

**Computational Methods for Understanding RNA
Catalysis: A Molecular Approach**

**A DISSERTATION
SUBMITTED TO THE FACULTY OF THE GRADUATE SCHOOL
OF THE UNIVERSITY OF MINNESOTA
BY**

Brian K. Radak

**IN PARTIAL FULFILLMENT OF THE REQUIREMENTS
FOR THE DEGREE OF
Doctor of Philosophy**

Darrin M. York

September, 2014

© Brian K. Radak 2014
ALL RIGHTS RESERVED

Acknowledgements

My path to graduate school started many years ago thanks to the inspiration and support of my peers and mentors in the Integrated Science Program at Northwestern University. Although there are too many to list, there are by no means too many in my life and I am continually enriched not only by our friendships but also those that they bring into my life.

The work put into this dissertation would not have been possible without the support and guidance of my advisor, Darrin York, as well as that of numerous members of his research group, especially (but in no particular order) Tai-Sung Lee, Timothy Giese, and George Giambaşu. I have also benefited immensely from discussions with our collaborators, particularly Michael Harris, Shantenu Jha, Emilio Gallicchio, David Case, and Jason Swails.

Lastly, I would have been hopelessly lost and likely not graduating at all were it not for the patience, skill, and all-around good nature of several departmental administrators, especially Nancy Thao and Janice Pawlo.

Dedication

To my editor, proof-reader, critic/skeptic, practice audience, pastry chef/taster, co-conspirator, co-pilot, dancing partner, and partner in all things - After seven years, four states, two time zones (not even counting summers!), and more plane, train, and automobile trips than I care to count, our next chapter is finally in sight.

Abstract

Molecular simulation is a powerful technology for providing a detailed picture of a wide range of chemical phenomena. The results of simulation studies are now increasingly used in supplementing experimental studies both as a predictive tool and as a lens through which to interpret results and generate new hypotheses. This dissertation describes several advancements in the development and application of molecular simulation methods to the study of RNA catalysis. Such reactions are representative of a broad class of chemistry associated with important biological functions including storage of genetic information, metabolism, and cell signaling and replication. Furthermore, the existence of naturally occurring RNA sequences that catalyze these reactions has significant implications for the origins of life and the potential design of new RNA based technologies. The work presented here offers new insights into these problems and contributes to a detailed, molecular understanding of the fundamental chemical principles that are in action.

Contents

Acknowledgements	i
Dedication	ii
Abstract	iii
List of Tables	vii
List of Figures	ix
1 Introduction	1
2 Background	3
2.1 RNA Enzymes	3
2.1.1 2'- <i>O</i> -Transphosphorylation	4
2.1.2 Comparisons with Protein Enzymes	8
2.1.3 Some Important/Unusual RNA Structural Motifs	10
2.1.4 The RNA Catalysis “Problem Space”	12
2.2 A Few Topics in Statistical Mechanics	14
2.2.1 Free Energy and Free Energy Landscapes	15
2.2.2 Integral Equation Formalisms	18
2.3 Molecular Simulations - Statistics and Sampling	21
2.3.1 Multistate Sampling	22
2.3.2 Enhanced Sampling	25

3	Methods Development for Molecular Simulation	29
3.1	Estimators for Free Energy Landscapes	30
3.1.1	Non-Parametric Estimators	30
3.1.2	Parametric Estimators	33
3.1.3	Comparison of Free Energy Landscape Estimation Methods	38
3.2	Asynchronous Replica Exchange	46
3.2.1	Synchronization Modes and Resource Utilization	48
3.2.2	Application to Multi-Dimensional Free Energy Manifolds	50
4	Computational Investigations of Phosphoryl Transfer in Non-Ribozymatic Systems	62
4.1	Molecular Simulations of RNA 2'- <i>O</i> -Transesterification Reaction Models in Solution	63
4.1.1	Introduction	64
4.1.2	Computational Methods	66
4.1.3	Results	71
4.1.4	Conclusion	86
4.2	Experimental and Computational Analysis of Ribonuclease A	87
5	Computational Investigations of Phosphoryl Transfer in Ribozymes	91
5.1	The Hepatitis Delta Virus Ribozyme	92
5.1.1	Background	92
5.2	A Framework for Assessment of Metal-Assisted Nucleophile Activation in the Hepatitis Delta Virus Ribozyme	96
5.2.1	Introduction	97
5.2.2	Computational Methods	100
5.2.3	Results	102
5.2.4	Discussion	109
5.2.5	Conclusion	116
	References	117

Appendix A. Supporting Information for: Molecular Simulations of RNA 2'-<i>O</i>-Transesterification Reaction Models in Solution	146
A.1 Umbrella Sampling	146
A.1.1 Detailed Equilibration Protocol	146
A.1.2 Overlap of Reaction Coordinate Distributions	147
A.1.3 Lennard-Jones Parameters	149
A.2 Kernel Density Estimation	152
A.2.1 Effect of Bandwidth Choice	152
A.2.2 Comparison with Histograms	153
Appendix B. Supporting Information for: A Framework for Assessment of Metal-Assisted Nucleophile Activation in the Hepatitis Delta Virus Ribozyme	155
B.1 Molecular Dynamics	155
B.1.1 Detailed Solvent Equilibration Protocol	155
B.1.2 Cluster Analysis	156
B.2 RISM and NLPB Calculations	159
B.3 Semiempirical Model Validation	160
B.3.1 Quantum Chemical Calculations	160
B.3.2 Benchmark Calculations on a Model for RNA Cleavage	161
B.4 Free Energy Surface Stationary Point Analysis	163

List of Tables

3.1	Examples of kernel density estimator forms.	32
4.1	Potential energy barriers and select geometric quantities at the transition state at different levels of theory using multiple reaction coordinates as in Figure 4.2.	73
4.2	Free energy profile extrema and select average geometric quantities from Figure 4.3.	76
4.3	Free energy profile extrema and select average geometric quantities from Figure 4.4.	78
4.4	Free energy profile extrema and select average geometric quantities from Figure 4.5.	80
4.5	Heavy atom KIEs for the leaving group, nucleophile, and non-bridge phosphate oxygens of UpG and model RNA substrates during catalysis by both solution and RNase A.	88
5.1	MM/3D-RISM and MM/NLPB/SA results for binding free energies and pK_a shifts in the hepatitis delta virus ribozyme under varying background concentrations of NaCl.	104
5.2	Analysis of free energy profiles of the HDVr catalytic reaction including normal mode analysis of stationary points and comparisons to previous studies.	107
A.1	Lennard-Jones parameters for atoms in the quantum region.	151
B.1	Results from cluster analysis of molecular dynamics trajectories of the HDVr in multiple states.	157
B.2	Quantum chemical calculations at various levels of theory for phosphoryl transfer in a model system in both the gas phase and solvent.	162

B.3 Parameters used for reduced mass determination of atom transfer coordinates.	163
--	-----

List of Figures

2.1	Chemical scheme for 2'- <i>O</i> -transphosphorylation of RNA via acid/base catalysis.	5
2.2	More O'Ferrall-Jencks diagram describing the mechanistic space of phosphoryl transfer reactions.	6
2.3	Chemical classification of standard protein backbone and side chain components with some biologically relevant pK_a values.	8
2.4	Chemical classification of ribonucleic acid backbone components and standard nucleobases with some biologically relevant pK_a values.	9
2.5	Nomenclature of nucleobase pairing by hydrogen bonding via interactions of the Watson-Crick, Hoogsteen, and/or sugar faces.	11
2.6	Some canonical and non-canonical RNA base pairs and their Leontis/Westhof classifications.	12
2.7	Schematic representation of the three-dimensional "problem space" of RNA catalysis.	13
2.8	Schematic representation of umbrella sampling via a localized bias potential to enhance sampling in high energy regions.	27
3.1	Performance of vFEP with different sample sizes when calculating the free energy profile of simple chemical reactions.	40
3.2	Performance comparison of several free energy methods with different sample sizes on a dense sampling grid when calculating the two-dimensional free energy landscape of alanine dipeptide.	42
3.3	Performance comparison of several free energy methods with different sample sizes on a sparse sampling grid when calculating the two-dimensional free energy landscape of alanine dipeptide.	43

3.4	Performance comparison of several free energy methods with different sampling densities when calculating the free energy profile of simple phosphoryl transfer reactions.	44
3.5	Performance comparison of several free energy methods with different sampling densities when calculating the two-dimensional free energy landscape of alanine dipeptide.	45
3.6	Evaluation of smoothing artifacts in non-parametric free energy methods versus vFEP for the two-dimensional free energy landscape of a phosphoryl transfer reaction.	47
3.7	Schematic of resource utilization in traditional REMD in which both MD and exchange occur <i>synchronously</i>	49
3.8	Schematic of resource utilization in a variation on traditional REMD in which MD occurs <i>asynchronously</i> but exchange occurs <i>synchronously</i>	50
3.9	Schematic of resource utilization in a REMD scheme in which both MD and exchange occur <i>asynchronously</i>	51
3.10	Schematic of resource utilization in a REMD scheme in which MD occurs on a fixed interval in real time (as opposed to simulation time) and exchange occurs <i>synchronously</i>	52
3.11	Schematic of dihedral angles used as bias coordinates during umbrella sampling of uracil.	55
3.12	Three-dimensional free energy manifold for solvated uracil along the Z_x , Z_y , and χ coordinates.	57
3.13	Free energy profiles for solvated uracil along its χ torsion using multiple schemes to reduce dimensionality.	58
3.14	Boltzmann weighted average free energy surface for solvated uracil along the Z_x and Z_y sugar pucker coordinates.	60
4.1	Reaction models for base-catalyzed phosphoryl transfer.	65
4.2	Gas phase potential energy profiles for the 2'- <i>O</i> -transesterification of 2-hydroxy ethyl methyl phosphate	72
4.3	Free energy profiles from simulation of an abasic RNA dinucleotide in different solvent environments.	75

4.4	Free energy profiles from simulation of an abasic RNA dinucleotide with different force field models.	77
4.5	Free energy profiles from simulation of several phosphoryl transfer models in solution.	79
4.6	Snapshots from simulations of an abasic RNA dinucleotide near the transition state.	83
4.7	Transition state structures for a model reaction for non-enzymatic RNA cleavage and enzymatic cleavage by ribonuclease A.	89
5.1	Proposed mechanism and schematic model for the self-cleavage reaction of the hepatitis delta virus ribozyme utilizing a Mg^{2+} ion bound at G1.	93
5.2	Proposed mechanism and schematic model for the self-cleavage reaction of the hepatitis delta virus ribozyme utilizing a Mg^{2+} ion bound at a G-U wobble pair.	94
5.3	Schematic diagram of the catalytic mechanism for the HDVr proposed by Chen, <i>et al.</i>	99
5.4	Reaction scheme for (de)protonation of the HDVr at U-1:O2' and Mg^{2+} binding events needed to attain catalytically active states with and without Mg^{2+}	101
5.5	Na^+ pair distribution functions around the HDVr computed via NLPB and 3D-RISM-PSE3 compared with peak positions from MD and volmap.	103
5.6	Free energy surfaces of the HDVr catalytic reaction starting from an activated (post base) state with and without Mg^{2+} bound at the position hypothesized by Chen, <i>et al.</i>	106
5.7	Reaction profiles and key geometric quantities along the minimum free energy paths for the HDVr catalytic reaction starting from an activated state with and without the presence of Mg^{2+} at the position hypothesized by Chen, <i>et al.</i>	108
A.1	Representative plots of reaction coordinate distributions from simulations of different phosphoryl transfer reaction models.	148
A.2	Plots demonstrating the effect of increasing/decreasing the bandwidths on the calculated free energy profile.	152

A.3	Plots comparing the use of a normal kernel density estimator versus a more traditional histogram estimator when calculating a free energy profile.	153
B.1	RMSD distributions of molecular dynamics trajectories of the HDVr in multiple states.	158
B.2	Comparison of quantum chemical models on a model system for RNA backbone cleavage.	160

Chapter 1

Introduction

The goal of this dissertation is to develop and apply molecular models to gain new insights into the biological problem of unraveling the mechanism of RNA catalysis. More specifically, it describes the use of computational modelling via molecular simulations, quantum chemical calculations, and numerical integral equation methods. The ultimate goal of these approaches is to provide a detailed structural and dynamical description of chemical processes at a molecular level and to use this as both a lens for interpretation of experimental data as well as a predictive tool. Although these techniques are widely useful (and used) in many areas of chemical research, the main problem of interest in the present work is the elucidation of mechanistic details in RNA enzymes, a broad class of non-coding RNA which have significant implications for the function of RNA in human biochemistry, the origins of life (the “RNA world hypothesis”), as well as for the possibility of engineering new biomolecules for a myriad range of applications (ribozyme engineering).

A rough description of this work is as follows. Chapter 2 summarizes the theoretical background that serves as the foundation for many of the methods developed and applied later in the thesis, as well as provide a review of the relevant literature. Section 2.1 describes the general space of chemical reactions occupied by RNA enzymes (specifically phosphoryl and proton transfer) as well as existing theoretical frameworks for discussing such mechanisms. This is followed by a brief review of the chemical functional groups present in standard RNA constructs as well as some important structural motifs and nomenclature that may not be as widely known. This discussion is concluded by a

summary of the main open questions that mechanistic studies of RNA enzymes seek to answer. In Sections 2.2 and 2.3 a number of important theoretical and practical aspects of molecular simulation are discussed with a specific focus on the models, methods, and techniques utilized in the applications that follow. As before, the intention is to briefly, but formally, introduce concepts that may not be widely known.

In Chapter 3, several original developments and contributions of both primary (Refs. 1, 2 and 3) and secondary (Refs. 4 and 5) authorship are presented. These works, in whole or in part, concern methodological problems in molecular simulation, principally those involving statistical analysis and sampling. This includes explorations of the problem of constructing estimators for free energy landscapes, both by non-parametric and parametric means, as well as novel, asynchronous approaches to replica exchange simulation schemes that are particularly suitable to large numbers of replicas and problems requiring the estimation of multi-dimensional free energy manifolds. Some proof of concept applications are also included as well as the results of benchmarking and performance studies.

In Chapters 4 and 5 the mechanism of phosphoryl transfer is studied in detail, and again several original contributions on the subject are presented. These begin with some foundational computational studies on non-enzymatic phosphoryl transfer reactions[1] followed by some collaborative experimental work on the protein enzyme ribonuclease A[6]. These studies lay a strong foundation on which to benchmark and validate the models and methods used in Chapter 5 to investigate phosphoryl transfer in the hepatitis delta virus ribozyme[7] and to establish a broader perspective on general catalytic strategies.

Overall this work represents a broad, multiscale molecular modelling approach to an important class of problems in biocatalysis. The tools and approaches employed and developed herein are significant and widely transferable into other biochemical problem spaces. The published works described have been well received in the peer-reviewed chemical literature and signify a modest, but firm step towards a deeper understanding of chemical systems in general.

Chapter 2

Background

2.1 RNA Enzymes

RNA enzymes or “ribozymes” were a relatively surprising discovery in the mid-1980’s[8, 9]. The continual discovery of new and ever more ubiquitous ribozyme sequences is now a testament to their enduring impact on modern biochemistry. The fact that such molecules could perform chemistry was a direct contradiction to the standing “central dogma,” which held that RNA was predominantly a messenger molecule, acting as a temporary mediator between the more robust nuclear DNA (the site of long-term data storage) and the protein machinery (the presumed “business end” of biochemistry). Perhaps even more astoundingly, this paradigm shift continued to manifest itself in the discovery that several major cell complexes, including the ribosome[10, 11] and the spliceosome[12, 13], contained catalytic cores largely or entirely composed of RNA. Although still only representative of a small portion of non-coding RNA sequences, ribozymes are now known to be widely distributed amongst nature[14, 15, 16] (including the human genome[17, 18]) and serve critical roles in important biochemical processes. This realization has generated significant interest in the “RNA world hypothesis,” which posits that the early biological world may have consisted entirely of RNA[19, 20, 21].

Naturally occurring ribozymes appear to exclusively deal with reactions involving phosphoryl transfer. Although efforts in *ribozyme engineering* have moved passed this apparent restriction[22], the majority of biochemical interest remains in the realm of

phosphoryl transfer, and with good cause, such reactions are critical in several biological processes, including metabolism, replication, and cell signaling[23]. The largest ribozymes (and the ones most firmly integrated into eukaryotic biochemistry) may or may not be autocatalytic. That is, phosphoryl transfer may or may not occur within the same RNA strand containing the catalytic core.¹ Nonetheless, this intramolecular catalysis has considerable parallels with intermolecular catalysis and so the former is often considered as a simpler mechanistic model for the latter. Conveniently, the smallest ribozymes, usually associated with viruses or prokaryotes[24, 25], fall in this category and are amenable to detailed mechanistic studies, both experimental and computational. As such, these systems will be the primary focus of this work.

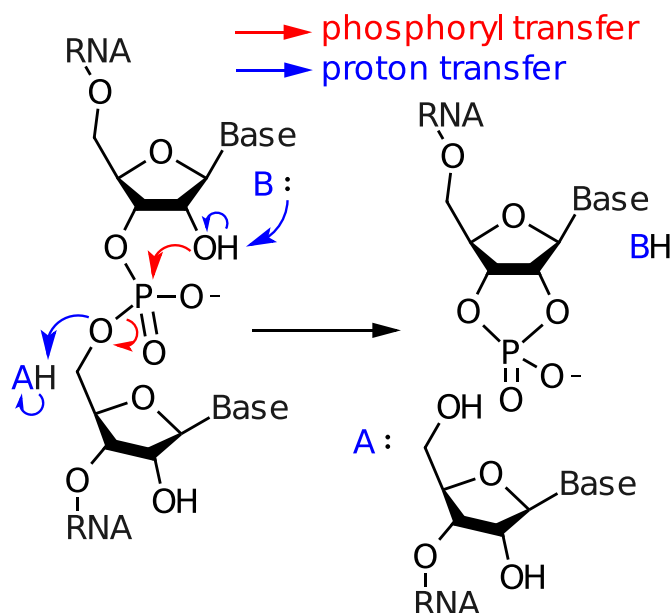
2.1.1 2'-*O*-Transphosphorylation

The primary form of intramolecular RNA catalysis is 2'-*O*-transphosphorylation. In this reaction the 3',5'-phosphate RNA linkage undergoes cleavage via nucleophilic attack by the O2' position (Figure 2.1). The immediately resulting product is a 2',3'-cyclic phosphate (although this may undergo further reaction). On either end of this phosphoryl transfer reaction are two proton transfer reactions, one to activate the nucleophile and another to stabilize the leaving group. The critical tasks for a successful ribozyme are thus to facilitate these two proton transfers and promote the intervening phosphoryl transfer. As might be expected, the specific requirements for the activation and stabilization steps are highly correlated with the energetics and geometry of the transition state for phosphoryl transfer. Thus, before considering these processes in too heavy of detail, it is worth laying out the possibilities for what these structures might look like and how they might be characterized.

A convenient method for analyzing related transition state geometries for a variety of chemical reactions is the use of a More O'Ferrall-Jencks (MOJ) diagram[26, 27, 28]. Such a diagram compares the relative progression of the bond forming and breaking processes thereby allowing for a graphical representation of the early/late transition state

¹ Strictly speaking, this may not be considered proper autocatalysis, as completion of the reaction may not regenerate the original molecule. Nonetheless, some ribozymes in this class can be engineered into proper catalysts with only modest changes. In both cases, however, the kinetic and thermodynamic character of the underlying chemical reaction is modified and so this minor abuse of terminology can be overlooked.

Figure 2.1: In the general chemical scheme for 2'-*O*-transphosphorylation of RNA via acid/base catalysis, the core phosphoryl transfer reaction consists of nucleophilic attack by O2' on phosphorous with O5' as the leaving group. Both O2' and O5' are highly reactive when charged and so a proton transfer reaction with a base or acid is necessary to activate or stabilize them, respectively.

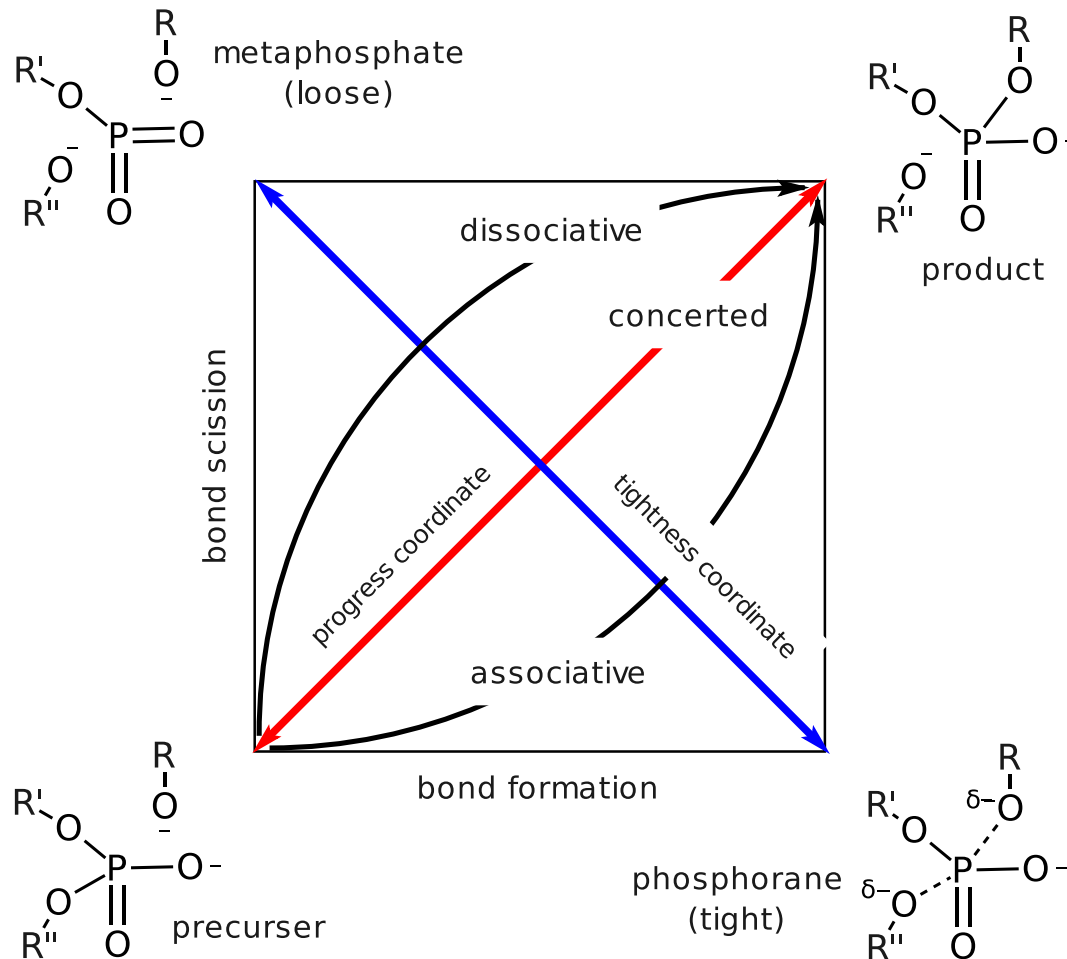


dichotomy from Hammond-Leffler type arguments. For 2'-*O*-transphosphorylation reactions in RNA the graph axes correspond to the progression of bond formation between O2' (the *nucleophile*) and phosphorous and the degree of bond scission between O5' (the *leaving group*) and phosphorous (Figure 2.2)[29]. Clearly, when one of these bonds is formed while the other is not, the reactant and product states are regained (Figure 2.2, bottom left and top right, respectively). The vector interpolating between these two corners is often called the *progress coordinate* while, for reasons that will be discussed below, the orthogonal direction is called the *tightness* or tightness coordinates. Any possible structure can thus be characterized by these two coordinates and a mechanistic pathway can be drawn as a parametric curve.

Various chemical reactions of the same general class will differ in their characteristic pathway on an MOJ diagram and it is useful to categorize these pathway types. If the bond forming and breaking processes occur in equal measure, then the reaction is said

to be *concerted* and the reaction progresses directly along the diagonal. Alternatively, if the bond breaking process proceeds without significant formation of the forming bond, then the pathway is called *dissociative* (Figure 2.2, top curved path). Finally, if bond formation occurs quickly without bond scission proceeding, then the reaction is called *associative* (Figure 2.2, lower curved path).

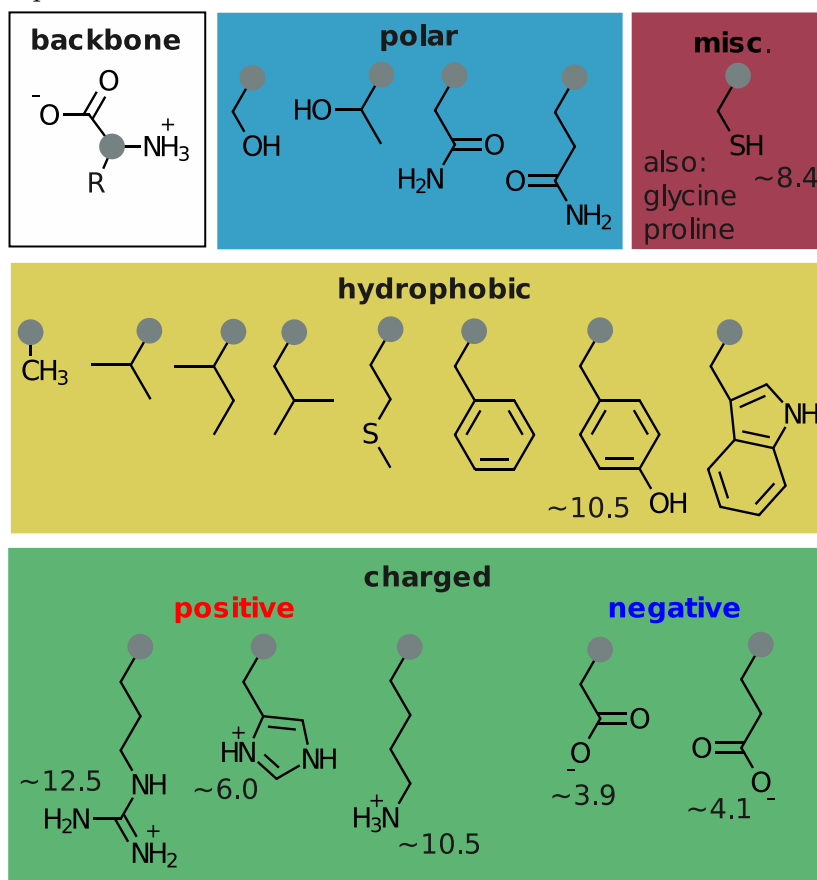
Figure 2.2: The mechanistic space of phosphoryl transfer can be neatly summarized by a More O'Ferrall-Jencks diagram. Reaction progress in this space is described by two coordinates, the progress and tightness coordinates, which are themselves combinations of bond formation and scission coordinates; all possible reaction paths and transition states can be identified by these two coordinates.



