistent set of high-level reference data that considers a wide range of properties and encompasses a broad set of molecules, complexes, potential energy surfaces and chemical reaction mechanisms. This was achieved by the construction of a database of quantum calculations for RNA catalysis (*QCRNA*) based on density-functional theory [32]. Subsequently, the *QCRNA* database was utilized to parameterize an AM1/d-PhoT [33] semiempirical Hamiltonian model for a wide range of phosphoryl transfer reactions and hydrogen bonding interactions. The *QCRNA* database and AM1/d-PhoT model are described in the following subsections, and together provide powerful tools to further the understanding of the mechanisms of ribozyme catalysis.

#### 14.2.1.1. QCRNA*: A Database of Quantum Calculations for RNA Catalysis*

The *QCRNA* database [32, 34] is an on-line resource of ab initio data relevant to phosphorus and RNA chemistry, and was specifically designed to act as the reference

data for the parameterization of new semiempirical models and refinement of existing semiempirical Hamiltonians. The database utilizes a strict computational protocol to determine optimized molecular structures, electronic structure properties, and thermodynamic quantities, including estimates of the solvation free energy and solvent-polarization using continuum models. For brevity we describe the protocol as B3LYP/6-311++G(3df,2p)//B3LYP/6-31++G(d,p); while a detailed explanation of the precise protocol is provided in Ref. 32. One of the database's advantages is that all of the data is obtained using the same computational procedure, whereas the literature data has often been inconsistent, in that disparate theoretical protocols have been employed by different groups that severely limit meaningful quantitative cross-comparison.

There are more than 2000 molecules, 300 molecular complexes, 250 chemical reactions, and 50 potential energy surfaces in the database. The main types of information that can be found in the database include: molecular geometries, electronic structure properties, vibrational frequencies, relative conformational energies, hydrogen bond energies, metal ion binding energies, proton affinities/gas-phase basicities, tautomerization energies, and chemical mechanisms. Many of molecules in the database correspond to stationary points, i.e., local minima and maxima, of chemical mechanisms; however there also exist simple potential energy surfaces involving non-stationary geometries resulting from constrained optimization. Of the mechanisms relevant to RNA catalysis are those that involve: acyclic and cyclic phosphates, phosphate mono-, di- and triester systems, different protonation and charge states, experimentally relevant thio effects, phosphorane pseudorotation reactions, metal-catalyzed reactions, and linear free energy relations that involve different nucleophiles and leaving groups. Data contained in the *QCRNA* database has been used as the basis for numerous publications on biological phosphates and phosphoranes, and phosphoryl transfer reactions [35–46].

The *QCRNA* database is viewable and searchable with a web browser on the internet and it is also contained as a MySQL database that is easily incorporated with parameter optimization software to allow for the rapid development of specific reaction parameters. Molecular structures can be viewed with the JMOL [47, 48] or MOLDEN [49, 50] programs as viewers for chemical MIME types. If the web browser is JAVA-enabled, then the JMOL software will automatically load as a web applet. Both programs allow the structure to be manipulated, i.e., rotated, scaled, and translated, and allow for measurement of internal coordinates, e.g., bond lengths, angles, and dihedral angles. Similarly, animations of the vibrational frequencies are available and can be viewed with either program.

### 14.2.1.2.  *AM1/d-PhoT Model for Phosphoryl Transfer Reactions*

The AM1/d-PhoT model [33] is a parameterization of a modified AM1/d Hamiltonian developed specifically to model phosphoryl transfer reactions catalyzed by enzymes and ribozymes for use in linear-scaling calculations and combined QM/MM simulations. The model is currently parametrized for H, O, and P atoms to reproduce

high-level density-functional results from the *QCRNA* database [32, 34], including geometries and relative energies of minima, transition states and reactive intermediates, dipole moments, proton affinities and other relevant properties. The model has been tested in the gas phase and in solution using molecular dynamics simulations with a QM/MM potential [33, 51]. The results indicate the method provides significantly higher accuracy than MNDO/d, AM1 and PM3 methods, and for the transphosphorylation reactions is in close agreement with the density-functional calculations at the B3LYP/6–311++G(3df,2p) level with a reduction in computational cost of 3–4 orders of magnitude. The model has recently been applied in QM/MM simulations of the hammerhead [52, 53] and hairpin ribozyme [54, 55] systems.

### 14.2.2.    Directions for Development of Next-Generation QM/MM Models

#### 14.2.2.1.    *Semiempirical Methods with Improved Charge-Dependent Response Properties*

Standard NDDO, MNDO, and tight binding semiempirical models systematically underestimate the polarizability of molecules, and this has been attributed mainly to the use of minimal basis sets [56]. The underestimation of the polarizability is most pronounced in the pathological case where all of the spin orbitals are occupied, such as $F^-$, in which case a minimal basis calculation is devoid of degrees of freedom and thus lacks the variational parameters to describe polarization. Molecules with physically relevant electronic configurations rarely approach the extreme limit of total occupation of the basis; however, the polarizability of neutral molecules has been observed to be underestimated by approximately 25% when comparing semiempirical calculations with experiment [57]. A more recent and much larger comparison between MNDO/d and B3LYP/6-31++G(3df,2p) polarizabilities for 1132 molecules taken from the *QCRNA* [32] online database indicates that the semiempirical polarizability can be underestimated by approximately 40% [56]. An obvious solution for correcting the polarizabilities is to use a larger orbital basis set; however, this has the disadvantage of increasing the dimensions of the Fock matrices and requires additional semiempirical parametrization. Giese and York [56] (GY) took the alternative approach of including an auxiliary basis of density which explicitly polarizes in response to the external potential. This approach is based on the principle of chemical potential equalization (CPE), which is derivable as a second order Taylor series expansion of the density functional in density response. GY opted to incorporate the CPE model as a post-SCF correction to the semiempirical method, whereby the semiempirical SCF density is used as the reference density in CPE expansion. The auxiliary basis used to describe the density response was chosen to be atom-centered primitive Gaussian dipole functions and the integrals produced from the CPE expansion where performed using a Coulomb approximation.

  GY noticed that the polarizability of an atom is approximately an exponential function of charge, and that the polarizability correction provided by the CPE expansion for an isolated atom was equal to the inverse of the Coulomb self energy of the

Gaussian dipole. From this, they developed an expression for the Gaussian exponent as a function of charge, thereby allowing for a charge dependent polarizability correction that directly results from the dilation of the response density. The resulting model (MNDO/d+CPE) was shown to reduce the error in the polarizability of the 1132 molecules down to 3.6% with a mean unsigned error of 0.1%. Furthermore, GY showed that the charge dependence form of the response basis allowed them to reproduce the polarizability of molecules in various charge states.

### 14.2.2.2.    *QM/MM Interactions with Coupled Many-Body Polarization, Exchange, and Dispersion*

The most commonly applied QM/MM methods utilize a quantum mechanical model combined with an empirical force field which are coupled within the QM Hamiltonian and through a QM/MM interaction energy consisting of electrostatic, bonded, and nonbonded terms. The nonbonded QM/MM interaction is modeled by a simple empirical van der Waals interaction such as a Lennard-Jones 6–12 potential, whose empirical parameters are based on the assignment of atoms to "atom-types". This term is purely empirical and completely neglects explicit coupling to the quantum mechanical electronic degrees of freedom. The lack of explicit quantum mechanical coupling of the van der Waals QM/MM interaction energy is at the root of many problems in QM/MM modeling. Unlike MM atoms, the association of a QM atom to a particular atom-type becomes ambiguous when the QM atom undergoes changes in charge state and/or chemical bonding environment during passage through a reaction coordinate or perturbation parameter. The profiles of reactions involving highly charged species, such as those encountered in phosphate hydrolysis and phosphoryl transfer reactions, can be skewed by the solvent effect if not properly treated. The solvent effect is sensitive to the van der Waals radii used, since these dictate the degree to which solvent can approach ionic substrates.