Within the framework of transition state theory, any chemical pathway can be largely understood in terms of its transition state structure (*i.e.* the structure of highest free energy). Even if two reactions follow the same or similar paths, it is possible for their transition state structures to be quite different in terms of geometry and thus the reaction characteristics and behavior could actually be quite different. Here it is useful to return to the tightness coordinate and examine the limiting structures at the other corners of the MOJ diagram for phosphoryl transfer. In the top left is a purely dissociative transition state structure where charge is highly localized on both the nucleophile and leaving group and phosphorous is weakly bound compared to the reactants and products (this potentially stable compound is called a *metaphosphate*). Such a structure is often characterized as “loose” because the atoms are less tightly held than at the endpoints; it represents one extreme of the tightness coordinate. At the lower corner Figure 2.2 the opposite extreme is observed. Phosphorous is bound to five oxygens (*i.e.* it is pentacoordinate) and negative charge is delocalized over this complex which, like a metaphosphate, can also potentially be stable, although usually with the aid of nearby positive charge or proton donating group. Such a structure is characterized as “tight” due to the large number of bonds holding the phosphorous atom in place. Lastly, it should be noted that the familiar Hammond-Leffler classification of “early” (reactant-like) and “late” (product-like) transition states still holds and, moreover, now has a distinct graphical interpretation in terms of distance along the progress coordinate.

A general vernacular for discussing 2'-*O*-transphosphorylation reactions has thus been thoroughly developed. A simple and elegant reduction of any such reaction can be obtained by characterizing its transition state within the MOJ diagram language. However, it may be sorely noticed that no specific route to obtaining information regarding either the pathway or the transition state has been provided. This is of course the great challenge of mechanistic studies. In the present work it will suffice to recognize that MOJ diagrams represent a specific coordinate space for free energy landscapes (discussed in Section 2.2.1). Conveniently, a wide array of molecular simulation methods are available for constructing such landscapes and this approach is heavily relied upon in Chapters 4 and 5. Any particulars are thus deferred until then.

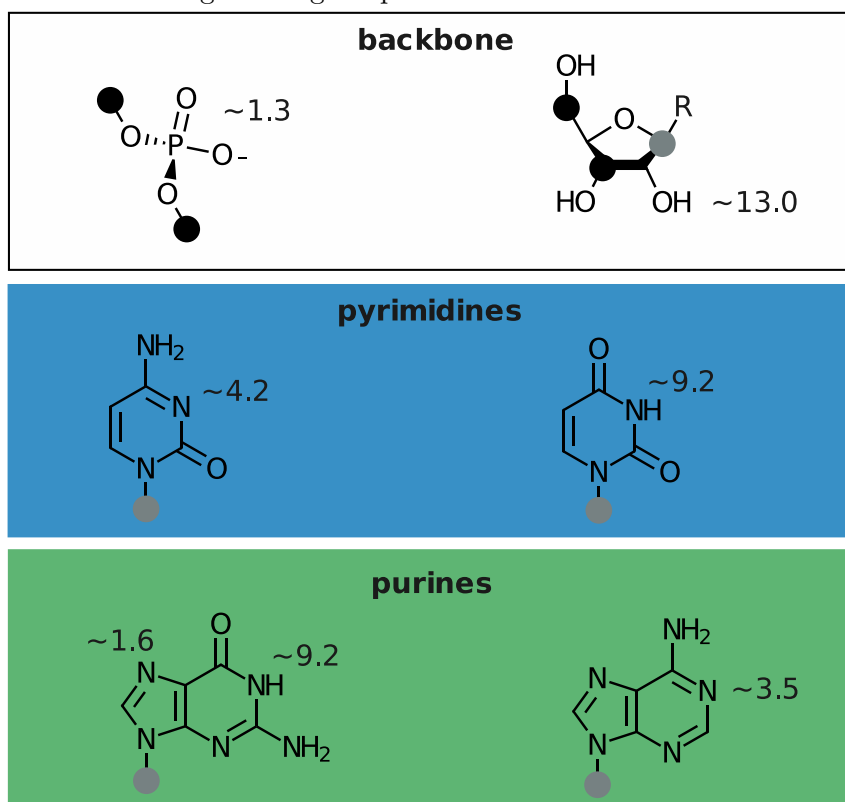
Figure 2.3: The backbone and side chain components of standard amino acids can be classified into several categories that highlight their chemical diversity. Across and within these groups, a wide range of functional groups possess biologically relevant pK_a values near pH 7-8.



2.1.2 Comparisons with Protein Enzymes

Given that the structure/function relationships of protein enzymes have a much longer history and are much better known than those of ribozymes, it is useful to build upon this body of knowledge when studying RNA enzymes. Specifically, it is instructive to assess the wide array of chemical functionality available amongst the unmodified amino acid side chains from which naturally occurring protein enzymes are built. In Figure 2.3 these components are broadly divided into four groups plus the ubiquitous amino acid backbone. It is rather plain to see that even within these groups, the side chains vary

Figure 2.4: The backbone and nucleobase components of standard nucleic acids are classified into fewer and narrower categories than their amino acid counterparts. In addition, only a few functional groups possess pK_a values between 1-14 and, as is, these are well outside the biological range of pH 7-8.



considerably in size, composition, and chemical characteristics, especially with regard to pK_a values. The last of these play a critical role in acid/base chemistry and are most chemically versatile when near biological pH values of $\sim 7-8$. Furthermore, this list does not include the additional wide range of protein cofactors such as vitamins and coenzymes, including small molecule metal binding motifs such as porphyrins. The ability of protein enzymes to catalyze a broad range of chemistry is thus not difficult to comprehend.

By way of contrast, in Figure 2.4 it can be seen that the components of nucleic acids display a relatively limited set of chemical characteristics. The two groups of nucleobases (purines and pyrimidines) have both extreme (*i.e.* far from pH 7) and

redundant pK_a values. Furthermore, there is the problem of the negatively charged phosphate backbone, which complicates higher order organization of long chains due to electrostatic repulsion.

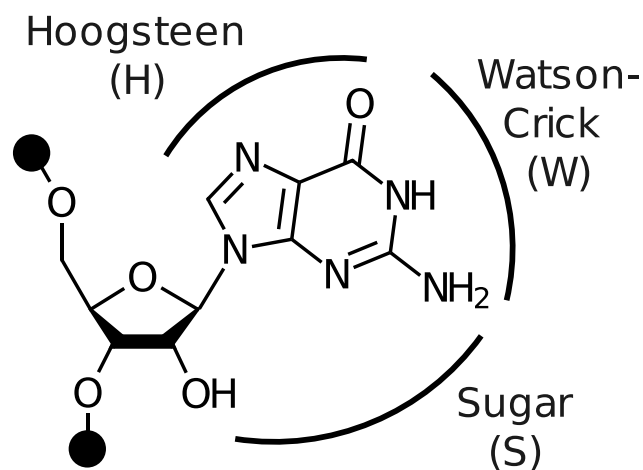
Nonetheless, since it is clear that RNA enzymes do exist, the question that these compounds should elicit is not “*can* they facilitate catalysis?” but rather “*how* do they facilitate catalysis?” Clearly the whole of RNA becomes greater than its parts and higher order organizations of these components must lead to altered and enhanced chemical functionality. As such, it is now broadly agreed that a number of catalytic strategies are in effect. Firstly, and this is true for nucleic acids in general, carefully maintained ionic environments are used to stabilize and fold the phosphate backbone. Second, specialized RNA folds must, in some fashion, shift the pK_a of select nucleobases in order to obtain biochemically relevant values closer to neutrality. Finally, carefully engineered regions must, either directly or by recruitment of metal ions, stabilize electrostatically strained structures, such as the phosphoranes found in Figure 2.2.

2.1.3 Some Important/Unusual RNA Structural Motifs

A number of unusual motifs found in RNA have only been codified in the last few decades and it is worth conveying these developments here[30, 31]. Firstly, the reader is expected to be familiar with the famous nucleobase pairing scheme (via two or three hydrogen bonds) utilized by Watson and Crick in the structure of the DNA double helix[32]. However, while these hydrogen bonding geometries are by far the most prevalent in RNA, they are by no means the only ones. Hydrogen bonds with at least three different “faces” have now been observed (via crystallography) for most or all of the different nucleobases and many possible combinations were enumerated by Leontis and Westhof[31]. In addition to the canonical Watson-Crick (W) face, the Hoogsteen (H) and sugar (S) faces are now frequently considered (Figure 2.5). Furthermore, the orientation of the nucleobase in either the *anti* (with the Watson-Crick face away from the phosphate) or *syn* positions (with the Watson-Crick face towards the phosphate) can give rise yet more complicated hydrogen bonding patterns between RNA strands.

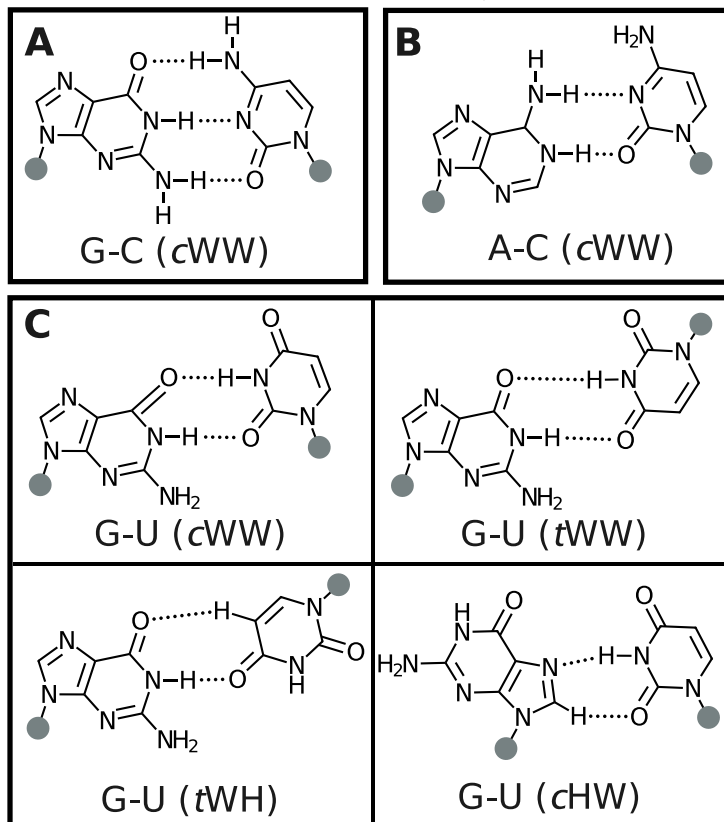
A small set of related base pairs are shown in Figure 2.6 and show how the canonical scheme (top left) for guanine and uracil relates to several non-canonical schemes. These schemes are labeled first by the identity of the pairs (here G-C, G-U, and A-C) and then

Figure 2.5: Nucleobases possess multiple faces capable of hydrogen bonding interactions, usually with other nucleobases. Although canonical interactions occur only with the Watson-Crick (W) face, other interactions can occur with the Hoogsteen (H) and sugar (S) faces as well.



by the involved hydrogen bonding face (W, H, or S). Since the *syn/anti* dichotomy loses specificity when considering two strands of potentially different orientations, the bases can be collectively referred to as either *cis* (*c*) or *trans* (*t*, on the same or opposite sides, respectively). Within this nomenclature it can be seen that the canonical base pairs are classified as being *cis* Watson-Crick/Watson-Crick or *cWW* base pairs. Of course, these are not the only *cWW* pairs, as evidenced by the non-canonical *cWW* A-C and G-U pairs (2.6). The range of bonding patterns available to G-U pairs are particularly impressive and include a number of strand orientations and hydrogen bonding faces (*N.B.* not all possible G-U pairs are shown here, see Ref. 31 for an exhaustive list). This can be especially important in loosely held regions of RNA folds, where interconversion between these forms may be possible. For example, the presentation of different faces to solvent can lead to molecular switching behavior when metal ions or cofactors are involved. Overall, it should be clear that the availability of non-canonical base pairs greatly adds to the flexibility of the RNA folding landscape and one must consider such possibilities when investigating overall RNA folds.

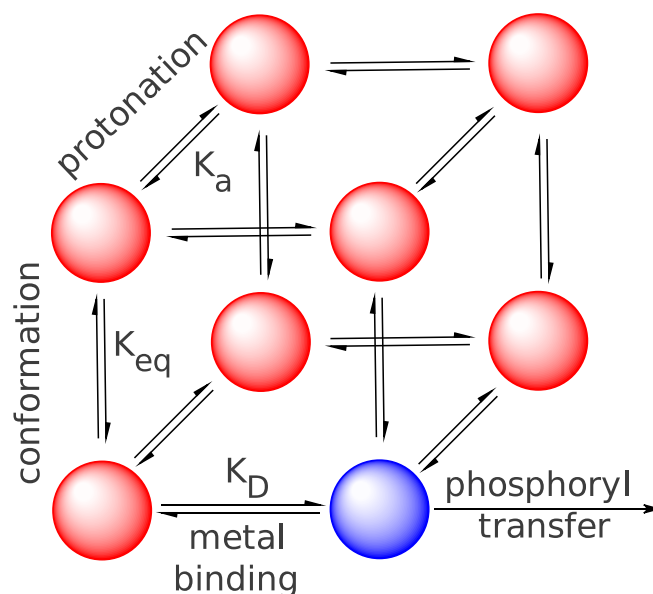
Figure 2.6: In addition to the canonical Watson-Crick/Watson-Crick (WW) base pairs, a large number of non-canonical pairs are also possible. These have been classified and enumerated by Leontis and Westhof[30, 31] by the relative orientation of the base attachments (on the same, *cis*, or opposite, *trans*, sides) and the nucleobase face involved (Watson-Crick, Hoogsteen, or sugar, see Figure 2.5).



2.1.4 The RNA Catalysis “Problem Space”

As a summary of the preceding sections, the “problem space” associated with RNA catalysis can now be discussed. The dominant issue is how such catalysts facilitate the underlying phosphoryl transfer reaction. In Section 2.1.1 it was seen that these reactions inherently pass through a highly charged and electrostatically strained transition state structure and that such a structure comes from similarly strained precursor and product structures that must be activated and stabilized, respectively, by at least a proton transfer event. However, in Section 2.1.2 it was shown that nucleic acid functional groups

Figure 2.7: The “problem space” of RNA catalysis can be schematically represented as an equilibrium in three-dimensions. That is, the question can be reduced to producing a catalytically competent state via proper alignment of changes in conformation, protonation state, and metal binding mode. It is only when these aspects, which may themselves be correlated, come together that the core phosphoryl transfer reaction can proceed.



span a rather limited set of chemistry for performing these proton transfers, especially compared to amino acids. Then, in Section 2.1.3 some examples of the unusual structural motifs available to RNA were discussed, especially with regards to non-canonical base pairings. These hinted at unusual conformations that might augment the character of the available functional groups. A final idea that is introduced now is the notion of interactions with electrolytes (*i.e.* ions and salt), as these are not only critical in providing structural stabilization of nucleic acids, but also may form specific interactions with phosphate and/or nucleobase functional groups. Taken together, all three of these aspects, chemically augmented functional groups, specialized conformations, and metal ion interactions, compose a broad space in which RNA can create circumstances favorable to catalysis.

A diagram of the RNA problem space is shown in Figure 2.7. Each axis represents each of the separate chemical processes potentially required to prepare the system for

catalysis of phosphoryl transfer. The most important point to be noticed is that there are multiple nodes on the diagram in which the system may be in a favorable state along on axis (*e.g.* conformation) but not along another (*e.g.* protonation state). In other words, although these processes may all need to be aligned simultaneously for a catalytically active state to be reached, their individual oscillations/fluctuations between active and inactive modes need not occur simultaneously. For example, a necessary metal binding event may occur rapidly, regardless of the conformation of the RNA, while the folding event(s) needed to initiate catalysis may occur slowly. This is a case of little/no correlation between degrees of freedom. However, it is also possible that two axes are correlated or anti-correlated. For example a conformational change may enhance the frequency with which a critical protonation event occurs or it may inhibit that event. Overall, catalytic RNA systems may have similar or distinct couplings in this problem space. Investigating and characterizing the specific signatures of a wide variety of systems (both ribozymatic and non-enzymatic) will lead to fundamental insights into how these reactions occur, how they can be intervened upon, and how they can be engineered (either repurposed or designed *ex nihilo*).

2.2 A Few Topics in Statistical Mechanics

Statistical mechanics is the branch of physics concerned with linking the behavior of *microscopic* systems with *macroscopic* observations of them (sometimes called the *Gibbs postulate*). The justification of this connection is complex and will remain outside the scope of this work.² Regardless, a basic result is that microscopic systems (here always assumed to be chemical/molecular systems) abide by probabilistic tendencies and their aggregate average behavior informs the observations we make. This notion can be distilled into a surprisingly concise expression, often known as the Boltzmann relation.

$$g(\mathbf{x}) \propto e^{-u(\mathbf{x})} \tag{2.1}$$

This statement relates the (dimensionless) probability density, $g(\mathbf{x})$ for various microscopic configurations of a system to a molecular model, $u(\mathbf{x})$, describing how the

² At least in the chemical literature, the textbook treatments by McQuarrie[33] and Hill[34] are often considered seminal, although many other excellent books are also available.

elements of that system interact.³ In practice $u(\mathbf{x})$ is a complicated expression and can take many forms depending on the qualities of the system (*e.g.* does it exchange energy with its surrounding, compose a fixed volume, chemically interact with a solvent, *etc.*). Here $u(\mathbf{x})$ is taken to be a physics-based (dimensionless) energy expression. This need not be true, but is general enough for most applications in which systems of atoms/molecules are being modelled. In any case, $u(\mathbf{x})$ will be referred to as the “reduced potential,” since it reduces many possible types of interactions (*i.e.* potentials) into a generic, compact expression. While this notation runs some risk of oversimplifying matters it has the advantage of retaining many of the essential features, namely that configurations with larger values of $u(\mathbf{x})$ (*i.e.* *higher* energies) are found with *lower* probability.

2.2.1 Free Energy and Free Energy Landscapes

Free energy is a foundational concept in both statistical mechanics and thermodynamics. In a basic sense, it encodes much of the information in the Boltzmann relation into a single value. For simplicity, we shall assume that systems have a continuous distribution of configurations. This is equivalent to the (often extremely good) approximation that the system behaves classically and does not experience discrete states due to quantum effects. In this case the probability of all configurations must integrate to one. This constraint on Eqn. 2.1 introduces a normalization constant which can be related to the (dimensionless) free energy, f .

$$\int g(\mathbf{x})d\mathbf{x} = 1 \quad \Rightarrow \quad \int e^{-u(\mathbf{x})}d\mathbf{x} \equiv e^{-f} \quad (2.2)$$

$$g(\mathbf{x}) \equiv e^{f-u(\mathbf{x})}$$

³ The informed reader may find this notation lacking, at least without detailed clarifications. First, $g(\mathbf{x})$ is not a proper probability density, which of course must bare units (in this case mass \times distance² \times time⁻¹). However, the proper density and this “pseudo-density” are related by a frequently uninteresting constant scale factor proportional to h^{3N} , where h is Planck’s constant and N is the number of particles in the system. Other constants may also be necessary to account for indistinguishable or identical particles. Second, here “configuration” is broadly meant to mean the way in which the atoms or molecules in the system are behaving at a point in time, not just their locations in space. The momenta of particles can often be canceled out and so equations are frequently written in terms of just the coordinates and *these* are more naturally referred to as configurations. Unfortunately, due to Liouville’s theorem, this disrupts the mathematically appealing aspect of having dimensionless quantities. As a compromise, this language is retained while clarifying that it is meant somewhat loosely.

This result clearly relates the free energy, reduced potential, and probability density. It should be noted that f is just a number, while $u(\mathbf{x})$ and $g(\mathbf{x})$ are functions of the configuration. Often f itself is of physical interest and can be related (via phenomenological arguments) to quantities describing, for example, chemical equilibrium and phase transitions[33, 35]. These connections are useful and interesting, but will only be elaborated upon here when necessary.

In the present work, we will mostly be interested in averages with respect to the distribution of configurations. In general, a microscopic observable $A(\mathbf{x})$, whose specific value depends on the configuration, has an average value (often called an “ensemble average”), $\langle A \rangle$, that can be written as an integral over the probability density.⁴

$$\langle A \rangle = \int A(\mathbf{x})g(\mathbf{x})d\mathbf{x} \quad (2.3)$$

The angle brackets $\langle \cdot \rangle$ is commonly used to indicate an integral weighted by the probability density and this notation is used here. As will be shown below, $A(\mathbf{x})$ can take a wide array of forms and represents one of the most powerful aspects of statistical mechanics.

An interesting variation on free energy is the construction of so-called free energy landscapes or manifolds describing important changes in a system. This insight is often attributed to Kirkwood[36]. The basic idea is to begin with Eqn. 2.2 and continue integrating both sides until only select degrees of freedom remain. In the simplest treatment, the configuration element, $d\mathbf{x}$, is assumed to derive from N atoms each labelled by an index, i . That is, $d\mathbf{x} = \prod_{i=1}^N d\mathbf{x}_i$. The (dimensionless) free energy landscape of an arbitrary atom (atom 1, say), $f(\mathbf{x}_1)$, is then, by definition, related to its *marginal* distribution, $g(\mathbf{x}_1)$.

$$g(\mathbf{x}_1) = \int d\mathbf{x}_2 \dots \int d\mathbf{x}_N g(\mathbf{x}) \equiv e^{-f(\mathbf{x}_1)} \quad (2.4)$$

Note that this is not an ensemble average as in Eqn. 2.3, because it does not include an integral over all configurations; the various configurations of atom 1 have been left out.

⁴ The term “ensemble” is unfortunately often abused. It variously refers to the set of all possible configurations of a system, a particular statistical sample of those configurations, or even the specific form of 2.1 that is being used.

However, it can be made into the form of an ensemble average by multiplying each side by a delta function and completing the integration.⁵

$$\int d\mathbf{x}'_1 \delta(\mathbf{x}_1 - \mathbf{x}'_1) g(\mathbf{x}_1) = \int d\mathbf{x}'_1 \delta(\mathbf{x}_1 - \mathbf{x}'_1) \int d\mathbf{x}_2 \dots \int d\mathbf{x}_N g(\mathbf{x})$$

$$g(\mathbf{x}_1) = \langle \delta(\mathbf{x}_1 - \mathbf{x}'_1) \rangle \quad (2.5)$$

Combining this manipulation with Eqn. 2.4 and rearranging yields an ensemble average definition for the free energy landscape.

$$f(\mathbf{x}_i) = -\ln g(\mathbf{x}_i) = -\ln \langle \delta(\mathbf{x}_i - \mathbf{x}'_i) \rangle \quad (2.6)$$

Here the index has been changed to the more generic i to emphasize that this derivation/definition holds equally well for any atom or group of atoms (indeed we could label the atoms in any arbitrary order). More complicated coordinates could also be chosen, such as functions of multiple atomic coordinates, but this adds some slight complications to the delta function step (*i.e.* the possible addition of a Jacobian term). The adjustments are subtle and can be quite important, but will not be discussed in detail here.

A deeper understanding as to why the free energy landscape was derived the way it is can be obtained by briefly returning to Eqn. 2.4, rearranging to get an expression in terms of $f(\mathbf{x}_1)$, and then differentiating both sides with respect to \mathbf{x}_1 .

$$f(\mathbf{x}_1) = -\ln \int d\mathbf{x}_2 \dots \int d\mathbf{x}_N g(\mathbf{x}) = -\ln \int d\mathbf{x}_2 \dots \int d\mathbf{x}_N e^{f-u(\mathbf{x})}$$

$$\nabla_1 f(\mathbf{x}_1) = \frac{\int d\mathbf{x}_2 \dots \int d\mathbf{x}_N \nabla_1 u(\mathbf{x}) e^{f-u(\mathbf{x})}}{\int d\mathbf{x}_2 \dots \int d\mathbf{x}_N e^{f-u(\mathbf{x})}} = \frac{\int d\mathbf{x}_2 \dots \int d\mathbf{x}_N \nabla_1 u(\mathbf{x}) g(\mathbf{x})}{\int d\mathbf{x}_2 \dots \int d\mathbf{x}_N g(\mathbf{x})}$$

Once again, the last line is almost an ensemble average except that integration over \mathbf{x}_1 is missing. The interesting point here is that the negative derivative/gradient of $u(\mathbf{x})$ with respect to atom 1 is simply the *force* on atom 1. Integrating over all other atoms weighted by the probability density then gives the *mean force*. The free energy

⁵ Unfortunately there is little room here to expand adequately upon the nature of delta functions. For our present purposes it can merely be seen as a device for performing an integration and returning a function instead of a number. An extensive exposition of this and other aspects of distribution theory can be found in many texts on partial differential equations, such as that by Stakgold[37].

landscape, as defined here, can thus be viewed as a potential whose negative gradient gives this mean force. As such, $f(\mathbf{x}_1)$ is often called the “potential of mean force” or PMF.⁶ This device is extremely useful as it $f(\mathbf{x}_1)$ encodes all of the effects from the other atoms in the system into a simple, compact expression in terms of \mathbf{x}_1 only. What is more, it can be calculated from a single ensemble average given by Eqn. 2.6. The widespread utility of this framework will be expanded upon in what follows.