Giese and York [58] (GY) modified their MNDO/d+CPE model [56] for QM/MM interactions to explicitly treat the charge dependence of the van der Waals forces. In their approach, called OPNQ, the van der Waals correction was applied as a post-SCF correction to the QM/MM energy; however, their approach can easily be incorporated within the SCF procedure. The OPNQ model consists of charge dependent repulsive and dispersion components. The repulsion energy is based on a density overlap model which was motivated in form from the observation that the Hartree-Fock repulsion of rare gas dimers could be reproduced from a parametrized overlap of the unperturbed atomic densities [59, 60]. In the OPNQ model, the atomic density is modeled by an atom-centered spherical Slater function whose exponent was optimized to reproduce the homonuclear dimer density overlap as a function of separation. The Slater exponent is then allowed to vary exponentially with respect to charge, and this introduces an empirical parameter to describe charge dependence. The dispersion model is a traditional multipole expression involving $C_6$ and $C_8$ dispersion coefficients which are damped in the short range by the Tang and Toennies function [61, 62]. The expressions for the dispersion coefficients were taken from the

work of Pellenq and Nicholson [63, 64] (PN), whom developed the equations from perturbation theory. The PN equations give heteronuclear dispersion coefficients using 1-body parameters only. GY modified the PN equations, which depend on such things as the dipole polarizability of the atoms, to include their charge dependence.

In their work [58], GY demonstrated that a standard Lennard-Jones model grossly over-predicted the well-depth of rare gas-halide ion dimer potential energy curves when they were parametrized to reproduce the neutral rare gas-halide dimer curves. They further showed that the OPNQ model performed just as badly when the charge dependence of the expressions were ignored, but the potential energy curves for both the neutral and ionic dimers could be simultaneously be reproduced if the charge dependence is considered.

### 14.2.3. Linear Scaling Electrostatic and Generalized Solvent Boundary Methods

#### 14.2.3.1. Linear Scaling Electrostatics

Linear scaling evaluation of the energy and forces is a prerequisite to the application of the new and improved model Hamiltonians described in the previous sections to simulations of large biomolecules. There are several methodological bottlenecks inherent within ab initio methods, such as diagonalization of the Fock matrix, that prevent linear scaling. The divide-and-conquer algorithm is one method for overcoming the diagonalization problem [65, 66]; however, we here assume that the quantum region is small and relatively independent of the size of the system and that the remainder of the system is composed of molecular mechanical atoms. In this case, only the long-range electrostatic interactions pose a major obstacle to achieving linear scaling. The two main approaches for overcoming this obstacle are Ewald methods for periodic boundary simulations and fast multipole methods for stochastic boundary simulations. The theory of these approaches are fairly standard and are not repeated here. Instead we outline some of our recent contributions in these areas below.

Recently, Nam, Gao and York (NGY) [67] reported a linear scaling semiempirical QM/MM Ewald method and its incorporation into the CHARMM simulation package. In order to take advantage of the optimized Ewald algorithms traditionally used for static point charge distributions, they performed a Mulliken partitioning of the electron density and treated the QM region as a series of point charges. The QM charges are not static, however, and therefore they also developed Fock matrix corrections to obtain a self-consistent wavefunction, which is required to compute the analytic gradients. NGY also examined how much error is introduced into the gradients when the Ewald Fock matrix correction was independent of the SCF cycles.

Giese and York (GY) [68] used the branch-free FMM algorithm of Watson et al. [69] and the recursive bisection ideas of Perez-Jorda and Yang (PJY) [70] to create an adaptive FMM for systems of particles composed of point multipoles, as opposed to the trivial case of point charges (monopoles). GY spent most of their effort in

developing an adaptive termination criteria, i.e., when to stop splitting the system, and evaluating the stability of the adaptive divisioning of the system. More specifically, GY chose to split the system either by: (1) a dividing plane in which the child subsystems lie on either side of the plane, i.e., "fluc-splitting" or (2) dividing the Cartesian rectangular box enclosing the system perpendicular to its largest edge, i.e., "box-splitting". Fluc-splitting gets its name from having chosen the location and orientation of the dividing plane from the center of distribution and the eigenvector of the largest eigenvalue from the $3 \times 3$ covariance matrix of the particle distribution, respectively. GY found that box-splitting was superior to fluc-splitting because it did not exhibit errors associated with the creation of subsystem shapes that are only adequately modelled accurately with a large multipole expansion. They also found that an adaptive termination criteria can be constructed to produce near-optimal performance for systems composed of point charges or point multipoles, and for small and large systems of various shapes.

### 14.2.3.2.    *Generalized Solvation and Solvent Boundary Methods*

In some instances, it is not computationally feasible to treat all of the water molecules explicitly such as in a large simulation cell with periodic boundary conditions. In such cases, recourse must be taken into alternative methods to treat the generalized solvation effects that integrate out the explicit degrees of freedom of the solvent using a continuum or linear-response approach. In this section, we briefly describe the current status of development of such methods in our group, including a smooth COSMO solvation model and a variational electrostatic projection (VEP) method for generalized solvation effects.

*Smooth COSMO solvation model.* We have recently extended our smooth COSMO solvation model with analytical gradients [71] to work with semiempirical QM and QM/MM methods within the CHARMM and MNDO programs [72, 73]. The method is a considerably more stable implementation of the conventional COSMO method for geometry optimizations, transition state searches and potential energy surfaces [72]. The method was applied to study dissociative phosphoryl transfer reactions [40], and native and thio-substituted transphosphorylation reactions [73] and compared with density-functional and hybrid QM/MM calculation results. The smooth COSMO method can be formulated as a linear-scaling Green's function approach [72] and was applied to ascertain the contribution of phosphate-phosphate repulsions in linear and bent-form DNA models based on the crystallographic structure of a full turn of DNA in a nucleosome core particle [74].

*Variational electrostatic projection method.* In some instances, the calculation of PMF profiles in multiple dimensions for complex chemical reactions might not be feasible using full periodic simulation with explicit waters and ions even with the linear-scaling QM/MM-Ewald method [67]. To remedy this, we have developed a variational electrostatic projection (VEP) method [75] to use as a generalized solvent boundary potential in QM/MM simulations with stochastic boundaries. The method is similar in spirit to that of Roux and co-workers [76–78], which has been recently

*Figure 14-1. Left:* Relative errors (RELE) in the force as a function of radial distance from the center of the active dynamical region for the VEP-RVM charge-scaling method [80] for the solvated hammerhead ribozyme at different discretization levels [151] of the $\omega$ surface. *Right:* The projected total electrostatic potential due to the fully solvated hammerhead ribozyme projected onto the VEP surface [80]

implemented into QM/MM simulations [79]. We have also developed a charge-scaling implementation of the method [80] that delivers high accuracy (Figure 14-1). Preliminary results suggest that the VEP method is more general and considerably more accurate than methods based on multipole expansions (Figure 14-1).

## 14.3.    SELECTED APPLICATIONS

In this section, we describe a sampling of applications that target different facets of the problems associated with RNA catalysis. The applications range from the study of small model phosphoryl transfer reactions in solution to chemical reactions, metal ion binding and conformational events that occur in ribozyme systems. First, phosphoryl transfer reactions in solution are examined, with the underlying goal of determining reliable model chemistries that capture the essential features of the reaction profile as characterized experimentally. Second, applications to ribozyme are explored that examine different aspects of RNA catalysis. The chemical steps of catalysis are explored in the hairpin ribozyme, a prototype ribozyme that does not have an explicit catalytic divalent metal ion requirement. Next the role of divalent metal ions are explored in a similar reaction catalyzed by the hammerhead ribozyme, whereby changes in conformation and metal ion binding mode have been implicated in proceeding from reactant to transition state. Finally, a very large-scale structural rearrangement is studied in the L1 ligase riboswitch in order to explain the role of conserved residues in stabilizing conformational intermediates. These applications tie together several important factors that provide a broader understanding of the

interplay between chemical and conformational steps, and how they are affected by metal ions and other solvent components to achieve catalysis.

### 14.3.1.    Case Study: Comparison of DFT Functionals on Model Phosphoryl Transfer Reactions

In recent years, density-functional theory has emerged as the computational quantum chemistry method of choice for biological problems of medium size range (up to a few hundreds of atoms) in applications that do not require extensive conformational sampling. The field continues to advance in the accuracy of new functionals, the improvement of algorithms and the functionality and computational performance of software [81].

In the case of the development of new density-functional exchange-correlation functionals, the current climate is one of rapid change. New functionals are being turned out at an increasingly feverish pace. This wave of new functionals makes it difficult for the community to assess their limitations and general reliability over a sufficiently broad range of chemistry. The strategy we have thus far taken is not to continually jump from one DFT functional to another as soon as a new functional appears to have made incremental improvement. Instead, we have remained largely with well-established functionals that have limitations that are well characterized and understood. Moreover, the data which we accumulate in this way has greater impact by allowing cross-comparison with other calculations such as those collected in the *QCRNA* database.