2.2.2 Integral Equation Formalisms

Section 2.2.1 introduced the concept of free energy landscapes in terms of marginal distribution functions describing a few degrees of freedom. Such expressions are considerably easier to manage than the high dimensional integrals needed to solve for the free energy directly and many theories have sought to exploit this. The main idea is to separate the many body integrals in expressions like Eqn. 2.2 into a hierarchy of lower dimensional integrals in a concerted fashion and via a select ansatz. There are of course many ways to proceed and a detailed overview of roughly half a century worth of attempts can be found in books such as those by McQuarrie[33] and Hansen and McDonald[38]. The simplest such integrals describe the configuration of only two atoms and are thus known as *pair distribution functions* (PDFs).⁷

$$g(\mathbf{x}_{12}) \equiv g(\mathbf{x}_1, \mathbf{x}_2) = \int d\mathbf{x}_3 \dots \int d\mathbf{x}_N g(\mathbf{x}) \quad (2.7)$$

A useful trick here is that one of the atoms can be arbitrarily taken as the origin and the PDF then viewed as a three-dimensional function in terms of the relative, rather than absolute, positions (denoted here as \mathbf{x}_{12}). In physical terms, the PDF can be nominally considered a measure of interaction or correlation between two atoms and its behavior generally consists of large and rapid changes at short distances and gradual decay to one at long distances. Another function is thus often defined by simply shifting this asymptotic behavior to zero and is called the *total correlation function*. This is because

⁶ Note that the wording here has been chosen quite carefully. When making more complicated coordinate choices beyond those for a single atom, it is easy to imperil the validity of this statement, especially when applying the definition of $f(\mathbf{x}_1)$ in Eqn. 2.6. In many cases an extra term appears that is separate from $f(\mathbf{x}_1)$ and it is only the combination of these two terms that is properly the PMF.

⁷ Once again the notation used here is not necessarily standard, but retains the most important concepts and is mathematically uncluttered.

it describes the total amount of correlation between two atoms which, intuitively, should vanish as those atoms become infinitely far apart.

$$h(\mathbf{x}_{12}) \equiv g(\mathbf{x}_{12}) - 1 \quad (2.8)$$

Reference Interaction Site Model

The reference interaction site model[39, 40, 41, 42, 43] (RISM) is one possible formalism for generating approximate solutions to the equations described above in Section 2.2. This is done via three main approximations/ansatzes. The first of these is the Ornstein-Zernicke equation[33, 38, 39], which proposes a decomposition of the PDF into two parts, a direct part and an indirect part.

$$h(\mathbf{x}_{12}) = c(\mathbf{x}_{12}) + \int c(\mathbf{x}_{13})h(\mathbf{x}_{13})d\mathbf{x}_3 \quad (2.9)$$

This separation, which is borne purely out of convenience, requires the definition a new function, $c(\mathbf{x}_{12})$, called the *direct correlation function*. The purpose of this function is to divide/classify the correlations of two atoms into those resulting from direct interactions depending on their separation (hence the name) and those which are mediated by a third particle (hence the appearance of the subscript 3 in Eqn. 2.9). Although not explicitly made obvious here, these new functions are mathematically nice and have some appealing properties for numerical solution. However, this formalism has increased the number of functions which need to be solved for from one ($g(\mathbf{x}_{12})$) to two ($h(\mathbf{x}_{12})$ and $c(\mathbf{x}_{12})$). This creates the requirement for a so-called *closure* relation, which makes the solution well-defined within certain requirements. An in depth discussion of closure relations is beyond the present work, but it should suffice to simply note that in order to proceed to a RISM solution a closure *must* be specified and that the accuracy of the solution depends upon it. Popular closures in historical or common use include the Percus-Yevick equation[44], the hypernetted-chain equation[38] (HNC), the mean sphere approximation (MSA), the Kovalenko-Hirata closure[42] (KH, a synthesis of the HNC and MSA closures), and the partial series expansion closure[45] (of which the KH closure is a special case).

The other important ansatz in the RISM formalism is the division of all atoms (or,

more generically, “sites”) into the categories of solute and solvent. This produces three kinds of distribution and correlation functions (essentially three cases of Eqn. 2.9) as described above: solute-solute, solute-solvent, and solvent-solvent. If the solute is rigidly held, then the solution for solute-solute correlations is trivial; there is no change and thus no correlation beyond the fixed coordinates. This assumption/approximation can be expanded upon, but is not immediately important nor unique to the RISM formalism and so will not be discussed further. Next, the highly interesting solute-solvent correlation functions (*e.g.* the interactions of an enzyme with water) can be solved, but only if the solvent-solvent correlation functions are already known (and *vice versa*). This is the same kind of impasse that was reached above that necessitated the closure relation and the solution introduced by RISM is its third and final defining aspect. The core approximation here is that solvent-solvent correlations can be uniformly orientationally averaged. That is, it is assumed that there is no orientational preference for interactions between solvent sites. This is exactly true for spherically symmetric sites (*e.g.* monoatomic ions), but is clearly a potentially severe approximation for larger solvents (*e.g.* water molecules). Unfortunately, such an approximation is unavoidable (at least without offering another one) and so it will not be dwelled upon. Going forward, Eqn. 2.9, which is clearly three-dimensional (and so it describes 3D-RISM), can now be written in one-dimension (*i.e.* the radial separation, and so it describes 1D-RISM). Such an equation can be solved given certain physical parameters (*e.g.* the pure solvent dielectric constant as in DRISM[46]) and used to produce a *solvent susceptibility* for use in a 3D-RISM calculation describing solute-solvent interactions.

Using the above, a practical workflow for RISM calculations can be developed providing a direct route to essentially any statistical mechanical relation. Once a molecular model has been chosen, the only assumptions/approximations needed are those described above. First, a 1D-RISM calculation is performed in order to establish the solvent-solvent correlations and the statistical mechanical nature of the pure solvent. Then, this information is combined with a three-dimensional structure of a solute to perform a 3D-RISM calculation that produces solute-solvent correlations. From here various quantities such as the free energy, local solvent densities, and other thermodynamic parameters can be obtained. Of course, computer software is generally needed to perform these calculations and it is generally non-trivial to do this in a fashion that

produces numerically-stable and robust results. Fortunately, much of this has already been done in the last few decades and will not be discussed further other than to state that it exists and enjoys widespread use[47].

2.3 Molecular Simulations - Statistics and Sampling

The general goal of molecular simulations is to take a model, $u(\mathbf{x})$, and explore what configurations, \mathbf{x} , arise and with what probabilities (see Section 2.2). A common strategy, first suggested in the literature by Metropolis, *et al.*[48], is to devise an algorithm that generates sample configurations with the correct probability weights. This has the distinct advantage that simple averages of a sample set will automatically correspond to ensemble averages. In the general case of a microscopic observable $A(\mathbf{x})$, whose specific value depends on the particular configuration, the ensemble average can be estimated from a statistical sample (of size N , say) from a simulation.

$$\langle A \rangle \approx \bar{A} = \frac{1}{N} \sum_{n=1}^N A(\mathbf{x}_n) \quad (2.10)$$

For example, $\langle A \rangle$ may represent the density of a system of water molecules at a given temperature and pressure. $A(\mathbf{x})$ could then be the volume of space that a number of water molecules subsume at a fixed point in time. Importantly, this microscopic density is not always the same, although the macroscopic density is. The two are only the same in an average, statistical sense. A key strength of molecular simulations is that they provide a nuanced description of these kinds of relationships.

Molecular Dynamics

The core of molecular dynamics (MD) is the connection of dynamics (*i.e.* how things move in time) and statistical mechanics (*i.e.* the probabilistic connection between the microscopic and macroscopic, see Section 2.2) by the *ergodic hypothesis*. Simply put, for very long times, the time evolution of a system will be such that it visits all possible configurations with the same probabilities dictated by the Boltzmann relation (Eqn. 2.1). The problem of statistical sampling outlined above is thus reduced to computation of

the dynamics of a system. This is enormously convenient, since the system dynamics are trivially given by Newton’s equations of motion and are greatly amenable to simple mathematical procedures that can be rapidly implemented in a computer program. It is worth noting that in practice this approach, at least naïvely, is not what one usually wants and the algorithm must be adapted to account for temperature and pressure for example. This can be done by the use of so-called thermostat and/or barostat algorithms[49, 50, 51, 52] as well as recourse to phenomenological force laws such as Langevin dynamics[53, 54]. Other modifications can be used to address changes in the environment such as pH[55, 56, 57, 58, 59, 60, 61].

2.3.1 Multistate Sampling

Up to this point simulations have been discussed in the context of generating a single statistical sample from a given molecular model. The idea of multistate sampling extends this concept in two ways. First, it provides a framework for using a simulation to generate statistical information about some *other* molecular model. As will be seen below, this is immensely useful in solving difficult sampling problems, as it allows for the modification of good physical models with poor sampling characteristics to become unphysical models with good sampling characteristics. Second, multistate sampling allows for multiple simulations with multiple models (which may or may not be good physical models) to be combined in a way so as to increase their statistical information. Both of these techniques sound incredibly powerful and have a small tinge of “something for nothing.” As would be expected, there are caveats. Nonetheless, the theoretical and statistical underpinnings are sound and actually quite simple.

Exact Relations and Historical Origins

The most basic multistate sampling expression is often attributed to Zwanzig (hence it will be referred to as the Zwanzig relation), although it was probably somewhat known before his 1954 paper[62].

$$f_i - f_j = -\ln \langle e^{-[u_i(\mathbf{x}) - u_j(\mathbf{x})]} \rangle_j = \ln \langle e^{[u_i(\mathbf{x}) - u_j(\mathbf{x})]} \rangle_i \quad (2.11)$$

Here an ensemble average for system i is indicated as $\langle \cdot \rangle_i$. The main result of this relation (which is fairly simple to derive and can be found in several modern textbooks[63, 64]) is that the free energy between two systems (labelled i and j , say) can be related to an ensemble average for *either* system, provided that the reduced potential for both states ($u_i(\mathbf{x})$ and $u_j(\mathbf{x})$) is known and can be applied to configurations from both systems. Furthermore, once this relative free energy is calculated, *any* arbitrary observable can be calculated in either system using a similar relation (often attributed to Torrie and Valleau[65]).

$$\begin{aligned} \langle A \rangle_i &= \frac{\langle A(\mathbf{x}) e^{-[u_i(\mathbf{x}) - u_j(\mathbf{x})]} \rangle_j}{\langle e^{-[u_i(\mathbf{x}) - u_j(\mathbf{x})]} \rangle_j} \\ &= e^{f_i - f_j} \langle A(\mathbf{x}) e^{-[u_i(\mathbf{x}) - u_j(\mathbf{x})]} \rangle_j \end{aligned} \quad (2.12)$$

Note that in the last line an explicit connection to Eqn. 2.11 has been made. Converting these equations to naïve sample estimators (Eqn. 2.10) gives a *weighted* average instead of a simple average (*i.e.* each sample receives a different weight that is probably not equal to $1/N$),

$$\begin{aligned} \hat{f}_i - \hat{f}_j &= -\ln \frac{1}{N_j} \sum_n^{N_j} e^{-[u_i(\mathbf{x}_{jn}) - u_j(\mathbf{x}_{jn})]} = -\ln \sum_n^{N_j} \frac{1}{N_j e^{[u_i(\mathbf{x}_{jn}) - u_j(\mathbf{x}_{jn})]}} \\ \hat{f}_i &= -\ln \sum_n^{N_j} \frac{1}{N_j e^{\hat{f}_j - [u_j(\mathbf{x}_{jn}) - u_i(\mathbf{x}_{jn})]}} \equiv -\ln \sum_n^{N_j} w_i(\mathbf{x}_{jn}) \end{aligned} \quad (2.13)$$

Here we have used $\hat{\cdot}$ to indicate that these are statistical estimates for f and have moved the ensemble average subscript, j , to label the sample configurations (so that \mathbf{x}_{jn} is sample configuration n from system j). Of course an equivalent estimator can be obtained by inverting i and j and taking samples from a simulation in system i . The reason for defining the sample weights $w_i(\mathbf{x}_{jn})$ is likely much more obvious in the context of observables (the manipulations are similar to the above).

$$\hat{A}_i = e^{\hat{f}_i} \sum_n^N w_i(\mathbf{x}_{jn}) A(\mathbf{x}_{jn}) \quad (2.14)$$

This expression has the form of a simple average if $f_i = 0$ and $w_i(\mathbf{x}_{jn}) = 1/N_j$ for each configuration. As should be expected, this is exactly the case when the samples are

drawn from system i instead of system j (there is no free energy difference between the systems and the samples are already generated with the correct probability weights).

Modern Developments and Improved Estimators

Unfortunately, in practice, the estimators described in the previous section are frequently not very accurate unless systems i and j are very similar; this has widely been known since the early work of Widom, for example[66]. In later work, Bennett showed that, for two systems, an optimal estimator (in the sense of having the lowest possible statistical error) equivalent to the Zwanzig relation can be obtained by generating samples for *both* systems and then combining the samples with a self-consistent estimate for the relative free energy[67]. This is done by modifying the sample weights used in Eqns. 2.13 and 2.14 to include information from both systems.

$$\begin{aligned}
 w_i(\mathbf{x}_n) &= \frac{1}{N_i e^{\hat{f}_i - [u_i(\mathbf{x}_n) - u_i(\mathbf{x}_n)]} + N_j e^{\hat{f}_j - [u_j(\mathbf{x}_n) - u_i(\mathbf{x}_n)]}} \\
 &= \frac{1}{N_i e^{\hat{f}_i} + N_j e^{\hat{f}_j - [u_j(\mathbf{x}_n) - u_i(\mathbf{x}_n)]}} \tag{2.15} \\
 \hat{f}_i &= -\ln \left[\sum_n^{N_i} w_i(\mathbf{x}_{in}) + \sum_n^{N_j} w_i(\mathbf{x}_{jn}) \right]
 \end{aligned}$$

The key difference in Eqn. 2.15 compared to those above is that \hat{f}_i appears on both sides of the equation and thus it must be solved for numerically in a self-consistent fashion. Despite the discovery of Eqn. 2.15 several decades ago, its utility was not fully appreciated into many years later. In particular, highly similar or identical generalizations to more than two states were repeatedly derived by various means[68, 69, 70, 71, 72]. The earliest of these often invoked a histogram approximation with various justifications after the fact[68, 70, 73, 74]. As such the method was called the weighted histogram analysis method or WHAM. Shirts, *et al.* re-discovered the analogies between WHAM and Bennett's method (often called the Bennett acceptance ratio or BAR)[75, 76] and so the revised formalism (including error estimators and extension to observables) was renamed the multistate Bennett acceptance ratio or MBAR[71]. In expressions highly analagous to Eqns. 2.14 and 2.15, K different systems are considered and the sample

weights are modified yet again. The self-consistent solution is thus for K different free energies and expressions for arbitrary observables follow immediately from them.

$$\begin{aligned}
 w_i(\mathbf{x}_n) &= \left[\sum_k^K N_k e^{\hat{f}_k - [u_k(\mathbf{x}_n) - u_i(\mathbf{x}_n)]} \right]^{-1} \\
 \hat{f}_i &= -\ln \sum_n^N w_i(\mathbf{x}_n) \\
 \hat{A}_i &= e^{\hat{f}_i} \sum_n^N w_i(\mathbf{x}_n) A(\mathbf{x}_n)
 \end{aligned} \tag{2.16}$$

Note that $N \equiv \sum_k^K N_k$ and summation over n now represents a summation over many (*i.e.* K) simulations.

Up to this point, it has been left somewhat vague as to what all of the different systems being simulated might be and how one might find themselves in the situation of having $K \gg 2$ simulations needing to be analyzed. Of course, the development of MBAR occurred right along these efforts and was frequently re-derived in order to address specific situations various researchers found themselves in. In the sections that follow some more specific examples of these will be given and it will be shown how they can be applied to a broad class of problems. In the remaining sections, specific problem applications in the study of catalytic reactions and thermodynamics will be presented and discussed.

2.3.2 Enhanced Sampling

Conventional molecular simulation techniques (such as MD) are often inefficient or inaccurate in solving specific problems. This can be for many reasons such as large kinetic barriers, phase transitions, or simply that the simulation generates a large amount of data that is not particularly relevant to the specific question that is meant to be asked. In these cases a general and powerful strategy is to focus or enhance the sampling in certain specific areas *a priori*. The wide array of such methods is almost too long to list and will not be attempted here. However, two particular, but quite general, methods will be discussed here.

Umbrella Sampling

The term “umbrella sampling” was first used in the literature by Torrie and Valleau[65]. However its use has changed somewhat over the years and is generally used to designate a number of strategies used together[77, 78]. The first of these methods is biased sampling.⁸ That is, a specific modification to the desired molecular model is made. This is usually done in the form of an additional, but non-physical, energy term. Clearly it is helpful (or even necessary) to know beforehand what coordinate(s) or quantity needs to be biased. This may be quite obvious if one wants to obtain information about high energy configurations relative to low energy ones, such as in the case of chemical reactions with large kinetic barriers. In this case an effective tactic is to identify a geometric coordinate aligned with that reaction (such as the length of chemical bonds that will break or form) and apply a localizing bias or restraint (Figure 2.8). That is, an additional term is added to force the system to sample configurations near a specific point of interest. A simple form for this kind of bias is a harmonic restraint, a virtual spring or rubber band holding various atoms in place. For an arbitrary coordinate, ξ , that tracks with the reaction progress (hence it is often called a “progress coordinate”) this takes the form

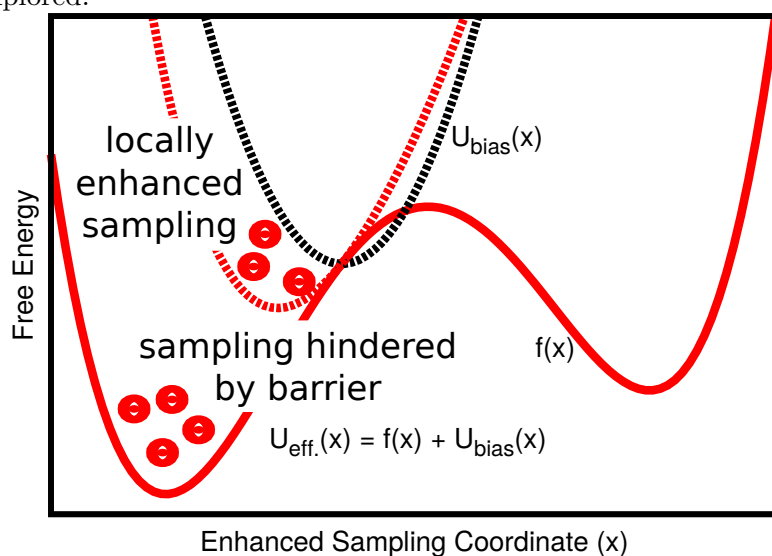
$$U_{\text{bias}}(\xi) = k(\xi - \xi_0)^2, \quad (2.17)$$

where $U_{\text{bias}}(\xi)$ is the bias energy, ξ_0 is the bias position (*i.e.* where in space sampling will approximately be localized), and k is the force constant (*i.e.* a measure of how strong the bias will be). It is also possible to employ many other more generic bias forms and on more complicated functions of the various coordinates. This is, for example, the approach used in accelerated MD[79]. In any event, the effect of the bias, which will otherwise cause the simulation results to be unphysical, can be removed by the techniques described in Section 2.3.1, although many other specific solutions to this problem have also been suggested[80, 81, 82, 83, 84, 85, 86], including, as will be discussed in Section 3.1.2, some recently published techniques presented here [4, 5]. A key drawback to this approach is that sampling may become *too* localized, in which case not enough data outside the localized region is obtained. This problem can be addressed by

⁸ Strictly speaking, this is the specific technique espoused by Torrie and Valleau in the original umbrella sampling paper.

stratifying or multi-staging the coordinate space by using multiple biases[87]. Again, after running multiple simulations with these biases, the results can be combined using multistate techniques. Thus, the two components of umbrella sampling, at least in modern parlance, have been described: 1.) biased sampling, usually with a harmonic localizing potential and 2.) stratified sampling using multiple simulations to span the sample space of interest.

Figure 2.8: The main concept behind umbrella sampling[65] is to augment the underlying free energy function, $f(x)$, of a specific coordinate whose sampling is impeded by the presence of large energy barriers. For example, it is straightforward to apply a harmonic bias, $U_{\text{bias}}(x)$, in order to localize sampling in a pre-specified region. The modified effective sampling potential, $U_{\text{eff.}}(x)$, then produces samples in a region not otherwise explored.



Replica Exchange

Replica exchange simulations, especially replica exchange molecular dynamics (REMD) variants, have become an extremely popular tool for improving the accuracy and efficiency of molecular simulations[88, 89, 90, 91, 92, 93, 94, 95, 96, 97, 98, 99, 100, 101, 102]. The seed of the idea dates as far back as the work of Bennett[67], but the specific algorithmic and implementational details for modern REMD were not formalized until the work of Sugita and Okamoto[103, 104]. The core idea of REMD builds off of earlier work

in *expanded ensembles*[105], which suggested the idea of running a single simulation of a system that changes between multiple sets of physical parameters, often referred to as *thermodynamic states*, or just *states*. As above, the resulting (unphysical) trajectory can be analyzed by multistate techniques. In theoretical terms, the transitions between such state are related to their relative free energies and thus an efficient simulation (*i.e.* one which efficiently samples all of the different states) requires at least a rough estimate of these free energies *a priori*.⁹ Unfortunately these can be difficult or tedious to determine in practice. The key development leading to replica exchange was that the statistical rules governing the *exchange* of states between two simulations do not require the relative free energies to be known[74, 103, 104] and that simulations performed in this way would produce virtually identical results within particular limits. Thus, by adding the complexity of running two or more simulations at the same time (a problem of significantly reduced difficulty on modern computing environments), the complexity of the state change move is reduced by making it an exchange move. More recent work has shown that such swap procedures are themselves special cases of permutation searching algorithms and have generalized the procedure to large groups of simulations[106, 107].

So far, the *exchange* part of REMD has been made apparent, but not the *replica* part. Taking the expanded ensemble view, each simulation can be seen as independent of the others, only being aware of the other simulations when an exchange move is made. Since all of the simulations have identical possible states to be in (they can always exchange to obtain one they have not yet occupied), they are formally identical in the long time limit (much like the ergodic hypothesis described in Section 2.3). In this way all of the simulations are copies (or replicas) of each other and have not only identical statistical properties but also identical dynamics (although this may only be in a sampling, rather than physical, sense)[74, 108].

⁹ It is somewhat astonishing that this insight was in fact recognized by Bennet more than 20 years before REMD was formally proposed[67]. However, he specifically did not advocate for actually attempting such an exotic simulation, although it is not exactly clear why. Perhaps his reluctance may have been due to the limits and scarcity of computing hardware.

Chapter 3

Methods Development for Molecular Simulation

Two main factors impact the quality of a molecular simulation, the accuracy of the model in describing the physical phenomenon of interest and the precision with which the simulation gathers information within the context of that model. The first of these is clearly important and a considerable amount of work has already been invested in developing models for the study of RNA and RNA catalysis including, but not limited to, molecular mechanical force field models[109, 110, 111, 112, 113, 114, 115], fast semiempirical quantum models[116, 117], long-range electrostatic treatments[118, 119, 120, 121], and linear-scaling electronic structure methods[122]. As such, this dissertation only presents tangential work in this are (specifically, Ref. 122). However the latter problem, that of producing precise statistical information, still represents a considerable challenge in the study of biological systems and RNA and catalytic RNA mechanisms in particular. Here we divide this problem into two categories and present two sections describing specific contributions two both of them.

In Section 3.1, the problem of optimal statistical estimation is considered. Again, considerable effort has already been invested in this problem and was covered in Section 2.3.1. Here the focus is on the optimal estimation of free energies and free energy landscapes (see Section 2.2.1) and the specific contributions that have been made in three recent publications[1, 4, 5]. In Ref. 1 a straightforward extension of the MBAR method

is described that incorporates modern developments in the field of non-parametric density estimation. These advances are reviewed and summarized in Section 3.1.1. In Refs. 4 and 5 a new free energy method based on maximum likelihood estimation is presented. The theory behind this work is reviewed in Sec. 3.1.2. Finally, in Section 3.1.3 the results of benchmarking studies comparing several methods, both old and new, are summarized using several results from Refs. 4 and 5.

3.1 Estimators for Free Energy Landscapes

In Section 2.3 molecular simulations were discussed within the context of statistical sampling. Several expressions, of various complexity, were presented that can be used to estimate arbitrary observables. However, specific forms or examples were not discussed. This section will address this omission, particularly by making connections with the quantities needed to estimate free energy landscapes described in Section 2.2.1. In particular, the problem of estimating marginal distributions will be considered. This is not the only approach to estimating free energy landscapes (see, for example, Refs. 86, 123, and 124), but it is the oldest and most widely used. Furthermore, the theoretical framework is quite flexible and accounts for a wide range of seemingly disparate estimators.

In what follows, marginal distribution estimators will be divided into two broad classes, non-parametric and parametric[125]. This is a non-unique distinction and others are possible. However, as will be seen, it draws a clear line between more traditional estimation methods and newer developments that will be presented here.

3.1.1 Non-Parametric Estimators

Kernel Density Estimation

Non-parametric density estimation methods make no assumptions as to the form of the underlying probability density. The main advantage of this is that the error of the estimate is, in the limit of large sample sizes, attributable entirely to statistical error[125, 126, 127]. That is, the systematic error vanishes. Such estimators are often known as being *asymptotically unbiased*, since they are guaranteed to approach the correct answer. A broad class of such estimators are known in the statistics literature

as *kernel density estimators* (KDEs) and are strikingly similar to Eqn. 2.5:

$$\rho(x) = \langle \delta(x - x') \rangle = \lim_{h \rightarrow 0} \left\langle \frac{1}{h} K \left(\frac{x - x'}{h} \right) \right\rangle \quad (3.1)$$

It should be noted that here we have not assumed that the configuration, x (written in one dimension for simplicity), is dimensionless (it has the same units as h) and so $g(\mathbf{x})$ has been replaced with $\rho(\mathbf{x})$ (which has units of $1/h$). The function K is often referred to as a *kernel* and the parameter h as the *bandwidth*. The kernel can have the form of any function that approaches a delta function in the limit that h is zero[37]. In practice, when finite samples (of size N , say) are used, this relation also requires that Nh approaches infinity[125, 126, 127, 128, 129].

For most readers, the most intuitive case of Eqn. 4.2 is likely that of a histogram, in which the space of all samples is divided into bins; an additional parameter for K is thus needed to define the location of these bins. The statement is now that this kind of estimator is exact in the limit that a large number of samples are sorted into the bins, which are in turn made ever smaller. Such an estimator is intuitive and simple to run on a computer (especially when the bins are evenly sized), but suffers from the fact that the resulting density is on a fixed, discrete grid. That is, the choice of bin positions affects the accuracy in a way that does not obviously vanish and the result is not rigorously smooth and differentiable, as the exact density is expected to be. This problem can be solved by replacing the histogram bins themselves with smooth, differentiable functions, such as those in Table 3.1. In multiple dimensions the one-dimensional argument $\frac{\mathbf{x}-\mathbf{x}'}{h}$ can be generalized by introducing a (real symmetric) bandwidth matrix $\mathbf{H}^{\frac{1}{2}}$ and instead using $\mathbf{H}^{-\frac{1}{2}}\mathbf{y}$, with $\mathbf{y} \equiv \mathbf{x} - \mathbf{x}'$.

Table 3.1: Examples of kernel density estimator forms generalized to n dimensions. It is convenient to work in the scalar quantity $u \equiv \mathbf{y}^T \mathbf{H}^{-1} \mathbf{y}$, where $\mathbf{y} \equiv \mathbf{x} - \mathbf{x}'$. This form translates the kernel position to the origin and normalizes the domain to the n -dimensional unit sphere via the bandwidth matrix, $\mathbf{H}^{\frac{1}{2}}$. Note that u is quadratic in \mathbf{x} and that normalization is obtained by dividing the kernel by its integral over all of space. In the table below, $\mathbf{1}_{\{\dots\}}$ is an indicator function, $V_n = \frac{n\pi^{n/2}}{\Gamma(\frac{n}{2}+1)}$ is the volume of an n -dimensional unit sphere, $\Gamma(n)$ is the gamma function ($\Gamma(n+1) = n!$ for $n \in \mathbb{Z}, n \geq 0$), and $B_{mk} \equiv \frac{m!}{k!(m-k)!}$ is a binomial coefficient.

name	kernel function, $K(u)$	$\int d\mathbf{x}K(u)$	1D plot
order m	$(1-u)^m \mathbf{1}_{\{ u \leq 1\}}$	$nV_n \sum_{k=0}^m \frac{(-1)^k B_{mk}}{2k+n}$	
Epanechnikov ($m=1$)	$(1-u) \mathbf{1}_{\{ u \leq 1\}}$	$\frac{2V_n}{n+2}$	
biweight ($m=2$)	$(1-u)^2 \mathbf{1}_{\{ u \leq 1\}}$	$\frac{8V_n}{(n+2)(n+4)}$	
triweight ($m=3$)	$(1-u)^3 \mathbf{1}_{\{ u \leq 1\}}$	$\frac{48V_n}{(n+2)(n+4)(n+6)}$	
Gaussian	$e^{-\frac{u}{2}}$	$(2\pi)^{\frac{n}{2}}$	

Application with the Multistate Bennett Acceptance Ratio

In the context of estimating free energy landscapes (Eqn. 2.6), it is convenient to use a multistate sampling framework, such as the MBAR method (see Section 2.3.1). However, in this case the negative logarithm of the observable is desired and this observable is now a function, potentially in multiple dimensions. A working estimator for the free energy landscape can be made by inserting Eqn. 4.2 as the observable:

$$\hat{f}(\mathbf{x}) = \hat{f} - \ln \sum_n^N w(\mathbf{x}_n) K(\mathbf{x}; \mathbf{x}_n, \mathbf{H}) \quad (3.2)$$

$$\mathbf{y} \equiv \mathbf{x} - \mathbf{x}_n \quad K(\mathbf{x}; \mathbf{x}_n, \mathbf{H}) = \frac{1}{|\mathbf{H}|} K(\mathbf{y}^T \mathbf{H}^{-1} \mathbf{y})$$

Depending on the context, \hat{f} may be zero and usually does not matter. In any event, \hat{f} (or any constant shift to $\hat{f}(\mathbf{x})$) will cancel when taking differences for any two values of \mathbf{x} . It should also be noted that Eqn. 3.2 has been written in the form of a *multi-dimensional* KDE, meaning it can be used for a free energy landscape or arbitrary dimension. In general, it is also possible to choose the bandwidth matrix, $\mathbf{H}^{\frac{1}{2}}$, to depend on n , but this was not done here (see Ref. 1 for a presentation of the equations where this is the case).

3.1.2 Parametric Estimators

The main principle behind parametric density estimation is that one frequently possesses (or hopes to possess) prior knowledge about the nature of the exact distribution. Properly leveraging this information against statistical data can then, in principle, improve the accuracy of the result. However, if the prior information is incorrect, imprecise, or simply not well-reflected by the data, then such an approach can lead to a biased result that does not look like the correct distribution. Moreover, this kind of bias will often not vanish as the amount of data increases (*i.e.* the method is *asymptotically biased*). Nonetheless, such a method may still be useful and preferable to an unbiased, non-parametric method if the bias error is comparable to the statistical error. Such a method is often called *weakly biased*.

Maximum Likelihood Estimation

A common and quite successful approach to parametric density estimation is the method of *maximum likelihood* or maximum likelihood estimation (MLE)[130, 131]. The basic idea behind MLE is to construct a model density for a sample data set, \mathbf{x}_n , composed of N samples in terms of a set of M parameters, $\boldsymbol{\theta}_m$. If each data point is independently and identically distributed (a non-trivial, but necessary assumption in much of molecular simulation), then this distribution can be written as a product distributions:

$$\rho(\mathbf{x}_n; \boldsymbol{\theta}_m) = \prod_n^N \rho(x_n; \boldsymbol{\theta}_m) \tag{3.3}$$

Note that on the right side each term depends on only *one* sample, but *all* of the model parameters. This expression is then defined, in the reverse sense, as the *likelihood function*, $\mathcal{L}(\boldsymbol{\theta}_m; \mathbf{x}_n)$:

$$\mathcal{L}(\boldsymbol{\theta}_m; \mathbf{x}_n) = \prod_n^N \rho(\boldsymbol{\theta}_m; x_n) \tag{3.4}$$

The goal of MLE methods is to maximize the likelihood with respect to the model parameters, $\boldsymbol{\theta}_m$, in terms of the observed data, \mathbf{x}_n . Since $\mathcal{L}(\boldsymbol{\theta}_m; \mathbf{x}_n)$ is always positive (probability densities cannot be negative) and the logarithmic function is monotonic, it is often convenient to re-express Eqn. 3.4 by taking the logarithm of both sides; such a function has the same extrema as the likelihood up to an arbitrary constant.

$$\hat{l}(\boldsymbol{\theta}_m; \mathbf{x}_n) \equiv \frac{1}{N} \sum_n^N \ln \rho(\boldsymbol{\theta}_m; x_n) \tag{3.5}$$

The function $\hat{l}(\boldsymbol{\theta}_m; \mathbf{x}_n)$, although technically a *log-likelihood* (or in this case, the *average log-likelihood* due to the multiplicative constant being chosen as $1/N$), can thus be used interchangeably with the likelihood in many instances. Again, the objective is to maximize $\hat{l}(\boldsymbol{\theta}_m; \mathbf{x}_n)$ by varying the parameters, $\boldsymbol{\theta}_m$. The optimal set of parameters can then be used to generate an optimal density estimate, under the assumption that the model is accurate. Clearly even a “good” parameter set for a poor model will still give rise to a poor density estimate.

The Variational Free Energy Profile Method

The premise behind the variational free energy profile[4, 5] (vFEP) method is to apply the MLE formalism above to the density estimation problem for free energy landscapes (Section 2.2.1) within the context of umbrella sampling simulations (Section 2.3.2). In this approach, an arbitrary (reduced) bias potential, $u_b(\mathbf{x}_1)$, is added to the reduced potential for the physical model, $u(\mathbf{x})$, in order to enhance sampling. This is usually done many times with many simulations with different *localizing* bias potentials (possibly within a replica exchange framework as in Section 2.3.2) until the sample space is well covered. In each case a *biased* free energy landscape, $f_b(\mathbf{x}_1)$, can be obtained from calculation of the biased marginal distribution, $g_b(\mathbf{x}_1)$. See Eqns. 2.2, 2.4, 2.5 and 2.6 for details.

$$\begin{aligned}
 e^{-f_b} &\equiv \int d\mathbf{x} g_b(\mathbf{x}) \\
 g_b(\mathbf{x}_1) &= \int d\mathbf{x}'_1 \delta(\mathbf{x}_1 - \mathbf{x}'_1) \int d\mathbf{x}_2 \dots \int d\mathbf{x}_N g_b(\mathbf{x}) \\
 &= \int d\mathbf{x}'_1 \delta(\mathbf{x}_1 - \mathbf{x}'_1) \int d\mathbf{x}_2 \dots \int d\mathbf{x}_N e^{f_b - [u(\mathbf{x}) + u_b(\mathbf{x}'_1)]} \\
 &= e^{f_b - f - u_b(\mathbf{x}_1)} \int d\mathbf{x}'_1 \delta(\mathbf{x}_1 - \mathbf{x}'_1) \int d\mathbf{x}_2 \dots \int d\mathbf{x}_N e^{f - u(\mathbf{x})} \\
 &= e^{f_b - f - u_b(\mathbf{x}_1) - f(\mathbf{x}_1)} \equiv e^{\Delta f_b - [f(\mathbf{x}_1) + u_b(\mathbf{x}_1)]}
 \end{aligned}$$

It is important to note that for the vFEP approach it is assumed (as is commonly done) that the bias potential is only a function of a few coordinates (the configuration of just one atom, say) and that these are the coordinates for which the free energy landscape is desired.¹ This assumption (which can be easily imposed) is critical to the manipulation in the third line above, which removes the bias potential from the integral. In the fourth and fifth lines, it has been shown that any biased free energy landscape can be related to the unbiased free energy landscape via the bias potential and a constant free energy difference, $\Delta f_b \equiv f_b - f$, between the two states. Such an expression is highly related to those discussed in Section 2.3.1.

¹ It is also possible to use bias coordinates that are linear or non-linear combinations of the desired free energy landscape coordinates. Such a transformation should not affect the expressions here, although it may impact the necessary Jacobian factor(s) when relating to a PMF (see Section 2.2.1 and footnotes therein).

Given a set of many umbrella sampling simulations with many bias potentials, a likelihood function can be constructed using all of the observed data points (*i.e.* summing over all of the K different simulations). Here the specific expression for the biased distribution functions (now indexed by k instead of b) derived above is inserted into Eqn. 3.5.²

$$\begin{aligned}
 \hat{l}(\boldsymbol{\theta}_m; \mathbf{x}_{11n}, \mathbf{x}_{12n}, \dots, \mathbf{x}_{1Kn}) &= \sum_k^K \hat{l}(\boldsymbol{\theta}_m; \mathbf{x}_{1kn}) \\
 &= \sum_k^K \frac{1}{N_k} \sum_n^{N_k} \ln g_k(\boldsymbol{\theta}_m; \mathbf{x}_{1kn}) \\
 &= - \sum_k^K \frac{1}{N_k} \sum_n^{N_k} \left[\Delta f_k + \hat{f}(\boldsymbol{\theta}_m; \mathbf{x}_{1kn}) + u_k(\mathbf{x}_{1kn}) \right] \\
 &= - \sum_k^K \left[\Delta f_k + \frac{1}{N_k} \sum_n^{N_k} \hat{f}(\boldsymbol{\theta}_m; \mathbf{x}_{1kn}) + u_k(\mathbf{x}_{1kn}) \right]
 \end{aligned} \tag{3.6}$$

The most important result here is that the exact free energy landscape, $f(\mathbf{x}_1)$, can now be estimated by optimizing, in a maximum likelihood sense, the function, $\hat{l}(\boldsymbol{\theta}_m; \dots)$ with respect to an arbitrary parameter set, $\boldsymbol{\theta}_m$. The target function only depends on the observed values of \mathbf{x}_1 from each simulation (such that \mathbf{x}_{1kn} is the n th sample of \mathbf{x}_1 from simulation k). This parameter set then leads immediately to an optimal estimate for the free energy landscape, $\hat{f}(\mathbf{x}_1; \boldsymbol{\theta}_m)$ (the order of the arguments has now been reversed to emphasize that this is now the important estimate).

A broad class of parameter dependent free energy landscape estimators can be defined as *spline functions*. The numerical details of such functions are largely uninteresting for the present application and will not be discussed in detail. However, some brief comments will clarify the motivation behind their use and some of the details of the implementation. Most splines are essentially piecewise polynomials constructed such that they are smooth and differentiable and can be evaluated analytically on a fixed domain. The ability of a spline to represent *any* function is limited by the number and location of spline nodes (*i.e.* how many polynomials are used) and the spline order (*i.e.*

² It is also possible to construct Eqn. 3.6 as a weighted sum of likelihoods for each simulation. However, there is no obvious *a priori* way to weight the simulations and, in any event, the effect of such a weighting will disappear in the large sample limit (see Ref. 4).

how many parameters define each polynomial). This can be thought of as a general and practical approach to Taylor series expansions, whereby any function can be represented exactly as an infinite series on some domain. An arbitrarily accurate representation can thus be obtained by including more and more terms in a finite series. Splines extend this by also adding more and more expansion points.

In physical terms, the use of splines to represent free energy landscapes is equivalent to very weak assumptions on their behavior. First, it assumes that the landscape is smooth and differentiable. It is not clear that this *must* be true for any arbitrary coordinate choice and it is likely possible to construct model systems in which the free energy landscape along one or more coordinates is discontinuous. Nonetheless, it is hard to imagine a clear physical process (such as a chemical reaction) being well described by such a surface. Therefore, this assumption can be taken as a simple extension of the choice of free energy coordinates. If the coordinates are well chosen then the surface is expected to be well-behaved and the assumption is likely justified. The second assumption, at least in current vFEP implementations, is that the free energy landscape can be locally approximated by a low order polynomial (cubic splines have been primarily used thus far). Clearly this approximation can be made arbitrarily accurate if the spline nodes become vanishingly close together (as in a Taylor series expansion). However, the placement of spline nodes is necessarily linked to the sample data set, as the vFEP likelihood function can only be meaningfully and uniquely optimized where data has been observed. Thus, the approximation of simple local behavior can be viewed as an interpolation technique, applying a minimal expectation of smoothness and differentiability in regions where data is sparse or non-existent. Of course, this is the exact kind of bias that parametric methods are known to create and can lead to potential errors. However, the vFEP framework can naturally adapt to additional data and modify this kind of assumption, essentially by changing the model. Furthermore, in the alternative, non-parametric framework, there is very little to be said in the way of interpolation, even if a local polynomial approximation is reasonable. As will be seen in several examples, this aspect of vFEP can be extraordinarily useful when analyzing sparse data sets and often leads to a more rich data analysis that can be used to guide further data collection.

3.1.3 Comparison of Free Energy Landscape Estimation Methods

In order to evaluate the performance of the vFEP method, several realistic model applications were recently examined and presented in the literature[4, 5]. As will be seen, a broad range of free energy landscapes with different characteristic shapes were considered and comparisons were made to other methods, both conventional and more recently developed. The two most well known of these are the WHAM and MBAR methods described in Section 2.3.1. Although these have been shown to be formally equivalent[71, 72], early implementations of WHAM (such as that by Grossfield used here[132]) strictly used histogram based estimates as opposed to more robust expressions.³ In order to cover a broader range of estimator performance and alleviate ambiguity, in what follows “WHAM” will be used to indicate this early form along with a histogram estimator for the free energy landscape. Conversely, “MBAR” will refer to the newer implementations (particularly the pymbar implementation by Shirts, *et al.*[133])⁴ along with a Gaussian kernel density estimator as in Section 3.1.1 and Ref. 1. Another, more recent parametric method known as umbrella integration (UI) will also be considered. UI can be seen as a specific case of vFEP, albeit with considerable modification of Eqn. 3.6 and slightly more restrictive assumptions (*e.g.* the bias potentials are assumed to be harmonic). Although implementations of UI in higher dimensions have been reported[83], these are not readily or widely available and so such results will not be presented here. Similarly, other alternative methods omitted here are left out not because of lack of awareness, but simply because they are non-trivial to implement and not publicly distributed.

In the comparisons that follow, although the specific protocols that were followed may vary, the general simulation strategies are identical. Umbrella sampling simulations using harmonic biases (see Section 2.3.2) were carried out on the system, usually along a simple coordinate such as a dihedral angle, bond length (or distance), or an atom transfer coordinate (*i.e.* a distance of two bond lengths). The reference positions for these biases

³ Interestingly, this was not actually recommended in the original WHAM paper by Kumar, *et al.*[68]. Instead, the “correct” expression that was re-discovered as MBAR was presented, although without explicit mathematical justification.

⁴ The cited software is actually a more recent version than was used in much of the published work presented here (see Refs. 4 and 5 for details). Most of the code modifications in that time were for performance purposes and, although some slight variations are possible/expected, the results with the newer versions are expected to be essentially identical.

were always spaced on a regular grid which is convenient, but unnecessary. Similarly, the strength of the biases (*i.e.* the harmonic force constant) were often uniform across all umbrella simulations (frequently called “windows” in the literature), although this is also unnecessary for the methods described here. Each of these simulations were then prepared and run with various molecular dynamics protocols for the same period of time.

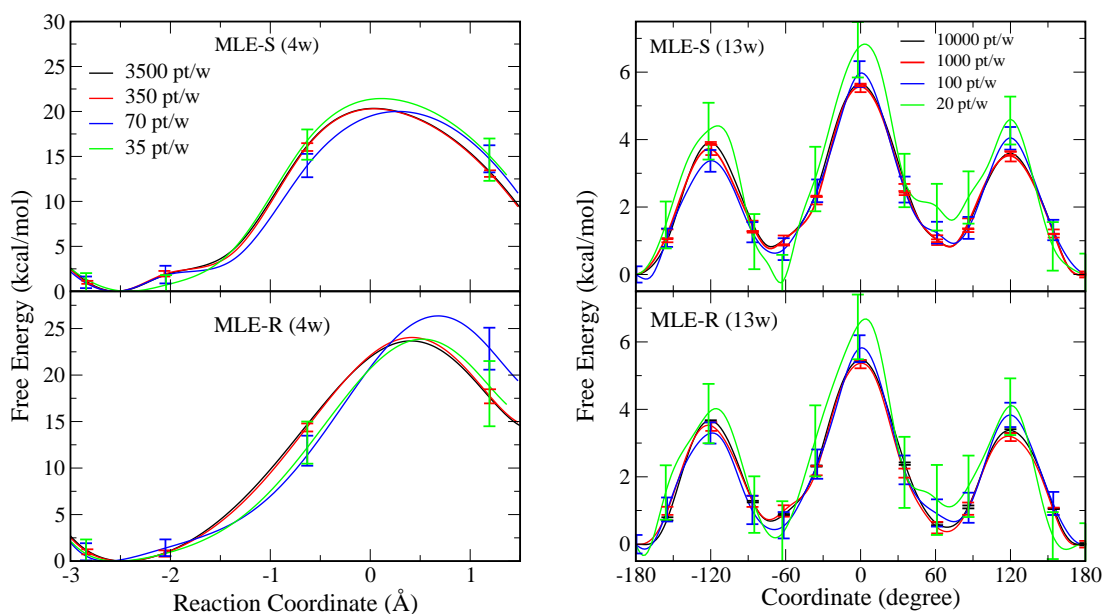
The two main aspects to be tested here are the effect of the size and sparsity of the data sets on performance. The first aspect describes *how much* data has been collected from each simulation and gives a measure of efficiency. That is, a method is considered to be efficient if it yields the same answer for a small data set as it does for a large data set under the assumption that more data will lead to a more accurate result. This is tested here by subsampling a full data set to create smaller data sets that, in principle, contain the same amount of information.⁵ The second aspect describes the effective length scale or *overlap* of the data. This has long been considered a critical point for multistate sampling methods[68, 70, 71, 72, 4, 5] and gives a measure of how well a method interpolates between regions in the free energy coordinate space. A method is considered to interpolate well if it yields the same answer for a sparse data set as it does for a dense data set. This is tested here by removing simulations to produce sparser and sparser grids of bias potentials and thus sparser data. Finally, for a related, but somewhat different, problem, the notion of data smoothing will be considered. This is primarily a concern for non-parametric methods, which often include an explicit smoothing degree of freedom. Parametric methods, on the other hand, often implicitly build smoothing into the statistical model. This will be gauged by the degree to which smoothing impacts the fidelity of the data analysis.

Data Convergence: Sample Size Performance

For the first examples of vFEP performance, some simple one-dimensional free energy landscapes will be considered. It is well known that standard approaches are capable of handling such data sets and so the focus here will primarily be on the performance of

⁵ Here data collection is always at regular intervals on the same period. This does require some assumptions regarding statistical correlations, but these are expected to affect all methods equally and will not be discussed further.

Figure 3.1: Performance of vFEP with different sample sizes when calculating the free energy profile of simple chemical reactions. vFEP is able to retain qualitative and quantitative profile features even with a small number of samples and this holds true for profiles of distinctly different shapes. The comparisons below show the free energy profiles of a simple chemical reaction (phosphoryl transfer of a simple RNA backbone model, left) and conformational transition (isomerization of butane, right). Reproduced with permission from Ref. 4. Copyright 2013 American Chemical Society.

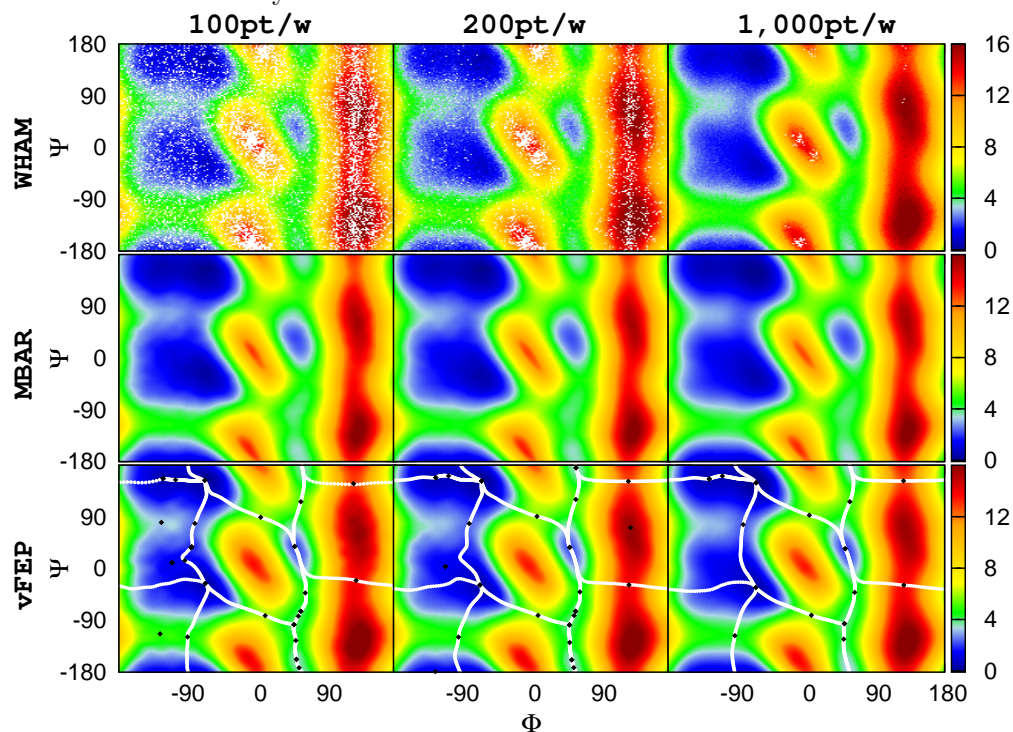


vFEP. Two different parametric functions were tested within the maximum likelihood framework of Eqn. 3.6, one based on spline interpolation (MLE-S) and one based on the rational interpolation function framework (MLE-R). Results with both methods are presented in Figure 3.1 for umbrella sampling simulations of a simple phosphoryl transfer reaction and the isomerization of butane, respectively. These two free energy profiles represent considerably different shapes. The former is effectively a double well potential with reactants, products, and a single transition state, while the latter is a periodic function with multiple extrema. It is reassuring, but rather uninteresting, that all of the methods surveyed produce statistically identical results for both examples when the full data set is used and, in general, produce results which gradually become identical to this result as the sample size increases (not all shown, see Ref. 4 and Figure 3.4 below for

details). More interestingly, vFEP retains this behavior (albeit quite differently for the MLE-S and MLE-R cases) even when many fewer simulations are used. Conventional non-parametric methods fail to even converge in these instances. Conversely, even in the extreme case of only dozens of data points (two orders of magnitude less than the full sample size), vFEP produces qualitatively converged results. This is not to say that these results are exactly correct, but rather that vFEP can be useful in obtaining coarse estimates from very little data when other methods would not offer any information. Moreover, estimates of the statistical error show that this convergence is fairly reliable and that the error bars decrease with the sample size and overlap with the results with large samples. As discussed above, since vFEP is a parametric method, it can be prone to bias and this is evident in Figure 3.1, where the two different parametric models converge to somewhat different results with different shapes and heights. Nonetheless, this bias decreases as more simulations are added (discussed below, see Figure 3.4).

As a more complicated example, the two-dimensional conformational landscape of alanine dipeptide is now considered. Without dwelling on specifics, such a compound is a model for the conformational space of the protein backbone and has been studied exhaustively via computational means, making it an excellent test for new simulation methods[134, 74, 135, 108]. As in the previous one-dimensional examples, all of the methods tested here identically converge within statistical error when the full data set is used (Figure 3.2). When a large number of simulations (or umbrella “windows”) are used, reducing the number of data points in the set does not strongly affect the results, except for when using WHAM, in which case the histogram estimator causes some loss of resolution. This is, however, regained when a smooth non-parametric estimator is used, as with MBAR. This overall consistency is similarly observed with vFEP and the improvements in smoothing are readily visible by analyzing the location of stationary points and zero gradient path points (*i.e.* the minimum free energy pathways, Figure 3.2, bottom row). As the number of data points increases, the statistical noise in the data diminishes and fewer false minima are observed, especially in the bottom right of the landscape and the basin at top left. More interestingly, when fewer simulations are used, this pattern is essentially only maintained by vFEP and the other methods fail to converge to the same answer (Figure 3.3. Although the performance of the other methods does improve with the sample size (most notably MBAR), this is still largely

Figure 3.2: Performance comparison of several free energy methods with different sample sizes on a dense sampling grid (576 simulations) when calculating the two-dimensional free energy landscape of alanine dipeptide. All three methods converge to essentially the same result. The rate of this convergence can be visually tracked by the disappearance of spurious stationary points (black dots) and minimum free energy paths (white lines) as the sample size increases. Reproduced with permission from Ref. 5. Copyright 2014 American Chemical Society.