Nonetheless, we find it is important to periodically assess the state of the art and reset our existing gold standard in order to progress. Up to this point, the majority of our database calculations have been performed using the well-established hybrid three parameter exchange functional of Becke [82, 83] coupled with Lee, Yang and Parr (LYP) correlation functional [84]. This functional performs relatively well for phosphate anions, but in general predicts barriers that are systematically too low, in contrast to Hartree-Fock methods that are usually systematically too high. Recently, new insights into the origin of the current limitations of density-functional theory have been discussed [81]. One of the most important caveats that developers have to deal with to improve B3LYP and other functionals of its generation is their poor description of medium and long range correlation that give rise to intermolecular dispersion interactions. Several different approaches to improve dispersion have been proposed such as perturbationally corrected functionals [85], the addition of semiempirical correction terms [86–88], and the use of an exchange-hole dipole moment model [89]. An alternate approach is to reparametrize existing functional forms so as to better model mid-range correlation. By construction these models can not reliably predict the long-range behavior of the dispersion energy, but if appropriately reparametrized, can considerably improve short-range non-bonded interactions. This latter approach has been adopted on a grand scale in the M05 and M06 suite of density functionals [90]. Many of these functionals have been praised for their

accuracy over a wide spectrum of applications [91, 92]. The current front-runner recommended for main group chemistry is the M06-2X functional [93]. In this section we compare B3LYP and M06-2X for the transesterification of a dinucleotide reaction model, and for the pathological case of dissociation of the *p*-nitrophenyl phosphate dianion.

### 14.3.1.1. *Transesterification of a Dinucleotide Model*

The transesterification reaction is at the core of the catalytic process in several prototype ribozymes such as the hammerhead and hairpin ribozymes discussed in the next section, and thus the accurate modeling of this reaction is critically important for the study of ribozyme catalysis (Figure 14-3). This reaction involves interaction between the $2'$ alcohol nucleophile of the RNA sugar with a highly negatively charged and polarizable adjacent 3' phosphate group. While in ribozymes it is generally believed that the reaction proceeds via a general or specific acid/base catalytic mechanism with indirect and, in some cases, direct chemical involvement of nucleobase functional groups and metal ions along the reaction coordinate [94, 95].

The transesterification of phosphates in a dianionic state is a concerted transformation sporting a single transition state along the reaction path [96]. Both B3LYP and M06-2X density functionals provide a similar description in terms of energetics and geometry on the 2-D potential energy surface. When modeling reactions occurring in biomacromolecules such RNA, the stationary points are an important but incomplete set of data to analyze. Accurate modeling of ribozymes requires capturing dynamic effects along the chemical step. These effects depend not only on the reactant, transition state and product species, but also on the specific shape and curvature of the multidimensional energy surface connecting reactant and product.

The first and simplest approach to explore the shape of the potential energy surface along the transesterification process is to follow the intrinsic reaction coordinate from the transition state downhill to reactant and product (Figure 14-4). From Figure 14-4 it is clear that the curvatures near the transition states are fairly different for each functional. B3LYP shows a steeper pathway uphill from the ligated dinucleotide model than M06-2X, however, on the other side of the saddle point, the descent to the cleaved product is less pronounced for B3LYP than for M06-2X. Despite the differences in slope shown by these functionals, the curvature of the potential energy surface at the transition state is similar for both profiles, and yield similar imaginary frequencies ($-153.22$ and $-168.42\,\text{cm}^{-1}$ for B3LYP and M06-2X, respectively). A more thorough approach to explore the potential energy surface can be taken into account by computing a two dimensional surface where the forming/breaking bond lengths are varied independently (Figure 14-5). Despite the overall qualitative similarity of the potential energy surfaces computed with B3LYP and M06-2X, there remain some quantitatively significant differences, emphasizing the need for careful selection of the density functional for chemical reactions where dynamical effects may be important to the reaction rate.

### 14.3.1.2.     Dissociation of p-*Nitrophenyl Phosphate*

It is often convenient experimentally to utilize small molecule substrates such as modified phosphates with enhanced leaving groups as reaction models. These compounds serve as models for RNA transesterification or phosphoryl transfer of phosphate monoesters in kinases and phosphatases. One particularly useful chemical probe for mechanistic studies is the molecule *p*-nitrophenyl phosphate (pNPP). Phosphoryl transfer in pNPP can be easily followed spectrophotometrically, and allows for kinetic isotope effects to be measured at primary (bridge) and secondary (non-bridge) phosphoryl oxygen positions, as well as at the exocyclic nitro group N position [97]. Sulphuryl substitution on the phosphate oxygen has been used to investigate the kinetic and stereochemical aspects of phosphoryl transfer [97, 98].

The experimental barrier for the hydrolysis of dianionic pNPP is estimated to be 29.5 kcal/mol at 39°C, while for the thio-substituted analog *p*-nitrophenyl thiophosphate (pNPTP) an approximate value of 27.9 kcal/mol has been observed [99, 100]. However, preliminary calculations of dianionic pNPP dissociation in solution using QM/MM methods based on the AM1/d-PhoT model predict (incorrectly) a much lower barrier. The origin of the problem can be traced back to the use of the B3LYP functional used to generate the reference data from which AM1/d-PhoT was developed. Despite the fact that B3LYP has been demonstrated to be nearly as accurate as much higher level methods for prediction of *relative* proton affinity values [45], the case of *p*-nitrophenol is somewhat of an anomaly having a proton affinity value in error of −4.2 kcal/mol, suggesting it is an even more enhanced leaving group to the extent that, in solution, dissociation is nearly barrierless.

Figure 14-2 compares the potential energy curve for dianionic pNPP dissociation using B3LYP (red) and M06-2X (blue) for dianionic pNPP (top) and pNPTP (bottom). The potential energy curve for dianionic pNPP dissociation using B3LYP indicates barrierless dissociation in the gas phase, whereas with the M06-2X functional, there is a barrier of 2.3 kcal/mol. The situation is similar comparing the dianionic pNPTP dissociation where again, B3LYP predicts a kinetically insignificant barrier and M06-2X predicts a 4.6 kcal/mol barrier. In solution it is expected that the barrier will be considerably increased due to solvent stabilization of the dianionic transition state. The striking feature is that B3LYP predicts barrierless dissociation in the gas phase, whereas M06-2X predicts a stable reactant species with activation energy barriers of 2.3 and 4.6 kcal/mol. Although further investigation needs to be made of this reaction profile in solution, the present results underscore the need to continue to assess new DFT functionals for their accuracy and predictive capability in order to determine the best, affordable quantum chemistry model from which high-level reference data for phosphoryl transfer reactions can be generated.

### 14.3.2.     Case Study: Chemical Steps of Catalysis in Hairpin Ribozyme

At first glance, it would seem that RNA enzymes, composed of fairly inert nucleobases connected by a sugar-phosphate backbone, are simply not equipped with

*Figure 14-2.* Potential energy curves (relative to separated monoanions) for the dissociation of *p*-nitrophenyl phosphate (pNPP) and *p*-nitrophenyl phosphorothioate (pNPTP) in the gas phase



*Figure 14-3.* Transesterification reaction of the dinucleotide model where the nucleophile-containing ribose sugar is modelled by a tetrahydrofurane structure, whereas the cleaving sugar is further simplified and modelled as a simple primary alcohol (ethanol)

*Figure 14-4.* Intrinsic reaction coordinate for the transesterification of the dinucleotide model with B3LYP and M06-2X functionals. Relative free energies of reaction and activation are provided in kcal/mol



*Figure 14-5.* Side by side comparison of the two-dimensional potential energy surface for the transesterification reaction computed with B3LYP (*left*) and M06-2X (*right*)

an adequate array of chemical functional groups for effective catalysis. This is in stark contrast to protein enzymes that have a fairly diverse repertoire of amino acids. The central question that has gripped the community that studies RNA enzymes is, simply, by what mechanisms can these molecules achieve catalysis [101]?

A number of factors have been implicated to be important for RNA catalysis, including the involvement of functional groups of the nucleobases or RNA backbone, divalent metal ions or other solvent components that might provide electrostatic

stabilization or act as general acid and base catalysts [102, 103]. Nonetheless, there currently exists no general consensus as to the origin of the catalytic proficiency exhibited by ribozymes, nor any detailed mechanism that has been unambiguously determined. In several small prototype systems, such as the hammerhead [104], hepatitis delta virus [105], and the L1 ligase [106] ribozymes, metal ions are essential for catalysis as well as RNA folding. The dual role played by metal ions in these systems complicates the identification of the chemical origins of catalysis and the unambiguous determination of detailed mechanism [102].

In contrast, the hairpin ribozyme (HPR) [107, 108], which catalyzes the reversible, site-specific phosphodiester bond cleavage of an RNA substrate, is unique in that the chemical steps of the reaction do not require involvement of a divalent metal ion [107–111]. This lack of an explicit metal ion requirement [112] makes the hairpin ribozyme an ideal target for theoretical studies aimed to characterize the contribution of "generalized solvation" provided by the solvated ribozyme on catalysis.

Here, we demonstrate with combined QM/MM simulations that the electrostatic environment provided by solvated HPR active site lowers the cleavage activation barrier up to 9 kcal/mol relative to the uncatalyzed transphosphorylation barrier in aqueous solution, accounting for the majority of the experimentally observed rate enhancement. Further work [54, 55] has gone on to explore in mode detailed mechanistic scenarios whereby A38 and G8 act as a general acid and base. The present results suggest that the electrostatic environment of the solvated ribozyme active site contributes significantly in achieving $10^6$ to $10^7$-fold rate enhancement of the phosphodiester cleavage [113–115] relative to the uncatalyzed, but spontaneous cleavage of RNA molecule in aqueous solution [116, 117]. Without the aid of a divalent metal ion, nor direct participation of nucleobase functional groups as a general acid or a general base, the majority of the observed rate enhancement can be realized through specific hydrogen bonding interactions (provided from G8 and other nucleobases) and non-specific electrostatic interactions of the solvated ribozyme active site. In the discussion that follows, the term "electrostatic solvation" is used to discuss the electrostatic component of the "generalized solvation" provided by the ribozyme environment.

The HPR in-line monoanionic mechanism considered in the present work are depicted in Scheme 14-1. This mechanism involves three reaction steps: (1) an initial intramolecular proton transfer from the 2′-hydroxyl group to either the pro-R ($O_{1P}$) or the pro-S ($O_{2P}$) non-bridging oxygen atoms of the scissile phosphate group, (2) a nucleophilic attack from the 2′-hydroxyl oxygen at the phosphate center, and (3) an exocyclic bond-cleavage of the leaving group from the phosphate center along with a second intramolecular proton transfer from the phosphate non-bridging oxygen to the leaving group. The proton transfer and nucleophilic substitution steps can occur either in a stepwise or concerted fashion. To explore these possibilities, two-dimensional reaction free energy profiles for the proton transfer and nucleophilic substitution reaction coordinates have been determined using molecular dynamics (MD) free energy simulations with a combined QM/MM potential along with density

*Scheme 14-1.* General in-line monoanionic mechanism of phosphodiester cleavage transesterification catalyzed by hairpin ribozyme; the first proton transfer ($PT1$), the nucleophilic attack ($Nu$), and the exocyclic cleavage ($Cl$) steps are shown, and the $O_{1P}$ and $O_{2P}$ pathways are indicated by *blue* and *red* colored hydrogens, respectively. For the uncatalyzed model reaction in solution, the $O_{1P}$ and $O_{2P}$ pathways are energetically equivalent

functional theory (DFT) corrections to the adiabatic pathways. In addition, the electrostatic solvation free energies are determined for the reactant state, intermediate states, transition states, and product states for both catalyzed and uncatalyzed reactions to address the effects of electrostatic environment provided by the ribozyme on the reaction.

### 14.3.2.1.    Two-Dimensional QM/MM Potential of Mean Force Profiles

Simulations are based on the second transition state analog crystal structure (PDB code 1M5O) [110] and performed using CHARMM [118] (version c32a2). Stochastic boundary MD was performed in a solvated 25-Å sphere centered at the scissile phosphate in the ribozyme active site using the all-atom CHARMM27 nucleic acid force field [119] and TIP3P water model [120], with the AM1/d-PhoT quantum model [33] and GHO method [121] for treatment of the QM/MM boundary. Full details are described elsewhere [55]. Two-dimensional potential umbrella sampling MD simulations [122] were performed, from which the potential of mean force profiles were constructed using the weighted histogram analysis method [123]. The reaction coordinates consist of a nucleophilic substitution coordinate, $\zeta_1 = R(P - O_{5'}) - R(P - O_{2'})$, and a proton transfer coordinate, $\zeta_2 = R(O_{X'} - H_{2'}) - R(O_{NB} - H_{2'})$, where $O_{NB}$ is either $O_{1P}$ or $O_{2P}$, and $O_{X'}$ is $O_{2'}$ for the first (nucleophilic bond formation) step, and $O_{5'}$ for the second (leaving group bond cleavage) step. Each umbrella window was run for 17 ps of equilibration and 50 ps of configurational sampling. The uncatalyzed model reaction consisted of a molecule

of 2-hydroxyethyl methyl phosphate solvated with a 40-$Å^3$ cubic box of 2038 water molecules and one $Na^+$ ion. Simulations were carried out using QM/MM-Ewald scheme [67] at 1 atm and 300 K. The computed free energy values were further refined by density functional adiabatic energies computed at the B3LYP/6-311++G(3df,2p)//B3LYP/6-31++G(d,p) level, in which "//" separates the level for the refined single point energy from the level for geometry optimization. The same geometry optimizations were carried out at the AM1/d-PhoT level in order to derive approximate correction factors. Figure 14-6 shows the two-dimensional free energy profiles for the uncatalyzed model transphosphorylation reaction in aqueous solution, and catalyzed by the HPR along $O_{1P}$ and $O_{2P}$ pathways (Scheme 14-1). Table 14-1 lists free energy values corresponding to stationary points along the minimum free energy path on the surface, determined from the QM/MM free energy simulations, along with corrections at the DFT level to the activation and reaction free energy values.



*Figure 14-6.* Two-dimensional free energy surfaces for in-line monoanionic mechanisms for the (**A**) un-catalytic model reaction in solution, and the catalytic (**B**) $O_{1P}$ and (**C**) $O_{2P}$ pathways in the hairpin ribozyme. $\zeta_1$ is defined as $R_{P-O_{5'}} - R_{O_{2'}-P}$, and $\zeta_2$ is $R_{O_{2'}-H_{2'}} - R_{O_{NB}-H_{2'}}$ for $\zeta_1 < 0.0$ Å and $R_{O_{5'}-H_{2'}} - R_{O_{NB}-H_{2'}}$ for $\zeta_1 > 0.0$ Å, where $O_{NB}$ is for the $O_{1P}$ proton transfer in (**B**), and for the $O_{2P}$ proton transfer in (**A**) and (**C**), respectively. The units for free energies and distances are kcal/mol and Å, respectively

*Table 14-1.* Calculated reaction free energies and barrier heights (kcal/mol) for uncatalyzed model and catalyzed transesterification reactions in solution and in the hairpin ribozyme[a]

| | GB[b] | TS$_{PT1}$ | INT$_1$ | TS$_{Nu}$[c] | INT$_2$ | TS$_{Cl}$ | Prod |
|---|---|---|---|---|---|---|---|
| Soln[d] | O$_{2P}$ | 19 | 19 | 32 | 31 | 37 | 0 |
| | | | | (33) | (33) | (38) | (2) |
| HPR[e] | O$_{1P}$ | 15 | 14 | 15 | 13 | 25 | −7 |
| | | | | (16) | (16) | (27) | (−6) |
| | O$_{2P}$ | 12 | 11 | 14 | 13 | 21 | −5 |
| | | | | (18) | (17) | (21) | (−4) |
| Expt | Soln[f] | | 21 | 32 | 25 | 34 | |
| | HPR[g] | | | | | ∼20–21 | ∼1 |

[a]The values are those estimated from the 2-D PMF profiles described in the text and given in Figure 14-6. A DFT correction (in parenthesis) of the semiempirical AM1/d-PhoT model is applied to the intermediate, product, and transition states, respectively, based on active site model calculations.
[b]General base (GB) activating the O$_{2'}$ nucleophile.
[c]A DFT correction on the TS$_{Nu}$ is based on the error of AM1/d-PhoT model at the TS$_{PT1+Nu}$, in which the nucleophilic attack (TS$_{Nu}$) is concerted with the proton transfer (TS$_{PT1}$) in the gas phase.
[d]Uncatalyzed model reaction in solution.
[e]Catalyzed reaction in the hairpin ribozyme.
[f]Values are estimated from references [152] and [153], which combine experimental and computational values for the reaction free energies and activation energies of relevant reactions in solution.
[g]Experimental values for the hairpin ribozyme are taken from references [113–115].