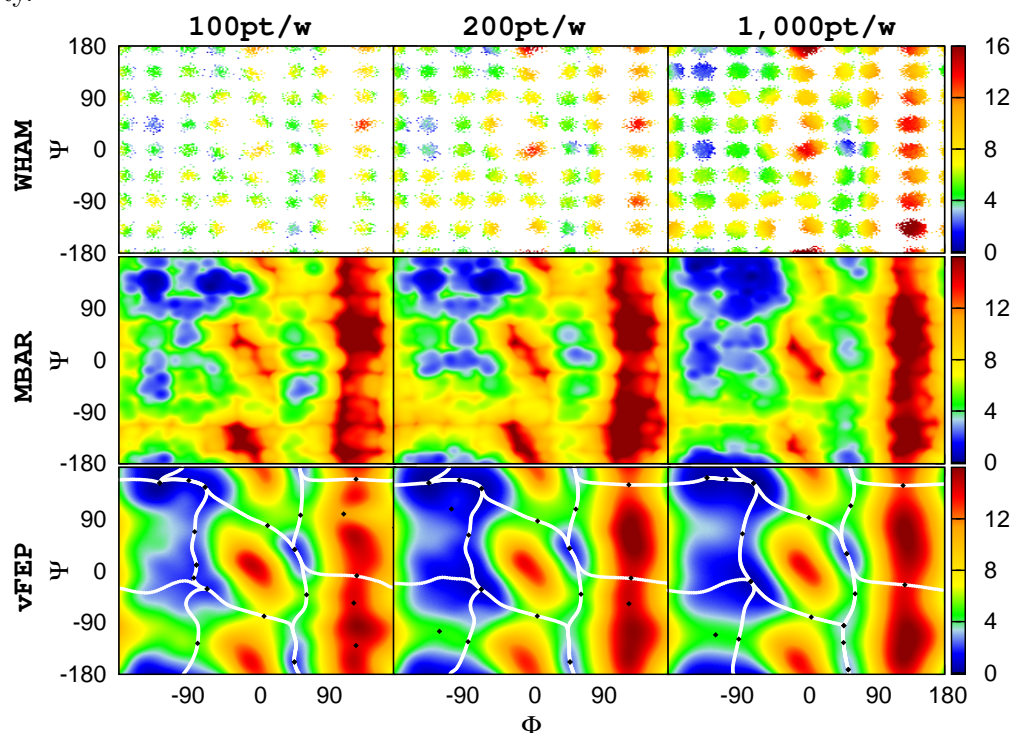


comparable to the performance of vFEP with a sample size an order of magnitude smaller.

Data Interpolation: Sample Sparsity Performance

In the examples of the previous section, it was seen that vFEP retains accuracy with smaller data sets as well or better than existing methods; this is especially true when data is taken from only a few simulations. The number of simulations is of particular importance when evaluating computational expense, as this is often the most significant factor in determining whether or not a simulation is feasible and/or worthwhile. In this

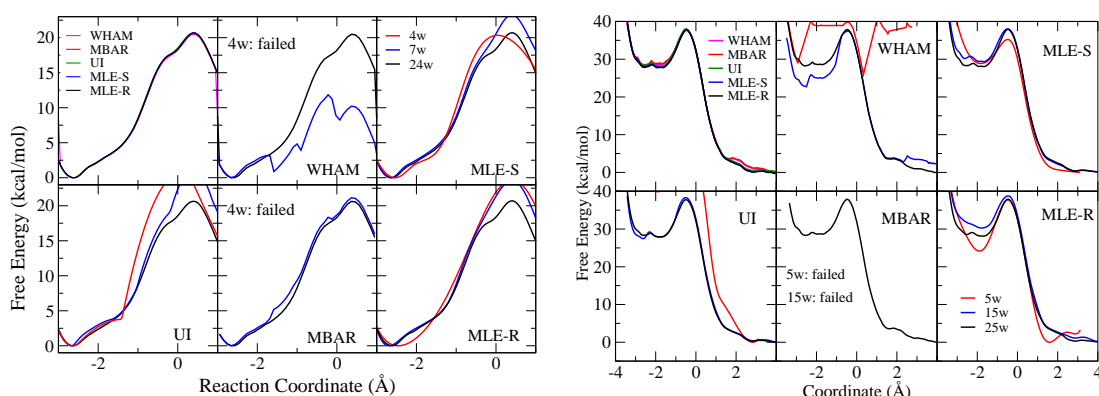
Figure 3.3: Performance comparison of several free energy methods with different sample sizes on a sparse sampling grid (64 simulations) when calculating the two-dimensional free energy landscape of alanine dipeptide. Only vFEP is successful in reconstructing a landscape similar to those in Figure 3.2. Furthermore, the rate of this convergence is impressive and can be visually tracked by the disappearance of spurious stationary points (black dots) and minimum free energy paths (white lines) as the sample size increases. Reproduced with permission from Ref. 5. Copyright 2014 American Chemical Society.



section many of the same simulations as above are revisited, but with the intention of determining a minimum threshold for the number of umbrella sampling simulations rather than sample size.

In Figure 3.4, free energy profiles for two distinct phosphoryl transfer reactions are presented. The first of these is identical to that discussed above and briefly analyzed in Figure 3.1. The first frame (both sets, top left plot) clearly shows that, when considering all of the simulations together, all of the free energy methods here yield statistically identical results. However, if the number of simulations is concertedly pared (from 24 to 7 and then 4 simulations or “windows”), a rapid degradation in quality is visible in

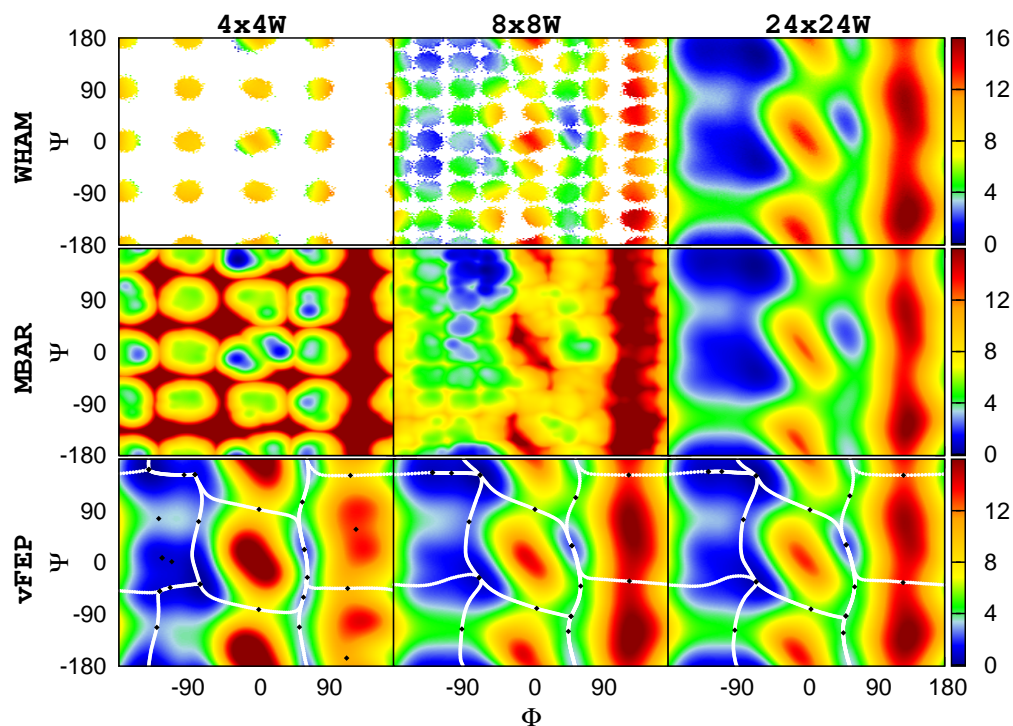
Figure 3.4: Performance comparison of several free energy methods with different sampling densities (*i.e.* number of “windows,” w) when calculating the free energy profile of simple phosphoryl transfer reactions. All of the methods ultimately converge to the same result (top left). However, at lower sampling densities (4w and 7w) only the two vFEP variants (MLE-S, MLE-R, and UI) produce useable results similar to that at the highest sampling density (24w). MBAR does show good performance, but is numerically unstable at extremely low sampling densities. Reproduced with permission from Ref. 4. Copyright 2013 American Chemical Society.



most methods, with the non-parametric estimators (*i.e.* WHAM and MBAR) failing to yield sensible results in the most extreme case. The parametric estimators, which bias the results towards a smooth, differentiable free energy profile fair considerably better in this limit. Most notably, both vFEP based methods (MLE-S and MLE-R, Figure 3.1, right plots) retain the same shape and nearly the same extrema in all cases. This is nearly true of umbrella integration as well (the only exception is the bottom left plot).

As in the previous section, all of the observed trends become more distinct when moving to higher dimensions. Figure 3.5 shows a similar comparison to those above for the same conformational free energy landscape of alanine dipeptide. Once again, with the full set of windows, all of the methods give essentially identical results. This indicates that there is little, if any, detectable bias or systematic error in the vFEP results. Conversely, at low numbers of windows the non-parametric methods show distinctly lacking results and provide little to no information in the regions where data is not collected. Indeed, this is an expected feature of such methods. On the other hand, vFEP, with only a few windows, gives results that are qualitatively comparable

Figure 3.5: Performance comparison of several free energy methods with different sampling densities when calculating the two-dimensional free energy landscape of alanine dipeptide. As in Figures 3.2 and 3.3, vFEP is most consistent across all sampling densities (*i.e.* number of “windows,” w). Furthermore, the rate of convergence is impressive and can be visually tracked by the disappearance of spurious stationary points (black dots) and minimum free energy paths (white lines) as the sample size increases. Reproduced with permission from Ref. 5. Copyright 2014 American Chemical Society.



to the full data set (Figure 3.5, bottom row). As above, the local shape of the free energy landscape can be helpfully visualized by examining the stationary points and zero gradient paths. Importantly, these are the aspects of the landscape that provide information regarding conformational transitions and are the result of principal interest in applications.

Data Smoothing: Non-Parametric Artifacts

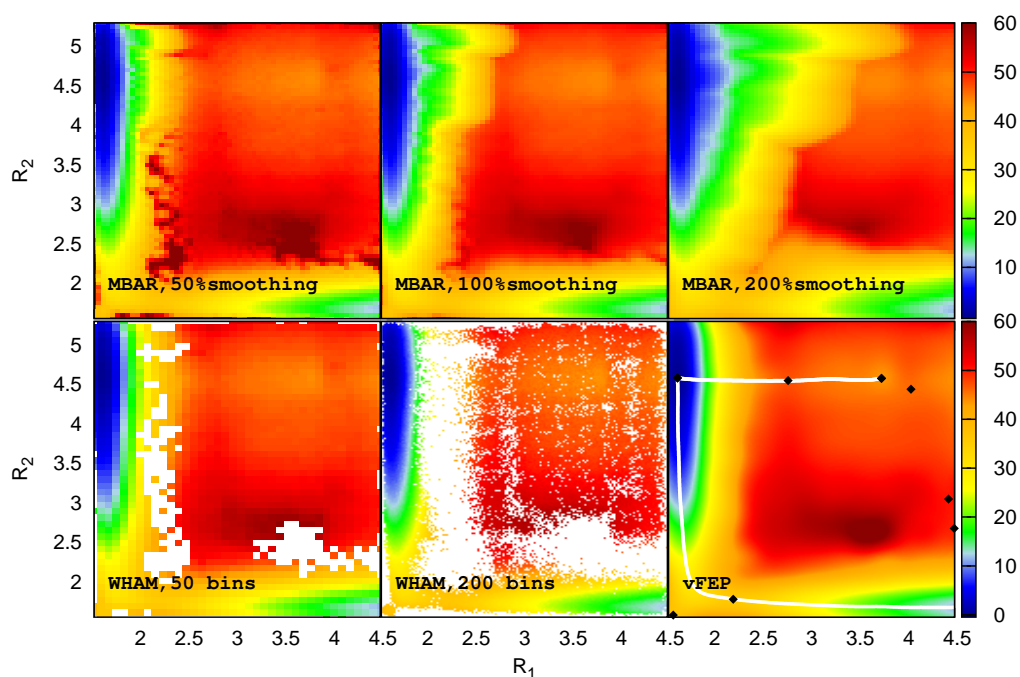
As a final comparison of non-parametric and parametric methods, we examine the degree to which data smoothing helps and hinders performance. In the previous sections

a fixed smoothing protocol was assumed using standard techniques[125, 126]. However, such procedures/algorithms are still a topic of active research and their effects are not always obvious. A simple probe of these effects is to make drastic changes in the smoothing parameters by scaling by a factor (one-half or two, say). In the case of histograms this is equivalent to a change in the bin size. Figure 3.6 shows several two-dimensional free energy landscapes for a phosphoryl transfer reaction using various scaling factors on top of standard smoothing procedures. As in previous examples, results with WHAM suffer from data sparsity, although the essential features (*e.g.* the basins and saddle points) appear to be mostly resolved. Using a kernel density estimator with MBAR provides some relief in these regions, but causes unusually shaped features in transitioning between the basins and high energy regions. Increasing the smoothing in a uniform fashion causes non-uniform changes in these features and apparently exacerbates the problem. Interestingly, vFEP, which implicitly optimizes smoothing within a parametric framework, does not suffer from these artifacts. Furthermore, at least to some degree, it resolves a number of stationary points in the high energy region and indicates the possible presence of an additional pathway. Again, this pathway information is one of the most important aspects of free energy landscape calculations in actual applications and the presence or absence of pathways is often a critical component of the purpose of the calculation.

3.2 Asynchronous Replica Exchange

In Section 2.3.2 the theoretical background and motivation behind replica exchange molecular dynamics (REMD) was introduced. However, little was said about specific exchange algorithms or practical implementations of this technique; this omission is addressed here. Early implementations of replica exchange consisted of multiple concurrent simulations running for a fixed period of simulation followed by a series of exchange attempts. These types of algorithms, commonly known as “nearest neighbor” exchanges, required all of the simulations to be sorted into pairs (*i.e.* each simulations with one of its “neighbors”). With an even number of simulations it was convenient to do this all at once, with all simulations paused while exchanges occurred. More recent

Figure 3.6: Evaluation of smoothing artifacts in non-parametric free energy methods versus vFEP for the two-dimensional free energy landscape of a phosphoryl transfer reaction. It is not always obvious as to how to identify or remove smoothing artifacts. The vFEP framework produces an inherently optimized smoothing choice based on the estimator type and the data and so provides a straightforward route to balancing out such artifacts. This is evidenced by the difficulty of non-parametric methods in resolving stationary points (black dots) and minimum free energy paths (white lines). Reproduced with permission from Ref. 5. Copyright 2014 American Chemical Society.



developments have shown this approach to be a specific case of a broader class of permutation searching algorithms[106, 107], some of which utilize global information about all of the replicas at once rather than as sets of pairs. Not only can this approach lead to more efficient sampling[106], it also highlights an interesting implementational detail of the exchange protocol. That is, the true objective of replica exchange is not to conduct exchanges, *per se*, but to sample different permutations of replicas amongst a set of thermodynamic states. Any algorithm (including pairwise exchanges) that accomplishes this is thus a valid scheme. A corollary to this statement is that the concurrency of simulations during this search is unnecessary, so long as all possible permutations

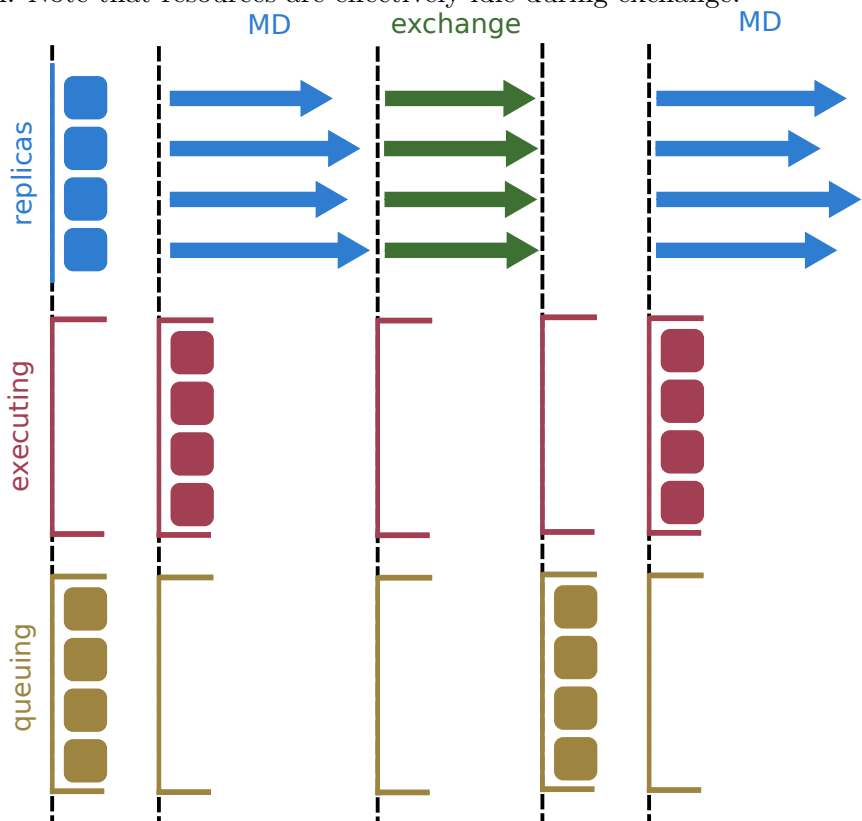
can eventually be sampled. That is, the simulations and exchanges need not occur *synchronously*, but can be performed *asynchronously*. This fact can be used to great effect in a variety of situations, for example, if not enough computational resources are available to simulate a large number of replicas. Another aspect is that the synchronization of resources implies a constant cycle of execution (when computationally demanding simulations are running) and idleness (when less demanding exchanges are occurring). This can be seen as an underutilization of the available computational hardware. As such, asynchronicity can be a valuable tool in enabling large scale simulations with large replica counts and in improving overall utilization of the available computational resources.

3.2.1 Synchronization Modes and Resource Utilization

In the following it is useful to establish a clear distinction between the replicas being simulated and the actual computational resources upon which they are being run. In the traditional REMD scheme, this distinction is not often made, as there is a straightforward assignment of any given replica to a specific computational resource. That is, each replica has a specific resource on which MD is run and this resource is always available. The exchange scheme, which is generally much less computationally expensive than MD, is often run on a single resource by utilizing output from all of the replicas. This scheme is illustrated in Figure 3.7. The important aspects here are that:

1. For a fixed amount of simulation time, not all replicas will take the same amount of real time to complete MD (represented by arrows of unequal length in Figure 3.7).
2. At each exchange all replicas must wait until all other replicas have exchanged in order to continue with MD (represented by arrows of equal length in Figure 3.7).
3. Each simulation effectively has no wait time in the queue, as the necessary resources are allocated at the outset (and it is assumed that this can be done).
4. Although all of the resources are fully utilized during the MD phase, effectively all of the resources are idle during the exchange phase.

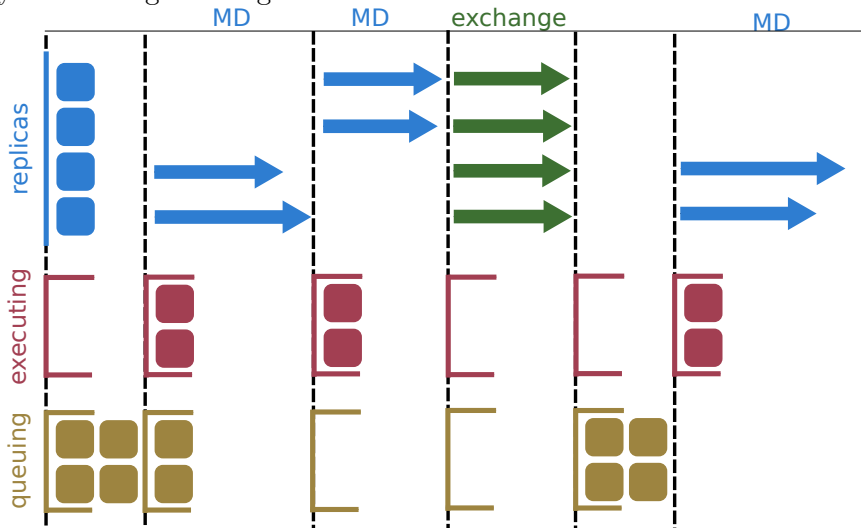
Figure 3.7: Schematic of resource utilization in traditional REMD in which both MD and exchange occur *synchronously*. The amount of computational resources is assumed to be equal to the number of replicas and so there is no queuing of replicas during execution. Note that resources are effectively idle during exchange.



A simple variation on the above scheme is possible if the criteria outlined in item three can not be satisfied, *i.e.*, if not enough resources are available for all replicas at once. In this case the MD phase can simply be performed in stages with some replicas now waiting in a queue before execution (Figure 3.8). The other items above remain unchanged.

In order to address the idle resources in the two schemes above, additional asynchronicity can be introduced. For example, if we continue with the assumption that not enough resources are available for all replicas, then the MD and exchange phases can be interleaved by not requiring that all replicas exchange at the same time (Figure 3.9). This is not problematic so long as the replicas performing MD are not run in

Figure 3.8: Schematic of resource utilization in a variation on traditional REMD in which MD occurs *asynchronously* but exchange occurs *synchronously*. The amount of computational resources is assumed to be less than the number of replicas and so some replicas must wait in the queue during execution. As in Figure 3.7, resources are effectively idle during exchange.



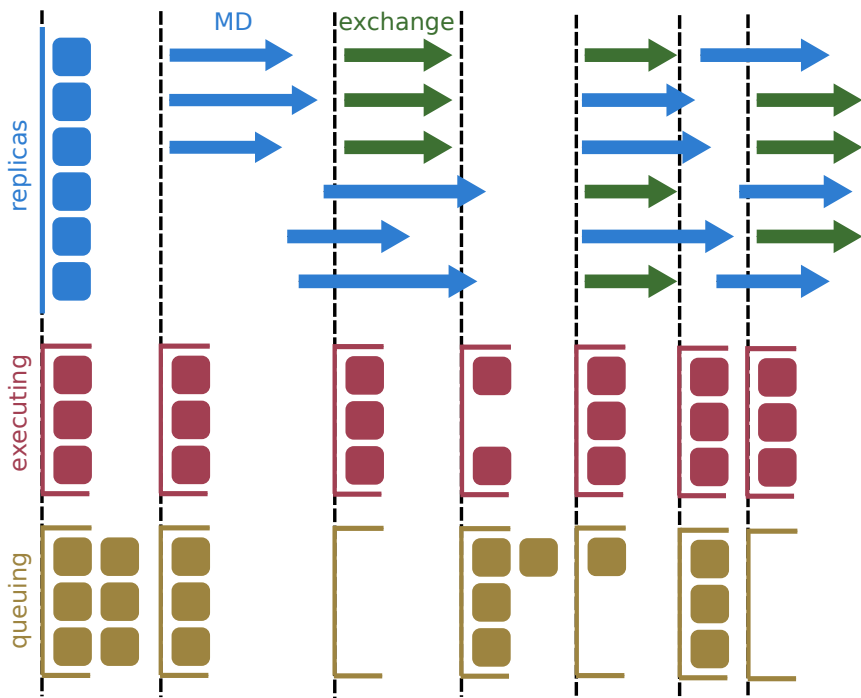
fixed groups. The main change here is that not all replicas will attain the same amount of simulation time between exchanges. However, this is not necessarily different from other schemes, since the exchanges are not always accepted and a rejected exchange is equivalent to no exchange attempt.

Further variations are also possible, but may require more persistent intervention in the MD phase. One such possibility is to enforce synchronization of the replicas in real time by forcibly stopping all of the replicas when the “fastest” replica finishes (*e.g.* the shortest arrow in Figure 3.7). This type of scheme is synchronous in real time, but potentially asynchronous in simulation time (*i.e.* not all replicas will perform the same amount of MD between exchanges). The exchanges can also be performed synchronously (as in Figure 3.10) or asynchronously if there are more replicas than resources.

3.2.2 Application to Multi-Dimensional Free Energy Manifolds

A portable and scalable software framework for implementing the replica exchange schemes described in the preceding section was recently reported by us in the literature[2,

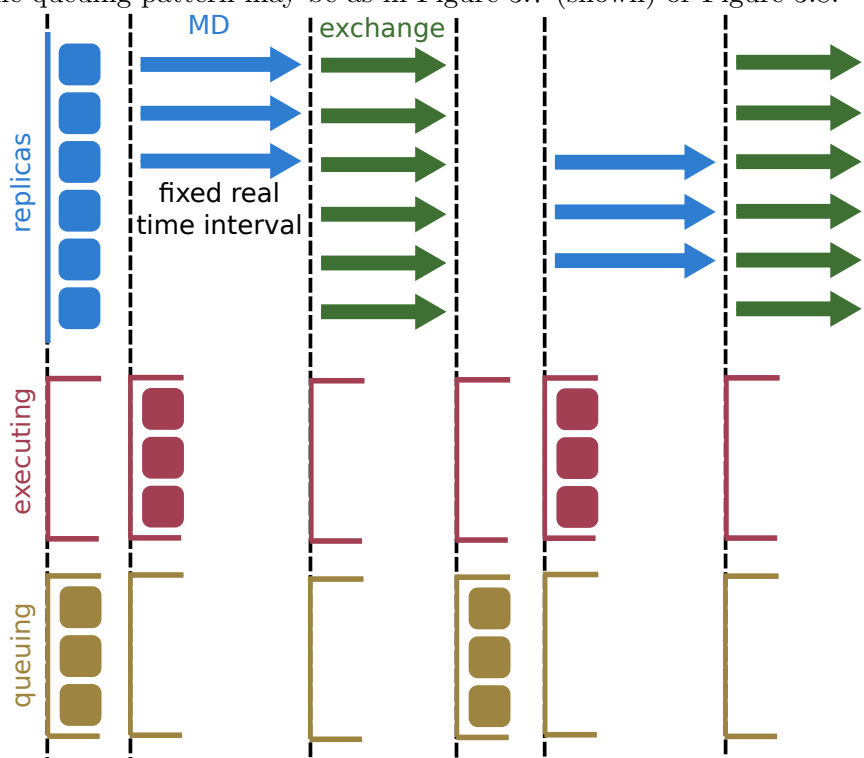
Figure 3.9: Schematic of resource utilization in a REMD scheme in which both MD and exchange occur *asynchronously*. The amount of computational resources is assumed to be less than the number of replicas and so some replicas wait in a queue during execution. However, this time can also be spent undergoing exchanges thereby making more replicas ready for execution and eliminating some of the idle resources during this phase.



3]. These publications include testing and benchmarking of the software as well as some first applications to specific problems in chemical biology. The following is reproduced with permission from the Journal of Chemical Theory and Computation (submitted).

Replica exchange molecular dynamics has emerged as a powerful tool for efficiently sampling free energy landscapes for conformational and chemical transitions. However, daunting challenges remain in efficiently getting such simulations to scale to the very large number of replicas required to address problems in state spaces beyond two dimensions. The development of enabling technology to carry out such simulations is in its infancy, and thus it remains an open question as to which applications demand extension into higher dimensions. In the present work, we apply asynchronous Hamiltonian replica exchange molecular dynamics with a combined quantum mechanical/molecular

Figure 3.10: Schematic of resource utilization in a REMD scheme in which MD occurs on a fixed interval in real time (as opposed to simulation time) and exchange occurs *synchronously*. No assumptions are made on the amount of computational resources and so the queuing pattern may be as in Figure 3.7 (shown) or Figure 3.8.



mechanical potential to explore the conformational space for a uracil ribonucleoside involving >3000 replicas. Two degrees of freedom are used to describe the sugar ring conformation, and one degree of freedom is used to describe the orientation of the nucleobase about the glycosidic bond. The resulting three-dimensional free energy manifold contains complex topological features and correlations that cannot be described in two dimensions. This demonstrates that multi-dimensional free energy manifolds are needed to describe the conformational space for even a simple, fundamental nucleic acid building block.

Introduction

In the past few decades replica exchange molecular dynamics (REMD) has become one of the primary tools with which to improve the accuracy and efficiency of molecular simulations[88, 89]. Examples of REMD now encompass a broad class of schemes ranging from temperature (including novel integration approaches)[103, 90, 91], to Hamiltonian (including alchemical and coordinate biasing)[92, 93, 94, 95] and pH spaces[96, 97, 98], as well as multidimensional combinations thereof[104, 99, 100, 101, 102]. The computational cost of these methods rapidly increases as the number of dimensions (*i.e.* number of variables in the state space being explored) increases. This is because the replica count (naively) increases as N^D , where N is the number of replicas per dimension, D . Hence, a broad practical challenge is to extend these methods so that they can efficiently scale to very large number of replicas, thereby enabling new applications to biological problems involving free energy manifolds of higher dimensions (*i.e.* beyond two). As progress is made toward surmounting this challenge, it is useful to identify fundamental biological problems that demand extension into higher dimensions and to characterize them with benchmark calculations.

Herein we present results from multi-dimensional replica exchange umbrella sampling (REUS) simulations of a single uracil ribonucleoside, applying localized biasing potentials to the key geometric coordinates that dictate the conformation of the ribose ring of the sugar-phosphate backbone (*i.e.* sugar pucker coordinates), and the orientation of the nucleobase about the glycosidic bond (*i.e.* χ torsion angle). The results are used to reconstruct the conformational free energy landscape that reveals a complex topology with a large number of minima subtly connected by correlated pathways. The correlations between these three dimensions are non-obvious *a priori* and non-trivial (or impossible) to recapitulate from lower dimensional surfaces. The lower dimensional surfaces display significant artifacts and bias toward lower energy states, even when the higher energy states should be appreciably populated and/or have biological relevance. The large number of replicas (3432) in this simulation is thus clearly justified in order to correctly characterize the targeted processes.

Further, we present an asynchronous replica exchange framework that provides a general approach for mitigating cost factors and enabling large-scale REUS on the order

of 10^3 replicas or more. The software is agnostic to the underlying molecular dynamics (MD) engine but is demonstrated here with the AMBER[136] package in order to utilize recent developments in quantum mechanical/molecular mechanical (QM/MM) models for accurately modeling sugar puckering modes [116, 137]. Taken together, the present work provides compelling support for the need to address free energy problems within a multi-dimensional framework, produces benchmark simulation results for the conformational free energy landscape of a fundamental nucleic acid building block, and demonstrates that asynchronous exchange is a promising route for taking on this challenge.

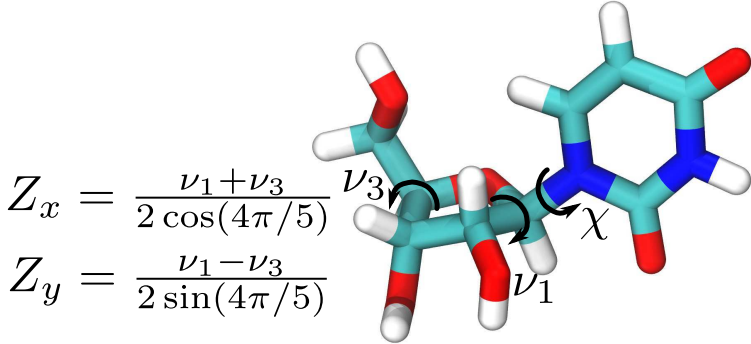
Computational Methods

Molecular Dynamics Each replica was realized as an instance of AMBER 14[136] describing a single, neutral uracil ribonucleoside solvated in a truncated octahedron composed of 1735 TIP4P-Ew rigid water molecules[138] and using periodic boundary conditions with the particle mesh Ewald method[118, 119, 120, 121]. The QM region (uracil) was described by the AM1/d-PhoT Hamiltonian[116] along with a recently developed sugar pucker correction[137] and Lennard-Jones parameters from the AMBER force field[109]. Langevin dynamics was performed at 300 K with a friction coefficient of 5 ps^{-1} and a 1 fs time step.

A total of 3432 replicas were defined by three separate harmonic biases on the χ and ν_1/ν_3 at 30° and 10° intervals, respectively (see Figure 3.11). The full simulation, however, used only 2000 CPU cores on the Stampede cluster at the Texas Advanced Computing Center. This coordinate basis has been shown to be convenient for applying stable, well-defined constraints in quantum chemical calculations, while an alternate basis using linear combinations produces coordinates more recognizably aligned with traditional sugar pucker coordinates[137]. In aggregate, >100 NS of simulation were produced roughly uniformly amongst the replicas, with each replica cycle (*i.e.* the time between exchange attempts) consisting of 500 fs.

Asynchronous Replica Exchange It is important to describe the different modes of *synchronicity*. In the present algorithm, both the MD and exchange protocols are asynchronous across replicas[2]. That is, these processes occur for different replicas at

Figure 3.11: Schematic of dihedral angles used as bias coordinates during umbrella sampling. The proper dihedrals ν_1 and ν_3 are more recognizable as traditional sugar puckering coordinates when taken as the linear combinations Z_x and Z_y as described in Ref. 137 (inset).



different times and never for all replicas at all times (a replica can run MD or exchange, but not both). However, the initiation of these processes is executed synchronously. That is, the controlling process concertedly submits replicas for MD and then coordinates exchanges amongst those that are not running. Since the latter process can be increasingly time consuming at large replica counts, we find it useful to oversubscribe replicas so that resources are taken up as they become available, even if the main process is busy coordinating exchanges.

Finally, the protocol here requires that not all replicas be available for exchange at all times. However, conventional nearest neighbor-type exchange schemes require the opposite scenario (*i.e.* all replicas must have an exchange partner at all times), but there is no physical or mathematical reason for this requirement beyond algorithmic convenience. Instead, we follow recent statistical developments[106, 107] and use the information from all available replicas (rather than a set number of replica pairs) to perform exchanges. This can be a more computationally expensive procedure, but the added cost is generally worthwhile given the vast improvement in the acceptance rate.

The primary conceptual advance underlying our implementation is the decoupling of the replica exchange algorithm details from the execution details of the replicas on high-performance resources[139]. This enables the efficient execution of a range of replica exchange schemes. An early prototype of the software system used for performing

current simulations using the current asynchronous protocol has been described in Ref. 2. Significant performance enhancements and improvements continue to be made; these enhancements and updated software are (and will be) publicly available online[140].

Results and Discussion

Umbrella sampling simulations in multiple dimensions are considerably complex, and analysis of the data required to construct a free energy manifold must be done carefully. We apply the multistate Bennett acceptance ratio[71, 133] (MBAR) method in tandem with a three dimensional Gaussian kernel density estimator (see Ref. 1 for details) along the χ , Z_x , and Z_y coordinates shown in Figures 3.11 and 3.12. Due to the large amount of data and memory restrictions (a full calculation required >20 GB of memory), the data was divided into two non-overlapping sets of states by taking every other χ value (*i.e.* they are segregated by placement of the biasing potential). These sets are not rigorously statistically independent,⁶ but the fact that the results from both data sets along this coordinate agree within statistical error (see Figure 3.13) provides some degree of confirmation of convergence. For simplicity, in the discussion that follows, data presented after Figure 3.13 are from only one of these data sets.

Low-Dimensional Free Energy Profiles Given Inadequate Representations

A general assumption used in the analysis of data from free energy simulations is that all degrees of freedom orthogonal to the chosen coordinate(s) are not strongly coupled to the process under investigation. If this is not true, significant artifacts can be encountered in the interpretation of the results. Generally it is assumed that the conformational states of a nucleoside can be enumerated as four discrete states based on a binary distinction between the sugar pucker mode (*C3'-endo* or *C2'-endo*) and the nucleobase orientation (*syn* or *anti*). A weak coupling between the coordinates connecting these states would imply that transitions between states are affected by, but do not directly involve, the orthogonal coordinate(s).

As an example, consider the univariate free energy profile obtained from the process of rotating the χ torsion between the *syn* and *anti* conformations (Figure 3.13), red

⁶ The analysis of replica exchange data by techniques such as MBAR generally assumes that each *replica* is statistically independent (see Ref. 74 for further discussion).

Figure 3.12: Three-dimensional free energy manifold for solvated uracil along the Z_x , Z_y , and χ coordinates (see Figure 3.11). Cross-sections in the Z_x, Z_y -plane are shown for multiple local minima (blue spheres) and indicate that many minima are connected out-of-plane by one or more saddle points (red spheres). Energies are in kcal/mol and axes are in degrees.

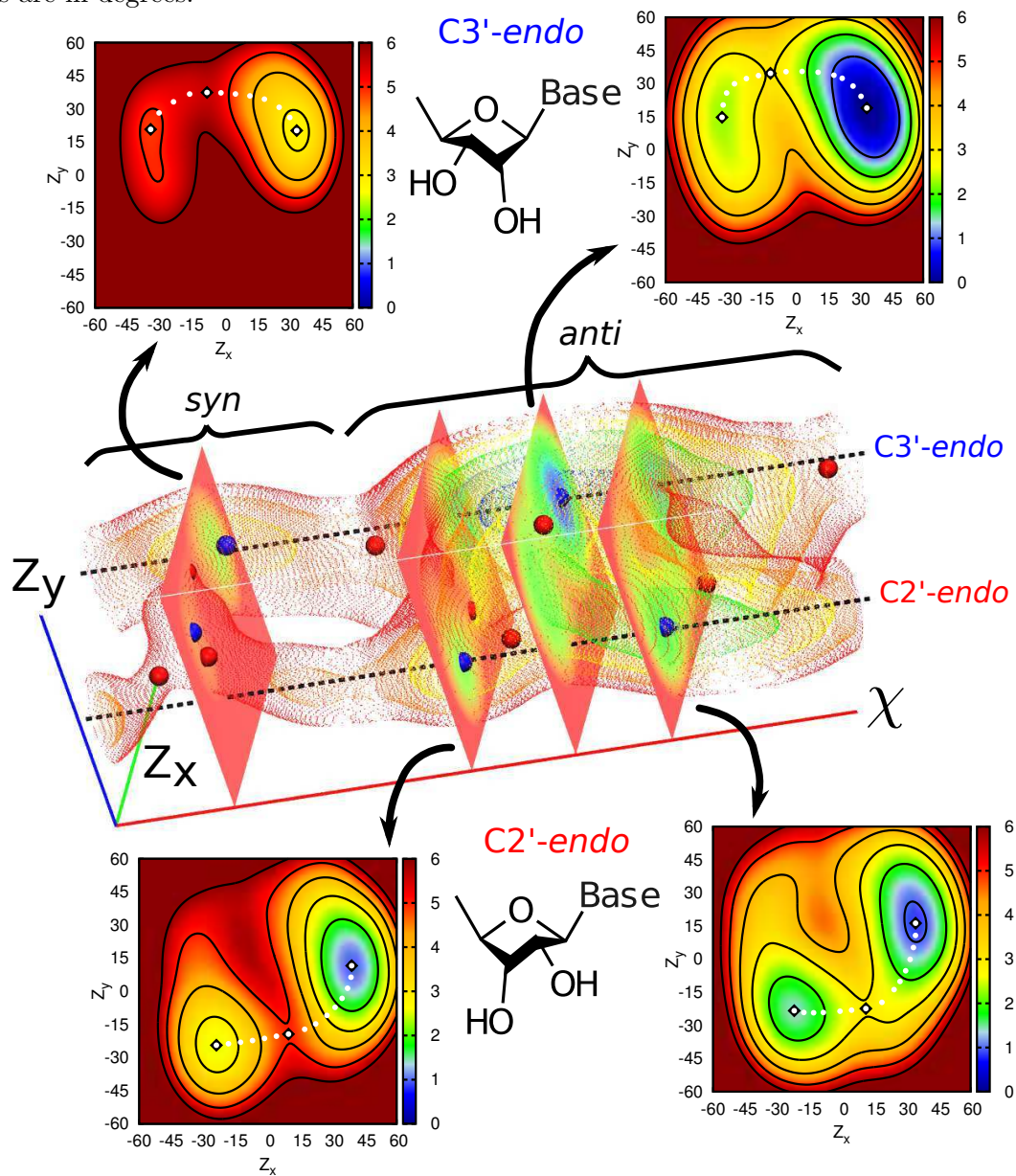
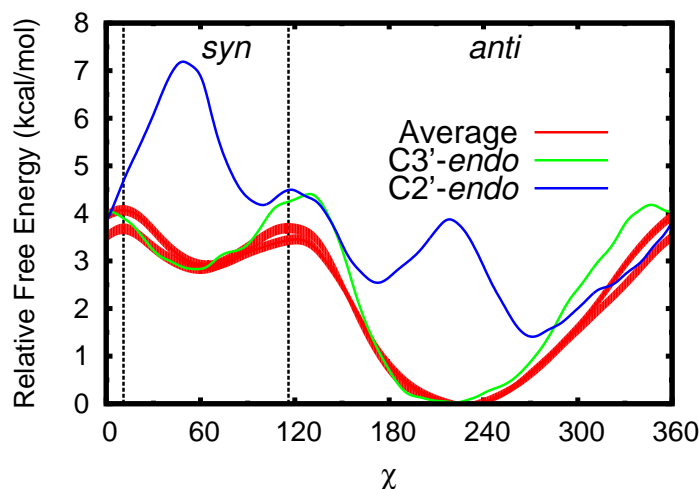


Figure 3.13: Free energy profiles for solvated uracil along its χ torsion using multiple schemes to reduce dimensionality. The whole data set can be used by taking the Boltzmann weighted average for all sugar pucker values (red, two curves corresponding to two statistically equivalent data sets with 95% confidence intervals). Conversely, fixed pairs of Z_x and Z_y can be followed, in this case corresponding to the average C2'-endo (blue) and C3'-endo (green) minima (see Figure 3.14).



filled curves). The average result, obtained from sampling in all possible states, is to be contrasted with those obtained in a localized sugar puckering mode, comparable to sample sets in which the sugar puckering mode never changes (Figure 3.13, blue and green curves). It is evident, even from brief inspection of Figure 3.13, that the one-dimensional χ profile depends considerably on the sugar pucker state. However, the average profile is dominated by contributions from the lower energy (by <2 kcal/mol) C3'-endo state, as borne out by the high degree of similarity between these two curves (Figure 3.13, red and green lines). The profile from the higher energy, but still significant, C2'-endo state has very different minima and maxima, and splits the main *anti* basin into two states that are nearly energetically degenerate.

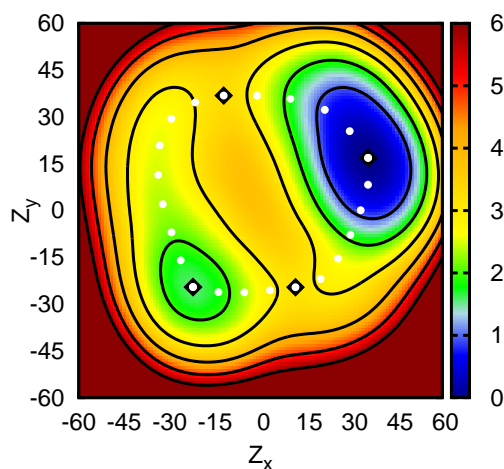
Thus, straightforwardly averaging along the single χ degree of freedom clearly removes a considerable amount of information from the full free energy profile. This may be considerably problematic, as the energetically low lying, but geometrically different, states in a C2'-endo sugar pucker could be of significant interest when studying larger

questions concerning RNA conformation. Also of note is that the relative shifts between the pucker-dependent χ profiles *cannot* be determined from separate simulations in those sugar pucker modes alone. This information is only available here because the sugar pucker coordinates were extensively sampled and connected by (and subsequently extracted from) a three-dimensional free energy manifold. Obtaining the two profiles separately could lead to erroneous interpretation.

A similar situation is encountered when analyzing the two-dimensional profile for sugar pucker coordinates alone. In this case the average profile contains two minima broadly recognizable as *C2'-endo* and *C3'-endo* states. These minima are connected by two transition states in a periodic fashion along the pseudorotation angle (Figure 3.14). However, as above, these minima and transition states only describe the full transformations in a coarse sense. If multiple two-dimensional surfaces are obtained departing from a specific conformation for the orthogonal χ coordinate (and subsequently not visiting other conformations within the time scale of the simulations), this average surface would look quite different. This is because the *syn* states, while clearly of interest in larger nucleic acid systems, are high enough in energy that their contributions to the sugar pucker profiles are small compared to the *anti* states.

Three-Dimensional Free Energy Manifold Unveils “Hidden” Pathways Lastly, the complete results for a three-dimensional free energy manifold for a ribonucleoside are discussed in detail. The most distinctive (and perhaps most unexpected) result is the presence of five, not four, stable minima in the complete coordinate space (Figure 3.12, blue spheres). This is because there are two *anti/C2'-endo* states, where only one might be expected based on a binary segregation of conformations. Interestingly, both of these degenerate minima are directly accessible from the *anti/C3'-endo* global minimum, but via different transition states (Figure 3.12, red spheres). Furthermore, in order to reach a true minimum, all of these transitions require motion in *both* the sugar pucker and χ coordinates and in potentially different orders. For example, the global minimum (*anti/C3'-endo*) can transition to the next lowest energy minimum (*anti/C3'-endo*) either by a rotation in χ followed by a pseudorotation of the ring or *vice versa*. An aspect of the free energy analysis that is meant to be especially emphasized here is that the extra “hidden” minimum and the accompanying pathways are not

Figure 3.14: Boltzmann weighted average free energy surface for solvated uracil along the Z_x and Z_y sugar pucker coordinates (see Figure 3.11). Two minima are observed roughly corresponding to C3'-*endo* and C2'-*endo* pucker states (black diamonds, also marking saddle points). An apparently periodic transition between these states along the pseudorotation angle ($P_\theta = \arctan \frac{Z_y}{Z_x}$) is also observed (white dots). Energies are in kcal/mol and axes are in degrees.



at all evident from the low dimensional analysis described above. Proper identification and characterization of the conformational transitions can only be performed with an exhaustive search of the collective coordinate spaces. The alternative is to make the, in this instance, quite incorrect assumption that these coordinates are uncorrelated. In general then, an attractive strategy is to expand the dimensionality of the coordinate search and determine uncorrelated degrees of freedom *a posteriori*. Since the only significant drawback to this approach is the potentially immense cost increase (in terms of processor hours), techniques that decrease that cost are of significant value. The asynchronous protocol described and used here provides such a tool for reducing the amount of computational resources needed and in future work will be extended and optimized for load balancing and multiple resource management.

Conclusion

REMD simulations are potentially powerful tools for improving the accuracy of molecular simulations. In this work, low-dimensional REMD simulations for a fundamentally

simple nucleic acid system, a single uracil ribonucleoside, are seen to provide an incomplete picture of the free energy landscape. Moreover, moving to higher dimensions reveals subtle correlations between the fundamental nucleic acid backbone and base orientation coordinates. Finally, the significant added cost of this multi-dimensional simulation is addressed by an asynchronous exchange framework, which provides a general approach for tackling multi-dimensional REMD on arbitrary software platforms. This software tool is readily extendable to other problems of conformational transitions as well as those addressing chemical reactions.

Chapter 4

Computational Investigations of Phosphoryl Transfer in Non-Ribozymatic Systems

The core concept underlying catalysis, of any kind, is that the kinetic rate of a chemical reaction can be enhanced by changes in the reaction pathway. These changes can be structural (*i.e.* the geometry of atoms is different as the reaction progresses) or energetic (*i.e.* new interactions or the extent of existing interactions are different as the reaction progresses). Identifying and characterizing the specific factors contributing to catalysis thus requires extensive knowledge of the *uncatalyzed* reaction, for it is only in this way that the aspects unique to the enhanced pathway can be separated from the intrinsic reaction character. For the 2'-*O*-transesterification reaction catalyzed by many small ribozymes (see Section 2.1) the core “uncatalyzed” reaction is often chosen as solvent catalyzed cleavage. This is because the completely uncatalyzed reaction is relatively slow so as to be experimentally inconvenient. In any event, the main catalytic step in solvent catalyzed cleavage is often assumed to be the facilitation of one or both of the proton transfer steps and so the key phosphoryl transfer step is largely unchanged. To this end, a number of simulations and experiments were performed by ourselves and our collaborators in order to better understand the most basic aspects of 2'-*O*-transesterification of RNA. These studies first explored the base catalyzed reaction in

solution[141, 1, 142] and were later extended to a similar analysis of protein enzymes that catalyze RNA cleavage, specifically ribonuclease A[6].

The specific impetus behind the new work presented here was a recently started collaboration with Michael Harris and Joseph Piccirilli. The aim of their planned experimental studies was to measure kinetic isotope effects (KIEs) for RNA cleavage in a variety of conditions, both non-enzymatic and enzymatic. This is challenging for several reasons. First, KIEs, especially the ^{18}O KIEs measured here, require extremely high precision measurements with specially synthesized, isotopically enriched compounds[143]. Second, although KIEs are one of the most sensitive probes to chemical mechanism, a straightforward understanding of such results is non-trivial and computational studies can be a crucial tool in such interpretations.

In what follows, measurements on a simple RNA dinucleotide are compared against molecular simulations of solution compounds. These simulations provide atomic detail for the measurements and are validated by close agreement with pH-rate measurements. In what follows this work is reproduced with permission from Ref. 1. Copyright 2013 American Chemical Society. The supporting information is also provided as Appendix A. Subsequent collaborative works have also been published[142, 6] but are somewhat outside the scope of this work. Accordingly, these are only briefly discussed in terms of the aspects that are most relevant to Section 5.

4.1 Molecular Simulations of RNA 2'-*O*-Transesterification Reaction Models in Solution

We employ quantum mechanical/molecular mechanical umbrella sampling simulations to probe the free energy surfaces of a series of increasingly complex reaction models of RNA 2'-*O*-transesterification in aqueous solution under alkaline conditions. Such models are valuable for understanding the uncatalyzed processes underlying catalytic cleavage of the phosphodiester backbone of RNA, a reaction of fundamental importance in biology. The chemically reactive atoms are modeled by the AM1/d-PhoT quantum model for phosphoryl transfer, whereas the aqueous solvation environment is modeled with a molecular mechanics force field. Several simulation protocols were compared that used different ionic conditions and force field models. The results provide insight

into how variation of the structural environment of the nucleophile and leaving group affects the free energy profile for the transesterification reaction. Results for a simple RNA backbone model are compared with recent experiments by Harris *et al.* on the specific base catalyzed cleavage of a UpG dinucleotide. The calculated and measured free energies of activation match extremely well ($\Delta F^\ddagger = 19.9$ - 20.8 versus 19.9 kcal/mol). Solvation is seen to play a crucial role and is characterized by a network of hydrogen bonds that envelopes the pentacoordinate dianionic phosphorane transition state and provides preferential stabilization relative to the reactant state.

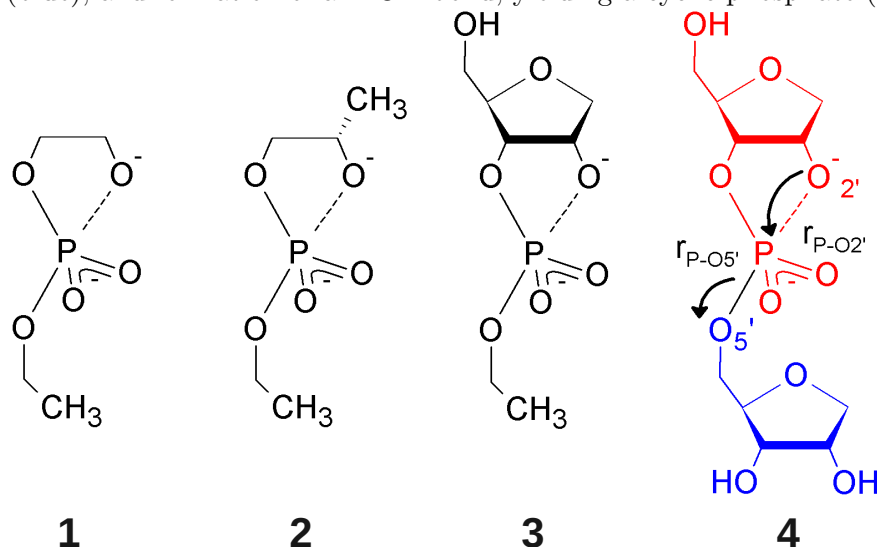
4.1.1 Introduction

Cleavage of the phosphodiester backbone of RNA is an essential reaction in biology that is fundamental to many important biological processes ranging from gene splicing and regulation to viral replication and cell signaling[23]. It is thus significant that several small RNA molecules, such as the hammerhead[9], hairpin[8], hepatitis delta virus[144], Varkud satellite[145], and *glmS*[146] ribozymes, catalyze the phosphoryl cleavage of their own backbones. While the secondary and tertiary structure of these ribozymes are all distinct and their optimal ion identity and concentration requirements differ significantly, they all catalyze the same intramolecular 2'-*O*-transesterification reaction and form a 2',3'-cyclic phosphate and 5'-hydroxyl as products[147, 29, 148].