### 14.3.2.2.    Active Site Structure and Mechanism

Figure 14-7 shows representative snapshots of the transition states for the nucleophilic substitution and exocyclic cleavage steps. The overall reaction may be characterized by a sequence of proton transfer, nucleophilic attack, exocyclic cleavage, and proton transfer steps. The rate-limiting step is the exocyclic bond-cleavage of the leaving group from the phosphorus atom, followed by a barrierless proton transfer to the departing O$_{5'}$ alkoxide anion. The same trend is found both for the catalyzed reaction in the ribozyme and the uncatalyzed model reaction in aqueous solution, but the hairpin ribozyme markedly lowers the reaction barriers for each of the three reaction steps and the free energies of the resulting intermediates. The greatest barrier reduction occurs in the nucleophilic attacking step (a net decrease of 17 kcal/mol for the O$_{1P}$ pathway and 18 kcal/mol for the O$_{2P}$ pathway relative to the uncatalyzed reaction). For the rate-limiting step, the free energy barriers are lowered by 12 and 16 kcal/mol along the O$_{1P}$ and the O$_{2P}$ pathways, respectively, while experimental estimation of barrier reduction is 13–14 kcal/mol. After the density functional correction at the B3LYP/6-311++G(3df,2p) level, the overall free energy barrier becomes 27 kcal/mol for the O$_{1P}$ pathway and 21 kcal/mol for the O$_{2P}$ pathway. These results are in accord with the experimental estimate of 20–21 kcal/mol [113–115]. Nonetheless, this does not preclude alternate mechanisms with explicit nucleobase involvement, such as A38 and G8 acting as a general acid and a general base catalyst [110, 124, 125], that could further lower the barrier.

*Figure 14-7.* Snapshots of the active site structures near the transition state of (*top*) the nucleophilic attack and (*bottom*) the exocyclic cleavage for the in-line monoanionic $O_{2P}$ mechanism of cleavage transesterification in the hairpin ribozyme. The yellow and red colored cartoon is for the substrate and ribozyme strands, respectively, and water molecules interacting with non-bridging oxygens and $O_{5'}$ are shown

These results explore the effects by which the change of the electrostatic environment provided by the hairpin ribozyme relative to that of aqueous solution affects the rate of the transphosphorylation reaction. Since the specific mechanisms explored here do not involve direct intervention of any nucleobases as a general base or general acid in the catalysis, the computed change in the free energy barriers is mostly due to the change of the heterogeneous electrostatic environment in the HPR active site relative to that of bulk solvation by water. The direct electrostatic solvation by the ribozyme and water lowers the overall free energy barrier by 7 and 9 kcal/mol for the two reaction paths corresponding to an initial proton transfer to either of the two non-bridging phosphate oxygen atoms. The results suggest that the non-specific interactions in HPR are sufficient to account for the majority of the observed change of barrier heights without the involvement of a metal ion and general acid-base catalysis by active site nucleobases. The in-line monoanionic mechanism establishes a baseline mechanism that invokes only the generalized solvation and specific hydrogen bonding interactions provided by the ribozyme environment, and provides a departure point for the exploration of alternate mechanisms where participation of nucleobases in the active site play an active chemical role.

### 14.3.3.    Case Study: Role of Divalent Metal Ions in Hammerhead Ribozyme Catalysis

The hammerhead ribozyme (HHR) catalyzes the same type of transesterification reaction as the hairpin ribozyme, which involve the site-specific attack of an activated $2'$OH nucleophile to the adjacent $3'$ phosphate, resulting in cleavage of the P-O5$'$ phosphodiester linkage to form a $2',3'$ cyclic phosphate and a $5'$ alcohol. However, unlike the hairpin ribozyme that has no divalent metal ion requirement for catalysis, the hammerhead ribozyme, under physiological conditions, requires divalent metal ions to promote its catalytic step [7, 104]. Recent crystallographic studies of a full length HHR have characterized the ground state active site architecture [126] and its solvent structure [127], including the binding mode of a presumed catalytically active divalent metal ion in the active site. These findings have reconciled a long-standing controversy between structural and biochemical studies for this system [128]. It is still not clear, however, what are the roles of the divalent ion in the active site. Recently, large scale molecular dynamics simulations using both molecular mechanics and hybrid QM/MM potentials have been used to explore the structure and dynamics at different states along the catalytic pathway in order to shed light on the possible role of divalent ions in the catalytic mechanism [52, 53].

### 14.3.3.1.    *HHR Folds to Form an Electrostatic Negative Metal Ions Recruiting Pocket*

One of the recent crystallographic structures of the full length HHR identifies five well-defined divalent ion binding sites, one of them (the C-site) being located in the catalytic pocket and being suggested to have a direct role in catalysis [127].

In order to probe cation occupation in the active site in the absence of $Mg^{2+}$ ions, we examined $Na^+$ distributions in the reactant and activated precursor (deprotonated 2OH' nucleophile) states. It has been noted in the recent literature that the modeling of ions in highly charged systems such as HHR affords tremendous challenges with regard to simulation time scales [129]. This section presents the results of series of five 300 ns simulations of the full length HHR, in both the reactant and activated precursor states, in order to ascertain the cation occupation requirement of the active site to maintain catalytic integrity.

The 3D density contour maps for the $Na^+$ ion distribution determined over the last 250 ns of simulation (Figure 14-8) show that the overall highest probability $Na^+$ occupation sites were concentrated in the active site for both the reactant and activated precursor. This suggests that the HHR folds to form a strong local electronegative pocket that is able to attract and bind $Mg^{2+}$ if present in solution, or recruit a high local concentrations of $Na^+$ ions in the absence of $Mg^{2+}$.



*Figure 14-8.* The 3D density contour maps (*yellow*) of $Na^+$ ion distributions derived from the activated precursor simulation. The hammerhead ribozyme is shown in blue with the active site in red. Only the high-density contour is shown here to indicate the electrostatic recruiting pocket formed in the active site

### 14.3.3.2.    The Bridging $Mg^{2+}$ Induces a Significant $pK_a$ Shift of the General Acid

The first published crystal structure of the full length HHR [126] in which there was no solvent or ions resolved showed A9 and the scissile phosphate in close proximity, consistent with the interpretation of thio effect measurements [130], and the $G8:O_{2'}$ and $G12:N_1$ poised to act as a general acid and base, respectively, as proved in previous photocrosslinking [131] and mutation experiments [132]. Given the strong evidence that $Mg^{2+}$ participates directly in the catalytic process together with the spatial proximity of the A9 and scissile phosphate, made the placement of an $Mg^{2+}$ ion in bridging position a reasonable assumption.

We have explored the role the $Mg^{2+}$ ion placed at the bridging position in the reactant state, the early transition state (ETS), and the late transition state (LTS) [52, 53]. In these studies we have used specifically designed molecular mechanics residues that are able to reproduce the geometry and charge distribution in the early and late transition states of the phosphoryl transfer reaction [46]. The $Mg^{2+}$ ion remained in the bridging position for the entire duration of the three simulations, displaying different potential roles at specific points of the reaction pathway. In the ETS simulation (Figure 14-9, left panel), the $Mg^{2+}$ ion is directly coordinating $G8:O_{2'}$ to induce a possible shift in its $pK_a$ to act as a general acid, while in the LTS



*Figure 14-9.* Snapshots from the simulations of the early transition state mimic (*left*) and the late transition state mimic (*right*), indicating the $Mg^{2+}$ ion direct coordination (*green lines*) and key hydrogen bonds and indirect $Mg^{2+}$ coordination (*dotted lines*). For clarity, the water molecules are not shown

simulation (Figure 14-9, right panel), the $Mg^{2+}$ ion acts as a potential Lewis acid catalyst to stabilize the leaving group, being poised to assist proton transfer from the $G8:O_{2'}$.

Hybrid QM/MM simulations were performed to further probe the role of the bridging $Mg^{2+}$ ion [52]. The spontaneous proton transfer from the implicated general acid, $G8:O_{2'}$, to the leaving group, $C1.1:O_{5'}$ was observed within the first ns of QM/MM simulation, confirming our assumptions about the role of the bridging $Mg^{2+}$ ion in the catalytic step of HHR. Thus, our simulation results supported the supposition that a single bridging $Mg^{2+}$ ion could assist in the cleavage step in HHR catalysis by acting to increase the acidity of the 2OH' of G8. The $Mg^{2+}$ at the bridging position also preserves the integrity of the active site structure, and may serve as an epicenter in the transition state that coordinates the A9 and scissile phosphates, $G8:O_{2'}$ general acid and $C1.1:O_{5'}$ leaving group.