As with all catalytic reactions, the mechanistic features of ribozymes and protein enzymes are inherently related to their rate enhancement relative to the background rate of the non-catalytic reaction, in this case the cleavage of an RNA backbone in aqueous solution. Hence, a logical starting point for determining the key mechanistic characteristics of self-cleaving ribozymes would be to first determine the characteristics of the uncatalyzed mechanism. Indeed, there is an established literature concerning model compounds for phosphate 2'-*O*-transesterification,[149, 150, 151, 152, 153] but these studies frequently focus on non-native, enhanced leaving groups[154, 155, 156], reactions perturbed by chemical markers needed for spectroscopic analysis[157], or temperature ranges far from normal biological conditions[154, 158]. Nonetheless, these studies provide a firm experimental baseline for comparisons between native uncatalyzed and catalyzed reactions.

Theoretical and computational approaches, particularly molecular dynamics (MD),

Figure 4.1: Reaction models for base-catalyzed phosphoryl transfer. Structural complexity increases from model **1** (2-hydroxy-ethyl ethyl phosphate) to model **4** (an abasic RNA dinucleotide) to bridge the gap between simple alkyl phosphates and the more complex RNA backbone. Ribose ring naming conventions are adopted for consistency such that all systems are said to undergo cleavage of the P-O5' bond, yielding a primary alkoxide (blue), and formation of a P-O2' bond, yielding a cyclic phosphate (red).



have emerged as a valuable tool in the study of chemical reactions because they allow access to full atomistic detail. However, the degree to which meaningful insights into mechanism can be gained from simulations relies on the accuracy of the models that are employed for the particular system under study. This leads to a natural synergistic relationship between experiment and theory, with experiment providing key benchmarks and theory providing detailed molecular level interpretations and testable predictions. This approach is well illustrated by the development of fast and accurate quantum mechanical/molecular mechanical (QM/MM) methods calibrated to specific experimental and *ab initio* data[116], which has opened the door for accurate tests of explicit reaction pathways, even for large biomolecules[159, 160]. The free energy surfaces of such pathways can be rigorously compared to experimental kinetics measurements and also have strong connections to highly sensitive mechanistic probes such as thio and isotope[142] as well as mutational[161] effects.

In order to rationally decompose the complexity of RNA backbone cleavage, the

present work focuses on mapping the free energy profiles of a series of molecules undergoing transphosphorylation to form cyclic phosphates under basic conditions (Figure 4.1). The systems are akin to the simple UpG dinucleotide recently studied experimentally by Harris, *et al.*[141] The dinucleotide UpG matches the sequence at the cleavage site of the self-cleaving hepatitis delta virus ribozyme[162] and is an active substrate for ribonuclease A[163]. Its cleavage mechanism in solution is therefore a valuable benchmark. In solution, the observed first order rate constant for UpG cleavage increases linearly from pH 10-13, becoming pH independent beyond that point; extrapolating to “infinite” pH gives an intrinsic rate constant of 0.06 s^{-1} [141]. Using novel techniques[143, 141], Harris, *et al.* also measured primary and secondary kinetic isotope effects (KIEs) for base catalyzed UpG cleavage, as well as the solvent D₂O effect. The lack of any significant solvent D₂O effect combined with an estimated correction for equilibrium isotope effects on the nucleophile confirms the conventionally accepted specific base mechanism. Taken together, the KIEs suggest that UpG undergoes a concerted mechanism with a “late,” product-like, transition state.

On the basis of QM/MM MD simulations, a theoretical free energy of activation is calculated here that allows direct comparison to the experiments of Harris, *et al.* The quality of the comparison begets significant confidence in also using the simulations to characterize the reactant and transition states structurally, as well as analyze the role of water and ions in the reaction. Additionally, because of the wide range of models available for (and commonly used in) QM/MM studies, the sensitivity of the results to different water models and Lennard-Jones parameters is compared and discussed.

4.1.2 Computational Methods

Molecular Dynamics

MD simulations were performed using the AMBER 12[164] suite of programs. An integration step of either 1 or 2 fs was used depending on whether or not the SHAKE[165] algorithm (tolerance = 1.0×10^{-8}) was used to constrain bonds with hydrogen in the solute; the SETTLE[166] algorithm was always used to constrain rigid water molecules. Temperature and pressure were regulated with the methods of Andersen[49] (310 K, “massive” collisions every 2000 steps) and Berendsen[167] (1 bar, time constant of 5

ps, compressibility = $44.6 \times 10^{-6} \text{ bar}^{-1}$), respectively. Long range electrostatics were treated using periodic boundary conditions in a rhombic dodecahedron cell and the particle mesh Ewald (PME) method for both molecular mechanical (MM)[118, 119] and QM/MM[120, 121] calculations. PME calculations employed 6th order B-spline interpolation with 50 grid points ($\approx 1 \text{ point}/\text{\AA}$) along each axis; the Ewald coefficient was chosen such that the estimated error in the direct space energy was on the order of 10^{-5} kcal/mol . QM/MM Ewald calculations used a reciprocal space defined by $k_{\text{max}} = 7$ and $k_{\text{max}}^2 = 98$; the QM/MM Ewald coefficient was separately chosen as $10V^{-1/3}$, where V is the cell volume in \AA^3 (see Supporting Information for details). Lennard-Jones and direct space Coulombic interactions were truncated at 10 \AA . For the QM/MM direct space, an atom based switching function was applied between 8 and 10 \AA .

All QM/MM simulations used the AM1/d-PhoT semi-empirical Hamiltonian[116], with the QM region defined as the entire solute. Lennard-Jones parameters were taken from either the AMBER FF10[109, 110, 111, 112] or CHARMM27[113, 114] nucleic acid force fields, with the exception of select interactions with sodium ions (see Discussion and Supporting Information). The solvent environment was modeled using either the TIP3P[168] or TIP4P-Ew[138] rigid water model and the associated alkali metal and halide ion parameters of Joung and Cheatham[169]. Simulations contained 2,640 solvent molecules (e.g., water molecules or water molecules in approximately 140 mM NaCl , see Supporting Information).

The selected model is appropriate for several reasons and similar QM/MM models have been successfully used elsewhere in studies of both enzymatic[159, 170, 160] and non-enzymatic [171, 172] phosphoryl transfer. First, AM1/d-PhoT is specifically parameterized to reproduce gas phase *ab initio* calculations of an extensive set of phosphate containing compounds and reactions (see Ref. 116 and 173 and the Supporting Information) and has also been shown to be the best choice for reproducing geometries and energies of penta-coordinated phosphorous systems amongst several common semi-empirical methods[174]. Second, a QM/MM approach to solvation (*i.e.* neglecting a QM description of the solvent) is well suited for the current application since chemical participation of water (*e.g.* via hydrolysis or proton transfer) is not expected to occur[152, 141]. An intermediate description including *some* water in the QM region in

an adaptive fashion (a necessary consequence of diffusion in a fully solvent exposed reaction) could potentially be advantageous, but such an approach is difficult to implement with smooth gradients suitable for dynamics and not yet widely available[175]. Therefore, at present, a QM/MM model is expected to be preferable to a full QM description due to the vastly decreased cost needed to obtain adequate sampling and since the bulk properties of MM water models are generally superior to both semi-empirical[176] and even certain *ab initio* quantum models[177].

Umbrella Sampling

QM/MM MD umbrella sampling[65] trajectories were performed along a mass weighted atom transfer coordinate, $\xi = \frac{1}{2}(r_{\text{P-O5}'} - r_{\text{P-O2}'})$, where $r_{\text{A-B}}$ is the distance between atoms A and B and the factor of one half arises from the leaving group and nucleophile masses being the same (it should be noted that, in the present case, this has no effect on the thermodynamics). Although such a coordinate has been widely used in the literature, especially with regards to phosphoryl transfer reactions [171, 159, 172, 170, 160, 178, 142], a recent study by Rosta, *et al.* has suggested that the calculated free energy barrier is potentially sensitive to this choice properly capturing orthogonal chemical events such as proton transfer[179]. This is not anticipated to be an issue here because the reactions take place in the high pH regime where the assumption of rapid, uncoupled deprotonation of the nucleophile is well-justified[141]. The orthogonal events are thus entirely structural and not chemical (*i.e.* they involve solvent rearrangement).

After an extensive initial equilibration protocol (see Supporting Information) production consisted of 1 ns (2 ns for the dinucleotide system) for each window, with sampling omitting the first 250 ps for relaxation/equilibration within the window. The value of the progress coordinate was stored at 0.5 ps intervals for analysis using the multistate Bennett acceptance ratio (MBAR)[71]. Approximately uncorrelated data sets were obtained by subsampling configurations at intervals equal to the statistical inefficiencies, which were estimated in each simulation by direct integration of the autocorrelation function of the progress coordinate using the fast, adaptive integration scheme of Chodera, *et al.*[74].

MBAR Analysis and Free Energy Profiles

MBAR provides a general formalism for re-weighting mechanical observables for estimation in arbitrary thermodynamic states provided that the relative statistical weight in those states is known. The MBAR estimator for the expectation of an observable, $\langle A \rangle$, that depends only on the configuration, \mathbf{x} , is given by[71]:

$$\hat{A} = e^{\hat{f}} \sum_{m=1}^M \sum_{n=1}^{N_m} w(\mathbf{x}_{mn}) A(\mathbf{x}_{mn});$$

$$w(\mathbf{x}_{mn}) \equiv \left[\sum_{l=1}^M N_l e^{\hat{f}_l - [u_l(\mathbf{x}_{mn}) - u(\mathbf{x}_{mn})]} \right]^{-1}, \quad (4.1)$$

where M is the number of states, N_m is the number of samples from state m , \hat{f}_l is the MBAR estimate of the free energy of state l relative to an arbitrary state (here $f_1 \equiv 0$), and $u_l(\mathbf{x}) \equiv \beta U_l(\mathbf{x})$ is the “reduced potential” characterizing state l [67, 71]. β and $U(\mathbf{x})$ are the inverse temperature (for simplicity assumed to be the same in all states) and potential energy. This expression is quite general, but in the present context of umbrella sampling the unindexed values refer to the unbiased state and the sample configurations, \mathbf{x}_{mn} , are drawn from M biased states.

One method of estimating the free energy profile, $F(\xi)$, is to estimate the marginal distribution, $\rho(\xi) = \langle \delta(\xi(\mathbf{x}) - \xi) \rangle$; the free energy profile, up to an additive constant, is then simply $F(\xi) = (-1/\beta) \ln \rho(\xi)$. Since the delta function is only a function in the distributional sense, an approximate estimator is needed for finite sampling. A broad class of such estimators are known in the statistics literature as kernel density estimators[128, 129, 125, 126], and may take the following form:

$$\langle \delta(\xi(\mathbf{x}) - \xi) \rangle = \lim_{h \rightarrow 0} \left\langle \frac{1}{h} K \left(\frac{\xi(\mathbf{x}) - \xi}{h} \right) \right\rangle, \quad (4.2)$$

The function K is often referred to as a *kernel* and the parameter h as the *bandwidth*. A case more common to the chemical literature is when K is an indicator function[68, 74, 70]; this returns the familiar histogram estimator and h is recognized as the bin width, with an additional parameter defining the bin center. The results obtained with a histogram estimator are often qualitatively, and even quantitatively, similar to those

obtained using a kernel density estimator. This general trend is confirmed in the present work (see Supporting Information). The marginal distribution could also be calculated using other estimators, such as those with a parametric form[81, 180], although this may require slight variation of the MBAR formalism.

In order to obtain a kernel based estimator for $F(\xi)$, Eqn. 4.2 is substituted into Eqn. 4.1 and the logarithm is taken:

$$-\beta\hat{F}(\xi) = \hat{f} + \ln \sum_{m=1}^M \sum_{n=1}^{N_m} w(\mathbf{x}_{mn}) K(\xi; \mathbf{x}_{mn}, h_m); \quad (4.3)$$

$$K(\xi; \mathbf{x}_{mn}, h_m) \equiv \frac{1}{h_m} K\left(\frac{\xi(\mathbf{x}_{mn}) - \xi}{h_m}\right)$$

As noted above, the additive constant \hat{f} can be arbitrarily set to zero (or any other convenient value). This work employs a standard normal kernel density estimator with the bandwidth chosen in each window as twice that given by the data based algorithm of Sheather and Jones[181] (see Supporting Information), hence the bandwidth is shown to vary amongst states.

A more unusual class of observables can be defined as expectations *along* a coordinate ξ :

$$\langle A \rangle_\xi = \frac{\langle A \delta(\xi(\mathbf{x}) - \xi) \rangle}{\langle \delta(\xi(\mathbf{x}) - \xi) \rangle} \quad (4.4)$$

Note that Eqn. 4.4 is expressed as a ratio of expectations in an unconstrained ensemble, rather than as an expectation in a constrained ensemble (*i.e.* the momentum conjugate to ξ is non-zero). Following a similar process as above and recognizing the denominator as being related to Eqn. 4.3, the following estimator is obtained:

$$\hat{A}(\xi) = e^{\beta\hat{F}(\xi) + \hat{f}} \times \sum_{m=1}^M \sum_{n=1}^{N_m} w(\mathbf{x}_{mn}) A(\mathbf{x}_{mn}) K(\xi; \mathbf{x}_{mn}, h_m) \quad (4.5)$$

Note that in this context \hat{f} is *not* arbitrary, although it could be made zero in certain contexts.

All of the estimators shown here were implemented in a locally modified version of

the Python MBAR implementation by Shirts and Chodera (pymbar v2.0)[71]. Visualization and other analyses were performed using Visual Molecular Dynamics (v1.8.7)[182], particularly the VolMap plugin (default, unscaled radii, 0.5 Å resolution).

4.1.3 Results

Validation of QM/MM Model and Reaction Coordinate

Both semi-empirical quantum methods and approximate reaction coordinates require some degree of caution when used in simulations; both can lead to significant deviation from physical behavior. In addition to the extensive validation and use in the literature of AM1/d-PhoT in conjunction with the simple atom transfer coordinate used here[159, 170, 160, 171, 172, 178, 142], we briefly present potential energy profiles comparing AM1/d-PhoT to standard B3LYP results as well as the reaction coordinate paths with the optimized stationary points. For simplicity, the gas phase 2'-*O*-transesterification reaction of 2-hydroxy ethyl phosphate (similar to model **1** in Figure 4.1) is considered. As is clear from Figure 4.2 and Table 4.1, AM1/d-PhoT provides excellent agreement with B3LYP/6-31+G(d), both in the energy barrier and geometry and location of minima and saddle points. In all cases, the approximate reaction coordinate correctly follows the reaction progress from reactant to product and predicts stationary points that are similar in geometry. Interestingly, B3LYP with a slightly smaller basis set (as one might consider using in QM/MM simulations due to lower cost) yields a substantial underestimate of the reaction barrier compared to both the higher basis set and AM1/d-PhoT, although the geometries are still comparable. Lastly, it is worthwhile to note that AM1/d-PhoT was actually trained and tested on even higher level results (B3LYP/6-311++G(3df,2p)//B3LYP/6-31++G(d,p)) [116, 173], which, at the very least, explains the existence of the minor deviations in energy and geometry visible in Figure 4.2.

Figure 4.2: Gas phase potential energy profiles for the 2'-*O*-transesterification of 2-hydroxy ethyl methyl phosphate (see inset) at the B3LYP/6-31+G(d) (red), B3LYP/6-31G(d) (blue), and AM1/d-PhoT (green) levels (as implemented in Gaussian 09[183] and AMBER12/AmberTools12[164, 121] respectively). The approximate reaction coordinate paths are obtained via constrained optimization at different coordinate values (using the DL-FIND library[184]). Crosses denote the location of optimized minima and transition states. All energies are relative to the optimized minimum at the relevant level of theory.

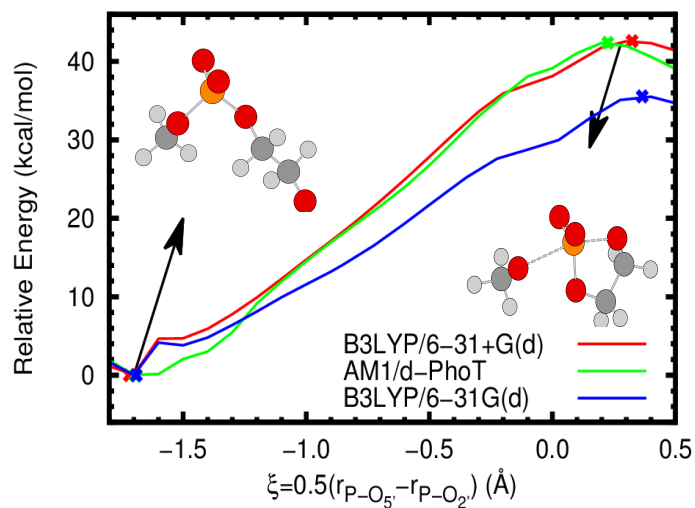


Table 4.1: Potential energy barriers and select geometric quantities at the transition state at different levels of theory using multiple reaction coordinates. The intrinsic reaction coordinate (IRC) is defined so as exactly to connect the minimum and first order saddle point. However, the approximate reaction coordinate (ARC) used in simulations is not guaranteed to connect either point on the potential energy surface. In this case the minimum and saddle point are defined by the reduced set of coordinates. It should be noted that the ARC values are obtained under the additional approximation that $m\Delta\xi = \frac{1}{2}(r_{\text{P-O5}'} - r_{\text{P-O2}'})$, with $\Delta\xi = 0.1 \text{ \AA}$ and m is some integer. Energies are in kcal/mol and lengths are in \AA .

		ΔE^\ddagger	$r_{\text{P-O5}',\text{TS}}$	$r_{\text{P-O2}',\text{TS}}$
B3LYP/6-31+G(d)	IRC	42.6	2.49	1.84
	ARC	42.5	2.44	1.84
AM1/d-PhoT	IRC	42.4	2.33	1.88
	ARC	42.3	2.27	1.87
B3LYP/6-31G(d)	IRC	35.5	2.56	1.83
	ARC	35.5	2.63	1.83

Free Energy Profiles and Mechanical Observables

The principle results of the present work are the free energy profiles, $F(\xi)$, for each of the specific base catalyzed reactions calculated from umbrella sampling simulations. The values of the progress coordinate, ξ , corresponding to the reactant and transition state (ξ_{R} and ξ^\ddagger , respectively) are determined as:

$$\begin{aligned} \xi_{\text{R}} &= \arg \min_{\xi < \xi^\ddagger} F(\xi) \\ \xi^\ddagger &= \arg \max_{\xi} F(\xi) \end{aligned} \tag{4.6}$$

The free energy barrier is then calculated as $\Delta F^\ddagger = F(\xi^\ddagger) - F(\xi_R)$. Additionally, the averages of select mechanical observables, $\langle A \rangle_\xi$, and their standard deviations, $(\langle A^2 \rangle_\xi - \langle A \rangle_\xi^2)^{\frac{1}{2}}$, were estimated at fixed values of the reaction coordinate (Eqn. 4.5).

Throughout this work the term “reaction model” is used to refer to a molecule that undergoes a reaction analogous to RNA transesterification (*i.e.*, contains a phosphodiester that reacts to form a cyclic phosphate). The reaction models used here are illustrated and numbered in Figure 4.1. In QM/MM simulations of these reaction models a “force field model” must also be chosen to describe the (non-bonded) MM and QM/MM interactions between the (QM) solute and solvent. Here the force field models for the solutes take parameters from either the AMBER (A) or CHARMM (C) force fields in conjunction with solvent (water or water + NaCl) defined by the TIP3P (3) or TIP4P-Ew (4) water models. A full simulation model is then given by both a reaction model and force field model. For example, an abasic RNA dinucleotide (reaction model 4 in Figure 4.1) with AMBER force field parameters in a simulation cell containing TIP4P-Ew water and sodium chloride is designated as 4-A4/NaCl. All models will hereafter be referred to with this nomenclature.

Abasic Dinucleotide Models

Solvent Environments We begin by examining several possible solvation models of an abasic RNA dinucleotide, for which a wealth of experimental data is available[149, 150, 152, 141]. In particular, we examine differences in the free energy profile due to variations in the water model and ion atmosphere. The results (Figure 4.3, Table 4.2) show no statistically significant variation in the reaction barrier or geometry when the water model is changed from TIP3P to TIP4P-Ew. The removal of ions (infinite dilution limit) appears to slightly lower the barrier for both water models by the same amount, although this difference is very similar in magnitude to the estimated error (0.8 ± 0.6 kcal/mol).

Figure 4.3: Free energy profiles for reaction model 4 (inset) in different solvent environments. The experimental value is given for a UpG dinucleotide (Ref. 141, 310 K and ionic strength of 1 M in NaOH/NaCl, dashed line). In order to aid visual comparison of barrier heights, the plots are shifted such that $F(\xi_R) = 0.0$. The average error bars (estimated 95% confidence interval, not shown for clarity) for all of the curves are less than 0.4 kcal/mol relative to the appropriate reactant state.

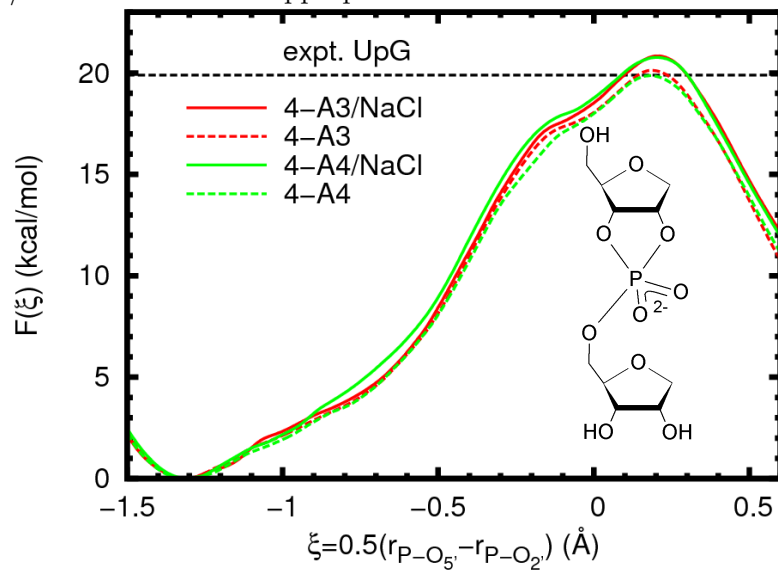


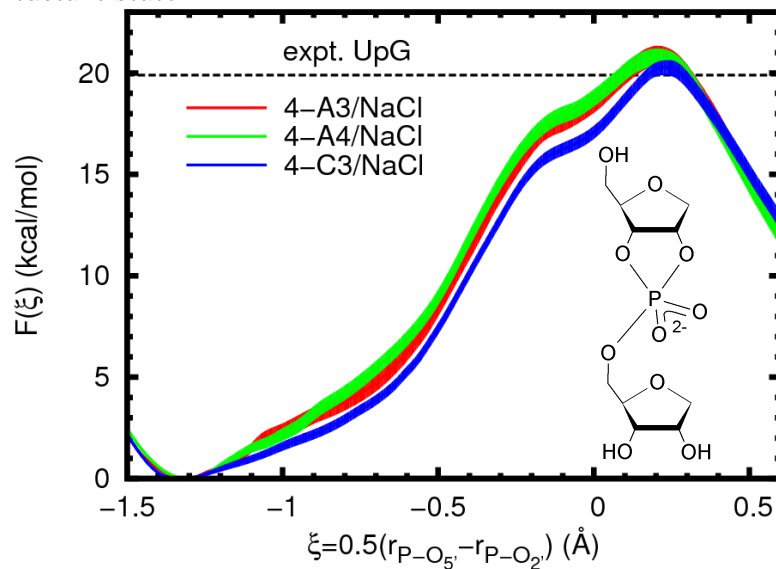
Table 4.2: Free energy profile extrema and select average geometric quantities at fixed values of the reaction coordinate for reaction model **4** in different solvent environments. For barrier heights, error bars represent approximate 95% confidence intervals, for all other quantities they represent twice the population standard deviation. Energies are in kcal/mol, lengths are in Å, and angles are in degrees.

reactant state					
model	ξ_R	ΔF^\ddagger	$\langle r_{P-O2'} \rangle_{\xi_R}$	$\langle r_{P-O5'} \rangle_{\xi_R}$	$\langle \theta_{O2'-P-O5'} \rangle_{\xi_R}$
4-A3/NaCl	-1.31	-	4.27 ± 0.11	1.65 ± 0.06	95 ± 28
4-A3	-1.31	-	4.27 ± 0.12	1.66 ± 0.06	97 ± 29
4-A4/NaCl	-1.30	-	4.26 ± 0.12	1.66 ± 0.06	94 ± 26
4-A4	-1.31	-	4.28 ± 0.10	1.66 ± 0.06	90 ± 21
transition state					
model	ξ^\ddagger	ΔF^\ddagger	$\langle r_{P-O2'} \rangle_{\xi^\ddagger}$	$\langle r_{P-O5'} \rangle_{\xi^\ddagger}$	$\langle \theta_{O2'-P-O5'} \rangle_{\xi^\ddagger}$
4-A3/NaCl	0.20	20.8 ± 0.5	1.78 ± 0.08	2.20 ± 0.22	161 ± 8
4-A3	0.18	20.1 ± 0.4	1.78 ± 0.08	2.18 ± 0.22	162 ± 9
4-A4/NaCl	0.20	20.8 ± 0.4	1.78 ± 0.09	2.19 ± 0.22	161 ± 9
4-A4	0.18	19.9 ± 0.4	1.78 ± 0.09	2.15 ± 0.18	162 ± 8

Solute-Solvent Interaction Models A necessary aspect of QM/MM simulations is to select a non-bonded, non-electrostatic interaction model for QM/MM interactions. Here, as is generally done, the choice is made from existing standard force field models. However, these models usually only aim to describe a fixed valence chemical structure and are thus not necessarily appropriate for describing chemical reactions[185]. As an investigation of the accuracy of the models used here, we compare the free energy profiles calculated with different force field models. Three solute/solvent combinations were examined using common parameters from the AMBER and CHARMM nucleic acid force fields and the TIP3P and TIP4P-Ew rigid water models.¹ A summary of

¹ For consistency with the ion parameters, the modified TIP3P model commonly used with the CHARMM force field is not employed here. Since the present purpose is simply to test whether and to what extent the models lead to different results (not to evaluate their relative quality), this potentially

Figure 4.4: Free energy profiles for reaction model 4 (inset) with different force field models in the presence of sodium chloride. The experimental value is given for a UpG dinucleotide (Ref. 141, 310 K and ionic strength of 1 M in NaOH/NaCl, dashed line). In order to aid visual comparison of barrier heights, the plots are shifted such that $F(\xi_R) = 0.0$. Filled curves represent estimated 95% confidence intervals relative to the appropriate reactant state.



these parameters is given in the Supporting Information. The free energy profiles largely display the same shape and only slight quantitative differences (Figure 4.4, Table 4.3). The change from AMBER to CHARMM Lennard-Jones parameters leads to a slight lowering of the profile between the reactant and transition states by roughly 1 kcal/mol, but without changing the reaction barrier to a statistically significant degree.