### 14.3.3.3. The Accumulation of the Negative Charge in the Precursor State Causes the Migration of the $Mg^{2+}$ Ion from the C-Site to the Bridging Position

Recently, a joint experimental/theoretical study has been reported of the full length hammerhead structure with resolved solvent and metal ions [127]. In this structure a resolved $Mn^{2+}$ ion in the active site was not positioned in a bridging position as postulated in our previous simulations; instead, it binded $G10.1:N_7$ and $A9:O_{2P}$ (the C-site). We performed simulations with a $Mg^{2+}$ initially placed at the C-site for different stages along the reaction pathway [52]. The results, shown in Figure 14-10, suggest that the $Mg^{2+}$ in fact migrates from the C-site to the bridging position in the transition states and the deprotonated reactant state. This migration is caused by the accumulated negative charge at the cleavage site after the general base step.

In summary, our simulation results draw a possible picture of the roles of $Mg^{2+}$ in supporting the catalytic step of HHR. First, HHR folds to form an electronegative cation recruiting pocket that attracts a $Mg^{2+}$ ion to the C-site. The ion moves to the bridging position between A9 and the scissile phosphate either upon deprotonation of the 2OH' nucleophile, or formation of the dianionic transition state. In this position, the $Mg^{2+}$ ion is poised to provide direct electrostatic stabilization of the transition state and the accumulating negative charge on the leaving group. Moreover, the $Mg^{2+}$ shifts the $pK_a$ of the general acid ($G8:O_{2'}$), and after the proton transfer to the leaving group, reverts back to stabilize the conjugate base. This mechanistic interpretation is supported by the present simulations, and is consistent with a considerable body of experimental work. First, the thio/rescue effect experiments [130] support a mechanism in which a single metal cation bound at the C-site in the ground state acquires an additional interaction with the scissile phosphate in proceeding to the transition state. Second, kinetic studies [133], photocrosslinking experiments [131] and mutational data [132, 134, 135] implicate G8 and G12 as possible general acid and base. Lastly, recent studies involving metal ion titrations suggest that the $pK_a$ of the general acid is down- shifted by around 4–7 $pK_a$ units in a metal-dependent

*Figure 14-10.* A schematic view of the possible migration of the $Mg^{2+}$ ion from the C-site to the bridging position. Spontaneous migration was predicted from the simulation for the transition states and the deprotonated reactant state, with the $Mg^{2+}$ ion initially placed at the C-site

manner, correlated with the metal $pK_a$ [94], and indicate that divalent metals may play a specific chemical role in catalysis [136].

These results represent an important first step in the detailed characterization of the structure, dynamics and free energy profile for the full length HHR catalytic mechanism. Together with experiment, it is the hope that a consensus will emerge that explains the detailed molecular mechanisms of hammerhead ribozyme catalysis, and in doing so may provide new insight into the guiding principles that govern RNA catalysis.

### 14.3.4.    Case Study: Conformational Transition in the L1 Ligase Ribozyme

RNA is characterized by a large and diverse ensemble of conformations that interchange on time scales that range from femtoseconds to microseconds [137]. This conformational variability allows RNA molecules to be designed to allow binding and catalytic activity (to be allosterically controlled) such as in the case of aptamers and aptazymes [138], or have the ability to regulate gene expression by binding small

molecules such as in the case of riboswitches [139]. However, due to RNA's rugged conformational landscape, structural biology methods such as X-ray crystallography and NMR face challenges to capture an accurate, complete picture of the range of important conformations, and their time scales, that might play important roles in function [140]. On the other hand, molecular simulation methods provide a wealth of detail into both structure and dynamics, and offer a powerful tool to complement structural biology, biochemical and biophysical experiments.

In this section, we report molecular simulations on the large-scale conformational transition of the L1 ligase ribozyme from an inactive to a catalytically active state. L1 Ligase ribozyme functions as an *allosteric molecular switch* or *aptazyme* and was iteratively optimized by *in vitro* selection to catalyze regioselectively and regiospecifically the 5′ to 3′ phosphodiester bond ligation (nucleotidyl transfer reaction) with the possibility to be controlled (activated) by small molecules, peptides and even proteins [141–144]. An unique feature of L1 ligase is its intrinsically flexible non-canonically base paired ligation site, a characteristic possessed only by two other ligase ribozymes [145, 146]. There is no naturally occurring ligase ribozyme and moreover, among the synthesized ligases ribozymes there are only five that accomplish their function in an regiospecific and regioselective way [141, 146–149].

L1 ligase's large intrinsic flexibility was revealed by the recent crystal structure of the ligation product of a reduced size variant with two vastly different conformers, differing by reorientation of one of the stems by around 80 Å, that were resolved in the same asymmetric cell [106]. Based on the presence/absence of specific contacts between distant conserved parts including the ligation site and a totally conserved residue, U38, one of the conformers was postulated to represent the catalytically active or *on* conformation, the other the inactive or *off* one [106].

We have explored the coupled on–off conformational switch in L1 ligase using large scale molecular dynamics simulations for more than 600ns departing from both active and inactive conformations (Figure 14-11). Based on the crystal structure of the two conformers, we identified a limited set of four virtual torsions (out of a total of 142) that can be used to distinguish between the active and inactive conformations found in crystal, They were denoted as $\theta_{18}$, $\theta_{37}$, $\theta_{44}$, and $\eta_{38}$, and were defined following Ref. [150]. These virtual torsions span two conserved and restricted regions located in the three-way junction and a loop that contains U38, the conserved residue that is postulated to be responsible for allosteric control of the catalytic step [106, 141–144].

The conformational rearrangement in the three-way junction and the U38 loop transition can both be mapped by monitoring the four virtual torsions, and occur on different time scales. The U38 loop transition occurs on the order of tens of nanoseconds, whereas the rearrangement of the three way junction is estimated to occur on a time scale longer than 0.4 $\mu$s. On this time scale of our simulations (several hundred of ns) the L1 ligase in its inactive conformation was predicted to cover approximately a third of the complete 80 Å conformational switch (Figure 14-11). Since this transition might correspond to the rate-controlling step of L1 ligase catalysis, it

*Figure 14-11.* Snapshots from the conformational switch path explored in the vicinity of the active conformation (unfolded docked conformation) and starting from the inactive (undocked) conformation found in crystal. Stems A and B were aligned (best fitted) and are shown in *yellow*, different instances of stem C are shown in stick representation with different colors. A schematic of the non-canonical binding scheme of the ligation site is shown in the *right panel* and the general mechanism of ligation in the *left panel*

is of interest to identify the dynamic hinge points and stabilization of conformational intermediates that allow the transition to occur.

The overall fluctuations in the active conformation of the reactant and product states were restricted to a reduced portion of the available conformations due to distant tertiary contacts between U38 loop and the ligation site (RMSD $\sim 4.2$ Å). In the absence of these tertiary contacts, the fluctuations are significantly larger (RMSD $\sim 6.7$ Å) in the vicinity of the active conformation or $\sim 7.9$ Å in the inactive form). The origin of these large variations were traced to fluctuation of the restricted region of the junction, with all the other structural elements remaining close to their starting structures. The large fluctuations observed in the simulations were accompanied by

the formation of new contacts not observed in any of the crystal structure between two conserved portions of the L1 Ligase: stem B and U19. Given their conserved nature we postulate that these contacts have to play an important role in stabilization of intermediary states along the *conformational switch – catalytic step* pathway.

The non-canonically base-paired ligation site shows a high degree of variability, and visits three distinct conformational states characterized by specific hydrogen bonding patterns between GTP1 on one side and G2:U50, U38:A51 and G52:U71 base pairs. The ligation reaction takes place between GTP1 and U71, and simulations were performed from two different initial arrangements of the reactant state differing in the conformation of the GTP1 triphosphate conformation and its ion coordination. Simulations of the ligation site predicted formation of conformational states where the $U71:O_{3'}$ atom, the nucleophile, makes close contacts to a potential general base, $GTP1:O_{2\alpha}$. The formation of these contacts was highly correlated with ligation site being in either the first or the third hydrogen binding pattern and absent when visiting the second one, suggesting that the L1 ligase catalysis might be facilitated by these specific hydrogen bonding patterns which are a direct result of the non-canonically base-paired ligation site and the intrinsic flexibility of the molecule.

The present simulation results have identified important hinge points in the conformational transition from inactive to active forms of the L1 ligase, and characterized interactions that stabilize intermediates along the transition pathway. The insights gained from these simulations are a first step toward a detailed understanding of the coupled catalytic/conformational riboswitch mechanism of L1 ligase that may ultimately enhance the future design and engineering of new catalytic riboswitches.

## 14.4. CONCLUSION

In this chapter, we present a description of a multiscale simulation strategy we have developed to attack problems of RNA catalysis. Ribozyme systems, due to the high degree of charge, strong interaction with ions and other solvent components, and large conformational variations that are coupled with the chemical steps of catalysis, present special challenges not present in typical protein systems, and consequently demand new methods. The main methodological components are herein summarized, including the assembly of the *QCRNA* database, parametrization of the AM1/d-PhoT Hamiltonian, and development of new semiempirical functional forms for improved charge-dependent response properties, methods for coupling many-body exchange, correlation and dispersion into the QM/MM interaction, and generalized methods for linear-scaling electrostatics, solvation and solvent boundary potentials. We then present results for a series of case studies ranging from non-catalytic reaction models that compare the effect of new DFT functionals, and on catalytic RNA systems including the hairpin, hammerhead and L1 ligase ribozymes. The ultimate goal of this work is to develop new multiscale computational tools and bring them to bear on the study of the mechanisms of RNA catalysis. The results may serve to aid in the interpretation of experiments, provide a deeper understanding of

ribozyme mechanisms, and unravel guiding principles for RNA catalysis that may facilitate the design of new technology.

## ACKNOWLEDGMENTS

## REFERENCES

1. Gilbert W (1918) Nature 319:618.
2. Scott WG (1996) Biochem Soc Trans 24:604.
3. Gesteland RF, Cech TR, Atkins JF (1999) The RNA World: The nature of modern RNA suggests a Prebiotic RNA, 2nd ed., Cold Spring Harbor Laboratory Press, New York.
4. Yarus M (1999) Curr Opin Chem Biol 3:260.
5. Chen X, Li N, Ellington AD (2007) Chem Biodivers 4:633.
6. Lilley DM (2008) Ribozymes and RNA Catalysis, chap. The Hairpin and Varkud Satellite Ribozymes, 66–91, RSC Biomolecular Series, RSC Publishing, Cambridge.
7. Scott WG (2007) Curr Opin Struct Biol 13:280.
8. Rubenstein M, Tsui R, Guinan P (2004) Drugs of the Future 29:893.
9. Vaish NK, Dong F, Andrews L, Schweppe RE, Ahn NG, Blatt L, Seiwert SD (2002) Nature Biotech 20:810.
10. Breaker RR (2002) Curr Opin Biotechnol 13:31.
11. Biopolymers (2007) Biopolymers 85:169.
12. Gao J, Furlani TR (1995) IEEE Comput Sci Eng 2:24.
13. Gao J (1995) Rev Comput Chem 7:119.
14. Hawkins GD, Zhu T, Li J, Chambers CC, Giesen DJ, Liotard DA, Cramer CJ, Truhlar DG (1998) Combined quantum mechanical and molecular mechanical methods, 201–219, ACS Symposium Series 712, Oxford University Press, New York.
15. Monard G, Merz KM, Jr (1999) Acc Chem Res 32:904.
16. Warshel A (2002) Acc Chem Res 35:385.
17. Warshel A (2003) Annu Rev Biophys Biomol Struct 32:425.

18. Senn HM, Thiel W (2007) Curr Opin Chem Biol 11:182.
19. Norberg J, Nilsson L (2002) Acc Chem Res 35:465.
20. Orozco M, Pérez A, Noy A, Luque FJ (2003) Chem Soc Rev 32:350.
21. Orozco M, Noy A, Prez A (2008) Curr Opin Struct Biol 18:185.
22. TE Cheatham III (2004) Curr Opin Struct Biol 14:360.
23. Chen S-J (2008) Annu Rev Biophys 37:197.
24. Auffinger P, Hashem Y (2007) Curr Opin Struct Biol 17:325.
25. Dewar MJ, Thiel W (1977) J Am Chem Soc 99:4899.
26. Dewar MJ, Thiel W (1977) J Am Chem Soc 99:4907.
27. Dewar MJS, Zoebisch E, Healy EF, Stewart JJP (1985) J Am Chem Soc 107:3902.
28. Stewart JJP (1989) J Comput Chem 10:209.
29. Thiel W, Voityuk AA (1992) Theor Chim Acta 81:391.
30. Thiel W, Voityuk AA (1996) J Phys Chem 100:616.
31. Giese TJ, Sherer EC, Cramer CJ, York DM (2005) J Chem Theory Comput 1:1275.
32. Giese TJ, Gregersen BA, Liu Y, Nam K, Mayaan E, Moser A, Range K, O Nieto Faza, C Silva Lopez, A Rodriguez de Lera, Schaftenaar G, Lopez X, Lee T, Karypis G, York DM (2006) J Mol Graph Model 25:423.
33. Nam K, Cui Q, Gao J, York DM (2007) J Chem Theory Comput 3:486.
34. QCRNA, http://theory.chem.umn.edu/Database/QCRNA.
35. Range K, McGrath MJ, Lopez X, York DM (2004) J Am Chem Soc 126:1654.
36. Mayaan E, Range K, York DM (2004) J Biol Inorg Chem 9:807.
37. CS López, Faza ON, Gregersen BA, Lopez X, AR de Lera, York DM (2004) Chem Phys Chem 5:1045.
38. Range K, Riccardi D, Cui Q, Elstner M, York DM (2005) Phys Chem Chem Phys 7:3070.
39. Liu Y, Gregersen BA, Lopez X, York DM (2005) J Phys Chem B 109:19987.
40. Xu D, Guo H, Liu Y, York DM (2005) J Phys Chem B 109:13827.
41. CS López, Faza ON, AR de Lera, York DM (2005) Chem Eur J 11:2081.
42. Liu Y, Lopez X, York DM (2005) Chem Commun 31:3909.
43. Liu Y, Gregersen BA, Hengge A, York DM (2006) Biochemistry 45:10043.
44. Lopez X, Dejaegere A, Leclerc F, York DM, Karplus M (2006) J Phys Chem B 110:11525.
45. Range K, CS López, Moser A, York DM (2006) J Phys Chem A 110:791.
46. Mayaan E, Moser A, Mackerell AD Jr, York DM (2007) J Comput Chem 28:495.
47. Jmol, http://www.jmol.org.
48. Willighagen E, Howard M (2005) CDK News 2:17.
49. Schaftenaar G, Noordik JH, Molden, http://www.cmbi.ru.nl/molden/molden.html.
50. Schaftenaar G, Noordik JH (2000) J Comput.-Aided Mol Des 14:123.
51. Nam K, Gao J, York DM (2008) Multiscale simulation methods for nanomaterials, Ross RB, Sanat M (eds) Wiley, New York, pp 201–218.
52. Lee T-S, Silva-Lopez C, Martick M, Scott WG, York DM (2007) J Chem Theory Comput 3:325.
53. Lee T-S, Silva Lopez C, Giambasu GM, Martick M, Scott WG, York DM (2008) J Am Chem Soc 130:3053.
54. Nam K, Gao J, York DM (2008) J Am Chem Soc 130:4680.
55. Nam K, Gao J, York D (2008) RNA 14:1501.
56. Giese TJ, York DM (2005) J Chem Phys 123:164108.
57. Matsuzawa N, Dixon DA (1992) J Phys Chem 96:6232.
58. Giese TJ, York DM (2007) J Chem Phys 127:194101.
59. Wheatley RJ, Price SL (1990) Mol Phys 69:507.

60. Piquemal J, Cisneros G, Reinhardt P, Gresh N, Darden TA (2006) J Chem Phys 124:104101.
61. Tang KT, Toennies JP (2003) J Chem Phys 118:4976.
62. Tang KT, Toennies JP (1984) J Chem Phys 80:3726.
63. Pellenq R, Nicholson D (1998) Mol Phys 95:549.
64. Pellenq R, Nicholson D (1999) Mol Phys 96:1001.
65. Yang W (1991) Phys Rev A 44:7823.
66. Khandogin J, Hu A, York DM (2000) J Comput Chem 21:1562.
67. Nam K, Gao J, York DM (2005) J Chem Theory Comput 1:2.
68. Giese TJ, York DM (2008) J Comput Chem 29:1895.
69. Watson MA, P Sałek, Macak P, Helgaker T (2004) J Chem Phys 121:2915.
70. Pérez-Jordá JM, Yang W (1995) Chem Phys Lett 247:484.
71. York DM, Karplus M (1999) J Phys Chem A 103:11060.
72. Khandogin J, Gregersen BA, Thiel W, York DM (2005) J Phys Chem B 109:9799.
73. Gregersen BA, Khandogin J, Thiel W, York DM (2005) J Phys Chem B 109:9810.
74. Range K, Mayaan E, LJ Maher III, York DM (2005) Nucleic Acids Res 33:1257.
75. Gregersen BA, York DM (2005) J Phys Chem B 109:536.
76. Roux B, Beglov D, Im W (1999) Simulation and theory of electrostatic interations in solution, Pratt LR, Hummer G (eds) vol. 492 of Proceedings of the Santa Fe Workshop on Treatment of Electrostatic Interactions in Computer Simulations of Condensed Media, AIP Conference Proceedings, Melville, New York, pp 492–509.
77. Im W, Bernèche S, Roux B (2001) J Chem Phys 114:2924.
78. Banavali NK, Im W, Roux B (2002) J Chem Phys 117:7381.
79. Schaefer P, Riccardi D, Cui Q (2005) J Chem Phys 123:014905.
80. Gregersen BA, York DM (2006) J Comput Chem 27:103.
81. Cohen AJ, P Mori-Sánchez, Yang W (2008) Science 321:792.
82. Becke AD (1988) Phys Rev A 38:3098.
83. Becke AD (1993) J Chem Phys 98:5648.
84. Lee C, Yang W, Parr RG (1988) Phys Rev B 37:785.
85. Grimme S (2006) J Chem Phys 124:034108.
86. Grimme S (2004) J Comput Chem 25:1463.
87. Grimme S (2006) J Comput Chem 27:1787.
88. Schwabe T, Grimme S (2007) Phys Chem Chem Phys 9:3397.
89. Becke AD, Johnson ER (2006) J Chem Phys 124:014104.
90. Zhao Y, Truhlar DG (2008) Theor Chem Acc 120:215.
91. Gu J, Wang J, Leszczynski J, Xie Y, Schaefer HF III (2008) Chem Phys Lett 459:164.
92. Cramer CJ, Gour JR, Kinal A, Wloch M, Piecuch P, Shahi ARM, Gagliardi L (2008) J Phys Chem A 112:3754.
93. Zhao Y, Truhlar DG (2008) Acc Chem Res 41:157.
94. M Roychowdhury-Saha, Burke DH (2006) RNA 12:1846.
95. Thomas JM, Perrin DM (2008) J Am Chem Soc 130:15467.
96. Perreault DM, Anslyn EV (1997) Angew Chem Int Ed 36:432.
97. Hengge AC (2002) Acc Chem Res 35:105.
98. Hengge AC, Cleland WW (1990) J Am Chem Soc 112:7421.
99. Kirby AJ, Jencks WP (1965) J Am Chem Soc 87:3209.
100. Catrina IE, Hengge AC (1999) J Am Chem Soc 121:2156.
101. Takagi Y, Ikeda Y, Taira K (2004) Top Curr Chem 232:213.
102. Lönnberg T, Lönnberg H (2005) Curr Opin Chem Biol 9:665.

103. Sigel RK, Pyle AM (2007) Chem Rev 2007:97.
104. Scott WG (1999) Q Rev Biophys 32:241.
105. Shih I-H, Been MD (2002) Annu Rev Biochem 71:887.
106. Robertson MP, Scott WG (2007) Science 315:1549.
107. Walter NG, Burke JM (1998) Curr Opin Chem Biol 2:24.
108. Lilley DM (1999) Curr Opin Struct Biol 9:330.
109. Doherty EA, Doudna JA (2001) Annu Rev Biophys Biomol Struct 30:457.
110. Rupert PB, Massey AP, Sigurdsson ST, AR Ferré-D'Amaré (2002) Science 298:1421.
111. Bevilacqua PC (2003) Biochemistry 42:2259.
112. Fedor MJ, Williamson JR (2006) Nat Rev Mol Cell Biol 6:399.
113. Nesbitt SM, Erlacher HA, Fedor MJ (1999) J Mol Biol 289:1009.
114. Fedor MJ (2000) J Mol Biol 297:269.
115. Kuzmin YI, Da Costa CP, Cottrell JW, Fedor MJ (2005) J Mol Biol 349:989.
116. Hertel KJ, Peracchi A, Uhlenbeck OC, Herschlag D (1997) Proc Natl Acad Sci USA 94:8497.
117. Li Y, Breaker RR (1999) J Am Chem Soc 121:5364.
118. Brooks BR, Bruccoleri RE, Olafson BD, States DJ, Swaminathan S, Karplus M (1983) J Comput Chem 4:187.
119. Foloppe N, MacKerell AD, Jr. (2000) J Comput Chem 21:86.
120. Jorgensen WL, Chandrasekhar J, Madura JD, Impey RW, Klein ML (1983) J Chem Phys 79:926.
121. Gao J, Amara P, Alhambra C, Field MJ (1998) J Phys Chem A 102:4714.
122. Torrie GM, Valleau JP (1977) J Comput Phys 23:187.
123. Kumar S, Bouzida D, Swendsen R, Kollman P, Rosenberg J (1992) J Comput Chem 13:1011.
124. AR Ferré-D'Amaré (2004) Biopolymers 73:71.
125. Cottrell JW, Kuzmin YI, Fedor MJ (2007) J Biol Chem 282:13498.
126. Martick M, Scott WG (2006) Cell 126:309.
127. Martick M, Lee T-S, York DM, Scott WG (2008) Chem Biol 15:332.
128. Scott WG (2007) Biol Chem 388:727.
129. Ponomarev SY, Thayer KM, Beveridge DL (2004) Proc Natl Acad Sci USA 101:14771.
130. Wang S, Karbstein K, Peracchi A, Beigelman L, Herschlag D (1999) Biochemistry 38:14363.
131. Lambert D, Heckman JE, Burke JM (2006) Biochemistry 45:7140.
132. Blount KF, Uhlenbeck OC (2005) Annu Rev Biophys Biomol Struct 34:415.
133. Han J, Burke JM (2005) Biochemistry 44:7864.
134. McKay DB (1996) RNA 2:395.
135. Wedekind JE, McKay DB (1998) Annu Rev Biophys Biomol Struct 27:475.
136. M Roychowdhury-Saha, Burke DH (2007) RNA 13:841.
137. Crothers D (2001) RNA, chap. RNA Conformational Dynamics, Elsevier Science & Technology, Amsterdam, pp 61–71.
138. Bunka DHJ, Stockley PG (2006) Nat Rev Microbiol 4:588.
139. Schwalbe H, Buck J, Frtig B, Noeske J, Whnert J (2007) Angew Chem Int Ed 46:1212. URL http://dx.doi.org/10.1002/anie.200604163.
140. Xia T (2008) Curr Opin Chem Biol 2:1.
141. Robertson MP, Ellington AD (1999) Nature Biotech 17:62.
142. Robertson MP, Knudsen SM, Ellington AD (2004) RNA 10:114.
143. Robertson MP, Ellington AD (2000) Nucleic Acids Res 28:1751.
144. Robertson MP, Ellington AD (2001) Nature Biotech 19:650.
145. Landweber LF, Pokrovskaya ID (1999) Proc Natl Acad Sci USA 96:173.
146. Ekland EH, Szostak JW, Bartel DP (1995) Science 269:364.
147. Rogers J, Joyce GF (1999) Nature 402:323.

148. Jaeger L, Wright MC, Joyce GF (1999) Proc Natl Acad Sci USA 96:14712.
149. Ikawa Y, Tsuda K, Matsumura S, Inoue T (2004) Proc Natl Acad Sci USA 101:13750.
150. Duarte CM, Pyle AM (1998) J Mol Biol 284:1465.
151. Gregersen BA, York DM (2005) J Chem Phys 122:194110.
152. Wolfenden R, Ridgway C, Young G (1998) J Am Chem Soc 120:833.
153. Glennon TM, Warshel A (1998) J Am Chem Soc 120:10234.

# INDEX