Table 4.3: Free energy profile extrema and select average geometric quantities at fixed values of the reaction coordinate for reaction model **4** with different force field models in the presence of sodium chloride. For barrier heights, error bars represent approximate 95% confidence intervals, for all other quantities they represent twice the population standard deviation. Energies are in kcal/mol, lengths are in Å, and angles are in degrees.

reactant state					
model	ξ_R	ΔF^\ddagger	$\langle r_{\text{P-O}2'} \rangle_{\xi_R}$	$\langle r_{\text{P-O}5'} \rangle_{\xi_R}$	$\langle \theta_{\text{O}2'\text{-P-O}5'} \rangle_{\xi_R}$
4-A3/NaCl	-1.31	-	4.27 ± 0.11	1.65 ± 0.06	95 ± 28
4-A4/NaCl	-1.30	-	4.26 ± 0.12	1.66 ± 0.06	94 ± 26
4-C3/NaCl	-1.32	-	4.29 ± 0.12	1.65 ± 0.06	95 ± 23
transition state					
model	ξ^\ddagger	ΔF^\ddagger	$\langle r_{\text{P-O}2'} \rangle_{\xi^\ddagger}$	$\langle r_{\text{P-O}5'} \rangle_{\xi^\ddagger}$	$\langle \theta_{\text{O}2'\text{-P-O}5'} \rangle_{\xi^\ddagger}$
4-A3/NaCl	0.20	20.8 ± 0.5	1.78 ± 0.08	2.20 ± 0.22	161 ± 8
4-A4/NaCl	0.20	20.8 ± 0.4	1.78 ± 0.09	2.19 ± 0.22	161 ± 9
4-C3/NaCl	0.23	20.3 ± 0.4	1.77 ± 0.08	2.24 ± 0.14	164 ± 10

Varying the Structural Environment of the Nucleophile and Leaving Group

In our final analysis, a series of reaction models that undergo phosphoryl transfer were established in order to systematically dissect levels of model complexity (Figure 4.1). In each case the general reaction scheme is identical to that of RNA cleavage. The nucleophile is either part of a simple alkyl chain or ribose ring and the leaving group is either ethoxide or 5'-deprotonated ribose. The calculated free energy profiles (Figure 4.5) show the barrier magnitudes clustering into three groups depending on whether or not the nucleophile is part of a ribose ring and the size of the leaving group ($\Delta F_1^\ddagger \gg \Delta F_2^\ddagger > \Delta F_3^\ddagger \gg \Delta F_4^\ddagger$, Table 4.4). A slightly different classification is obtained when comparing the location of the profile minimum (*i.e.* the reactant state); in this case the presence of a ribose ring is the most obvious factor.

Figure 4.5: Free energy profiles for several reaction models (n-A3/NaCl, where $n = 1-4$ as in Figure 4.1). Experimental values are given for a UpG dinucleotide (Ref. 141, 310 K and ionic strength of 1 M in NaOH/NaCl, dashed line, for comparison to model 4) and model 2 (Ref. 154, 353 K and 0.05 N in NaOH, dotted line). In order to aid visual comparison of barrier heights, the plots are shifted such that $F(\xi_R) = 0.0$. Filled curves represent estimated 95% confidence intervals relative to the appropriate reactant state.

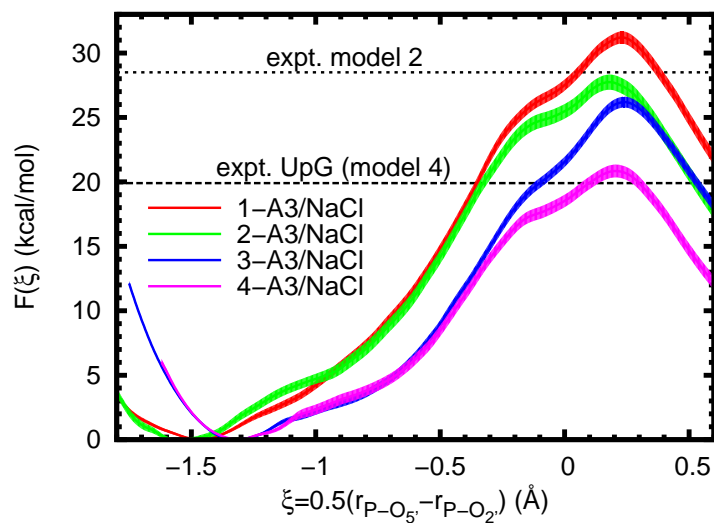


Table 4.4: Free energy profile extrema and select average geometric quantities at fixed values of the reaction coordinate for several reaction models (n-A3/NaCl, where n = 1-4 as in Figure 4.1). For barrier heights, error bars represent approximate 95% confidence intervals, for all other quantities they represent twice the population standard deviation. Energies are in kcal/mol, lengths are in Å, and angles are in degrees. ^a $\Delta G_{\text{expt}}^{\ddagger} = 28.5$ kcal/mol (Ref. 154) ^b $\Delta G_{\text{expt}}^{\ddagger} = 19.9$ kcal/mol (Ref. 141) See Figure 4.5 for details.

reactant state					
model	ξ_{R}	ΔF^{\ddagger}	$\langle r_{\text{P-O2}'} \rangle_{\xi_{\text{R}}}$	$\langle r_{\text{P-O5}'} \rangle_{\xi_{\text{R}}}$	$\langle \theta_{\text{O2}'\text{-P-O5}'} \rangle_{\xi_{\text{R}}}$
1-A3/NaCl	-1.47	-	4.58 ± 0.14	1.65 ± 0.06	130 ± 42
2-A3/NaCl	-1.55	-	4.74 ± 0.14	1.65 ± 0.06	112 ± 60
3-A3/NaCl	-1.32	-	4.29 ± 0.12	1.65 ± 0.06	100 ± 33
4-A3/NaCl	-1.31	-	4.27 ± 0.11	1.65 ± 0.06	95 ± 28
transition state					
model	ξ^{\ddagger}	ΔF^{\ddagger}	$\langle r_{\text{P-O2}'} \rangle_{\xi^{\ddagger}}$	$\langle r_{\text{P-O5}'} \rangle_{\xi^{\ddagger}}$	$\langle \theta_{\text{O2}'\text{-P-O5}'} \rangle_{\xi^{\ddagger}}$
1-A3/NaCl	0.23	31.2 ± 0.5	1.78 ± 0.09	2.22 ± 0.24	166 ± 9
2-A3/NaCl	0.18	$27.7^{\text{a}} \pm 0.6$	1.77 ± 0.08	2.20 ± 0.24	162 ± 9
3-A3/NaCl	0.24	26.2 ± 0.4	1.77 ± 0.08	2.26 ± 0.33	164 ± 10
4-A3/NaCl	0.20	$20.8^{\text{b}} \pm 0.5$	1.78 ± 0.08	2.20 ± 0.22	161 ± 8

Comparison with Experiment: UpG Dinucleotide

A primary motivation of this work was the recent publication by Harris, *et al.* of pH-rate and kinetic isotope effect data for the specific-base catalyzed cleavage of a UpG dinucleotide[141]. In that work a rate constant of $0.06 \pm 0.002 \text{ s}^{-1}$ was extrapolated at “infinite” pH near biological conditions (310 K, ionic strength of 1 M in NaOH/NaCl).²

Applying transition state theory (and standard error propagation) then gives a free energy of activation of 19.9 ± 0.02 kcal/mol. This is important because it provides optimal comparison to the constant protonation state simulations performed here. Other

² The fitting error on this quantity was not reported in the original publication.

experimental work has demonstrated small ($<5.0 \times 10^{-4} \text{ s}^{-1}$), but detectable, variations of the rate constant with respect to nucleobase sequence[153], but we do not consider these effects in the present work.

As seen in Figure 4.3 (and Figure 4.4), in all cases the agreement is very good (between 19.9 and 20.8 versus 19.9 kcal/mol), and within statistical errors. However, one must be wary that this agreement may, at least to some degree, be serendipitous. The model parameters related to solvation (water and solute Lennard-Jones parameters), which are known to influence barriers in reactions where local changes in charge state occur along the reaction coordinate as in the present example, have not been tuned for the specific reaction considered here.

The most prominent difference between the parameters of our simulations and those of experiment is the ionic strength. By necessity, experiments at high pH require high concentrations of NaOH or some other base, usually buffered with a salt[151, 153, 141]. Such conditions have been known to be problematic in periodic boundary simulations[186, 187] and it was only recently that models (employed here) were developed that robustly reproduce experimental bulk behavior[169]. However, “local,” solvation properties, such as binding coefficients[188], are often not well reproduced by many ion models, at least not on the time scales accessible in typical simulations. The modified sodium interactions used here (see Supporting Information) were designed to prevent direct binding of sodium to the phosphate in a minimally perturbative fashion so at least to enforce consistency across all umbrella sampling simulation, which might otherwise sample different bound conformations. This is justifiable based on the low binding coefficient between sodium and phosphates[188], but could be problematic at higher concentrations where the fraction of bound sodium ions is expected to be non-negligible. Li and Breaker have noted that the observed rate constant of base catalyzed RNA cleavage increased with increasing potassium concentration[153]. However, they were only able to hypothesize that this effect was primarily due to influence on the pK_a of the 2'-hydroxyl group and not on the intrinsic rate constant, as they were unable to establish simultaneously high pH and low potassium conditions. Although the system studied in that work was a DNA 22-mer with a single embedded RNA dinucleotide, the observed free energy barriers (21.5-22.5 kcal/mol, 296 K, 3.15 M K^+) are reasonably close to those measured and calculated for a simple dinucleotide. The results

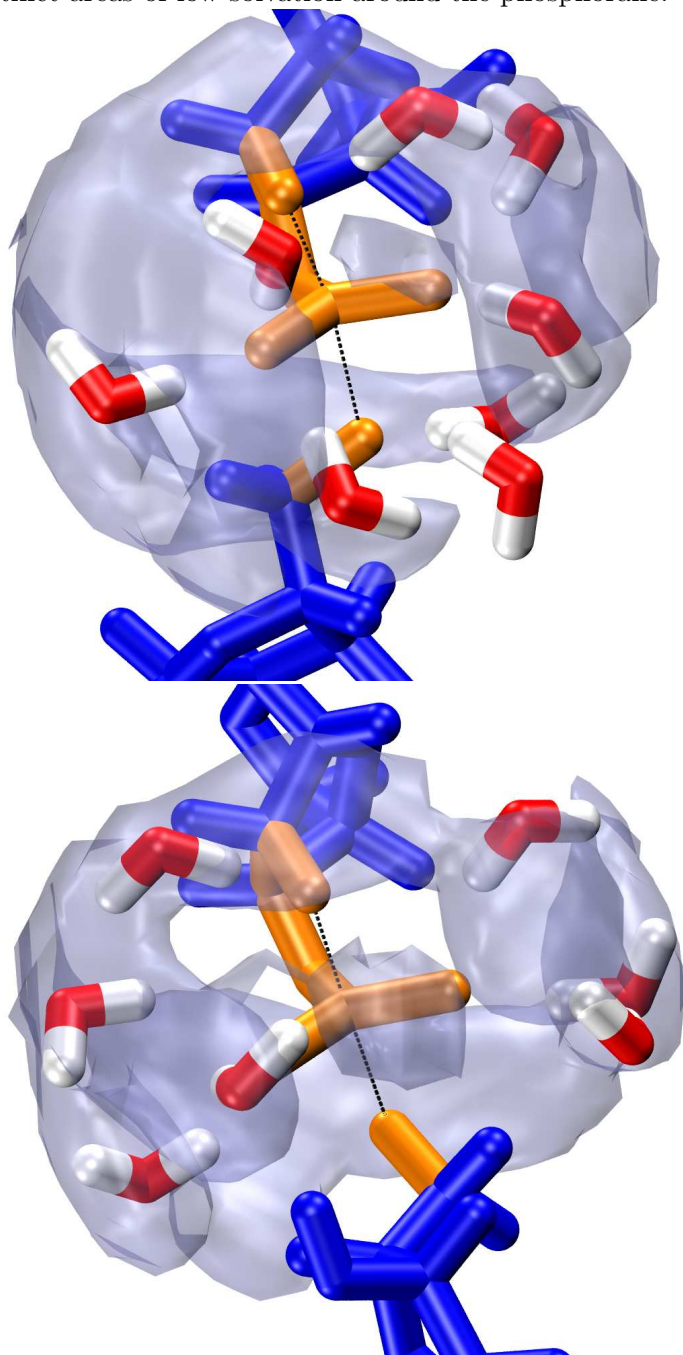
reported here are thus, at the very least, consistent with that hypothesis. Although it is tempting to suggest that the hypothesis is supported by the lack of a catalytic effect when sodium is removed, additional data would be required, for example (as suggested by a referee) examining the concentration dependence of the nucleophile pK_a or more rigorously considering sodium binding.

Free Energy Barriers and Solvent Structure

While the effects of solvation models on free energy profiles are obvious and easily compared, this does not necessarily make them an optimal metric for assessing solvation model quality. That is, there is not necessarily a direct, or even one-to-one, correspondence between empirical parameters that give the “correct” free energy barrier and those that are physically sensible. An ideal model would satisfy both criteria. As a qualitative check of the TIP3P and TIP4P-Ew water models, we examine the radial and three dimensional distribution of water molecules (or rather the water oxygens) around the phosphorane transition state (Figure 4.6). The results for both models are very similar, with distinct gaps of density around each of the non-bridge oxygen bond axes as well as parallel to the breaking and forming bonds (although a clear patch of density appears along the non-bridge oxygen angle bisector). Viewing a representative transition state-like configuration shows tetrahedral coordination of water around each of the non-bridge oxygens. As would be expected, the water molecules neatly reside on a density isosurface roughly corresponding to the first peak of the radial distribution function.

There have been some discussions in the literature concerning the choice of Lennard-Jones parameters for atoms in the QM region[189, 190, 185]. Mulholland and co-workers showed that modified parameters for nucleobases in the quantum region can improve hydrogen bond geometries with MM water molecules[191], and further demonstrated that changing between a point charge and *ab initio* electronic density for several common rigid water models can lead to unexpected (and potentially unsatisfying) results[192]. By examining extreme choices of radii and well depths, Riccardi, *et al.* demonstrated that both a reduction potential and proton transfer barrier, as well as hydrogen bonding interaction energies, display a systematic dependence on Lennard-Jones parameters[193].

Figure 4.6: Snapshots from umbrella sampling simulations near the transition state of model 4-A3 (top) and 4-A4 (bottom). Both transition state structures consist of a pentacoordinated phosphorane (orange) with advanced bond formation between the phosphorous and 2' oxygen (upper black line) and bond cleavage between the phosphorous and 5' oxygen (lower black line). This is indicative of a “late” transition state. Density maps of the water oxygens [transparent gray, isosurfaces correspond to $4\pi\rho_{\text{bulk}}g(r_{\text{max}})$] indicate three distinct areas of low solvation around the phosphorane.



However, their conclusion was that physical accuracy (in the sense of properly balancing enthalpic and entropic effects, rather than agreement with experimental data or *ab initio* calculations) would be more profitably improved through other aspects of the QM/MM model.

The present work appears to support the view that simple tuning of Lennard-Jones parameters is not a fruitful avenue to producing more physically accurate free energy profiles in QM/MM simulations. However, that is not to say that the parameters cannot non-trivially affect the profile or that the expected change cannot be predicted. In the present case, although not contrived to be so, the CHARMM potential for the highly charged non-bridge oxygens have somewhat weaker repulsive components (both Lennard-Jones ϵ and $A = \epsilon R_{\min}^{12}$ coefficients) than the related AMBER parameters (see Supporting Information). This would seem to correlate with the near systematic lowering of the CHARMM free energy profiles in comparison to the AMBER profiles (Figure 4.4). That is, with identical electrostatic interactions, the CHARMM parameters allow more preferential stabilization of the dianionic transition state relative to the monoanionic phosphate and nucleophile reactants since the weaker repulsion is overwhelmed by electrostatic attraction. It is not clear as to whether this systematic behavior is desirable. For example, a more physically correct profile might possess a steeper approach to a transition state that is lower in energy, a scenario not obviously attainable by simple modification of the Lennard-Jones potentials. This suggests the merit of a different approach. Developments in our group have sought to replace the empirical Lennard-Jones potential with a more physical model that inserts directly into the QM/MM self-consistent field calculation[194, 195]. This would allow for charge dependent exchange and dispersion interactions and would not rely on static atom type based parameters. The potential advantage would arise from non-systematic changes in the free energy profile, thereby inserting more of the physical behavior that Riccardi, *et al.* found lacking.

Dependence of the Barrier on Ground State

Amongst the free energy profiles of the four different reaction models studied here, the most obvious differences are the locations of the reactant state and height of the free energy barrier (Figure 4.5). A comparison of the average geometries at all four

reactant states (Table 4.4) shows nearly identical bond lengths between phosphorous and the O5' leaving group, but quite different distances to the O2' nucleophile. This is sensible, as the ring structure effectively prevents the secondary alkoxide from getting too far from the negatively charged phosphate, even though it would be electrostatically favorable to do so. However, when no ring is present (as in models **1** and **2**) much larger separation is possible. In this case the barrier to phosphoryl transfer is effectively increased by the added conformational change. The sizable difference in the reaction barriers of models **1** and **2** (3.5 ± 0.8 kcal/mol) is likely explained by differing solvation of the nucleophiles. That is, in principle, the additional methyl substituent on model **2** should stabilize the alkoxide more than in model **1**, thereby lowering the energy of the (deprotonated) reactant state and raising the barrier. However, this is clearly not the case. The dominating effect is the increased solvent exposure of the primary alkoxide due to the presence of fewer substituents. It therefore experiences a higher energetic cost when moving towards the transition state and thus a higher reaction barrier.

Near the transition state, all four models are strikingly similar both in the shape of the free energy profile and in the average geometry (Table 4.4). The exception to this is model **3**, which has a slightly longer bond breaking distance ($r_{\text{P-O5'}}$), as well as a "looser" bonding environment, as indicated by increased fluctuations. This variation is not enough, however, to change the classification of the transition state; the progressed bond breakage in all of the models clearly indicates a "late" transition state characterized by a nearly fully formed bond with the nucleophile, and a nearly fully broken bond to the leaving group. However, unlike the difference between the reaction barriers of models **1** and **2**, the difference between models **3** and **4** (5.4 ± 0.6 kcal/mol) seems more anomalous. The obvious departure point for examining this is the nature of the leaving groups. In the current AM1/d-PhoT QM model, the gas-phase proton affinities of these compounds are quite dissimilar ($\Delta\Delta H_{\text{PA}} = \Delta H_{\text{PA,EtOH}} - \Delta H_{\text{PA,5'-ribose}} \approx 18$ kcal/mol) and indicative of a higher cost to remove the ethoxide leaving group (if one is willing to assume the solvation properties are not too different). This rather dramatic difference also seems to be semi-quantitatively in line with high level density functional theory (DFT) calculations[196] ($\Delta\Delta H_{\text{PA}} \approx 10 - 16$ kcal/mol). However, a DFT model using continuum solvation similar to **3** was also recently shown to give a free energy barrier in very close agreement with that for an RNA dinucleotide[?], implying that

the present result is in fact anomalous, but for reasons other than the proton affinity difference. Because of the reasonable agreement of models **2** and **4** with experimental results, we suspect that there is an imbalance of solvent effects when the nucleophile is large and the leaving group is small. This is in line with the clear importance of solute/solvent and QM/MM interactions and allows for the results for model **3** to be seen as anomalous.

4.1.4 Conclusion

A significant aspect of understanding enzymatic mechanism is understanding the nature of rate enhancement over the native (solution) mechanism. In this work we have studied several variants of non-enzymatic phosphoryl transfer, a reaction catalyzed by a wide range of proteins and the most ubiquitous amongst known ribozymes. Using QM/MM MD umbrella sampling simulations, the free energy profiles were calculated along a simple atom transfer coordinate.

The calculated barrier for an abasic dinucleotide agrees almost exactly with the experimental result for the UpG dinucleotide under similar conditions. Analysis of the transition state structure indicates a “late” transition state, also in agreement with inferences from the experimental kinetic isotope effects. Although ionic conditions are a necessity of most high pH experiments, removing all monovalent ions did not lead to any significant change in the barrier. However, it should be noted that the model used here effectively precluded the possibility of direct coordination of cations to the phosphate. The obvious corollary to this result is that changes in the rate constant observed in conjunction with changes in ionic environment are likely due to the types of coordination neglected here. Such coordination is likely to be electrostatically favorable, but will come at the cost of disrupting extensive hydrogen bond networks surrounding the dianionic transition state if it occurs at the non-bridge oxygens.

The conclusions made here are strengthened by extensive testing of the available QM/MM models, including the commonly used TIP3P and TIP4P-Ew water models, as well as Lennard-Jones parameters from the AMBER FF10 and CHARMM27 force fields. Changing any one of these aspects of the model also does not lead to substantial change in the calculated reaction barrier. This does not suggest that these parameters are not important, as it has been demonstrated previously by our group and others that

free energy profiles involving local changes in charge are highly sensitive to solvation and Lennard-Jones parameters. Rather, the combinations of parameters examined here appear to be similarly balanced in terms of overall solvation.

The structural complexity of the leaving group and nucleophile has a subtle effect on the structure of the transition state, and a significant effect on the magnitude of the free energy barrier. This is most evident according to the presence or absence of a ribose ring, which acts to constrain the degree of separation between the deprotonated nucleophile and monoanionic phosphate in the reactant state. Lastly, improper balancing of model pK_a s for the nucleophile and leaving group can lead to abnormally disparate reaction barriers and should be considered when embedding QM regions into larger systems such as ribozymes. Next-generation models that are able to consider adjustments of non-electrostatic non-bonded interactions as a function of local charge are likely to considerably improve the robustness of QM/MM simulations for chemical reactions or processes that involve charge migration or change of local charge state.

4.2 Experimental and Computational Analysis of Ribonuclease A

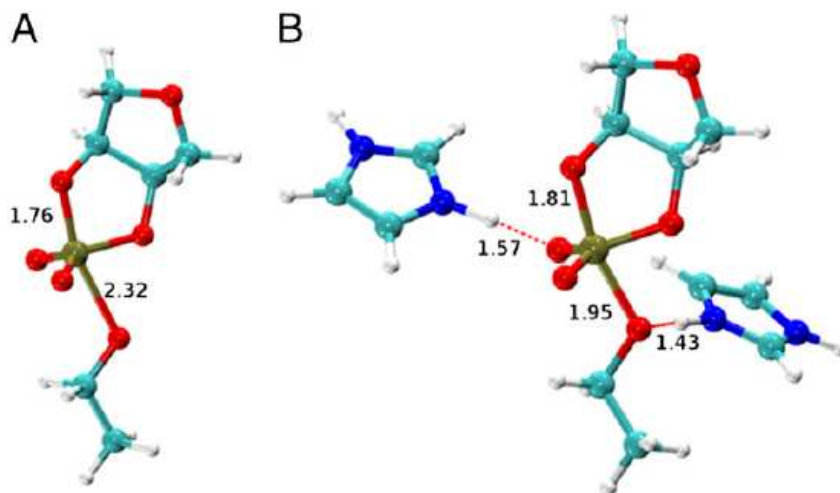
Ribonuclease A (RNase A) is one of the longest known and best studied ribonucleases (*i.e.* protein enzymes which catalyze RNA strand cleavage). As such, it is an excellent starting point for validating and calibrating new experimental and theoretical approaches. Furthermore, in the present work, it is noteworthy that RNase A can catalyze the exact same UpG RNA dinucleotide used as a substrate in the studies above as well as in the hepatitis delta virus ribozyme studied in Section 5. The difference, of course, is that RNase A provides catalysis through a protein environment while the other two situations utilize solution and RNA environments. A detailed analysis of RNase A can be found in a recent collaborative publication[6], the main focus here will be on the aspects relevant to non-enzymatic and ribozymatic reactions and the methods used here to study them.

Table 4.5: Heavy atom KIEs for the leaving group (LG), nucleophile (NUC), and non-bridge phosphate oxygens (NPO) of UpG during catalysis by both solution (acid or base) and RNase A. The reported values are measured unless indicated to be from quantum mechanical (QM) calculations on a model RNA compound. All values are as reported in Ref. 6 and references therein. Values in parentheses are standard deviations.

Catalyst	$^{18}k_{\text{LG}}$	$^{18}k_{\text{NUC}}$	$^{18}k_{\text{NPO}}$
RNase A, pH 7	1.014(3)	0.944(2)	1.001(1)
RNase A (QM)	1.026	0.998	1.006
Acid, pH 0	1.005(4)	0.990(4)	0.991(1)
Base, pH 12	1.037(2)	0.996(2)	0.999(1)
Base, pH 14	1.034(4)	0.984(3)	-
Base, pH 14 (QM)	1.046	0.973	1.002

Several KIE values, both measured and from quantum mechanical calculations, are presented in Table 4.5. Reflecting on the previous section, it will be noticed that the solution catalyzed reactions in high and low pH give different KIEs, especially for the leaving group ($^{18}k_{\text{LG}}$). The most straightforward interpretation of this observation is that bond scission with the leaving group is less pronounced in the acid catalyzed reaction, perhaps because additional proton transfers not present in the base catalyzed reaction stabilize the transition state (or an intermediate)[197]. In envisioning how this might be so, it is useful to consider the transition state structures identified via simulation in the previous section for the base catalyzed reaction. The leaving group bond lengths in these structures are considerably longer than those with the nucleophile, likely due to the fact that the structure is a dianionic phosphorane and the only available solvent stabilization is from water (monovalent ions appear not to play a strong role[153, 1]). It is considerably striking then that the KIE values for the RNase A catalyzed reaction are nearer to those of the base catalyzed reaction than to those of the acid catalyzed reaction.

Figure 4.7: Transition state structures for a model reaction for non-enzymatic RNA cleavage (A) and enzymatic cleavage by ribonuclease A (B). Reproduced (with permission) from Ref. 6.



In order to construct a structural model that coincides with the above KIE measurements, simplified atomic models were suggested and the KIEs were calculated (Figure 4.7). A description of these methods is outside the scope of this work, see, for example, Ref. 28. As can be seen in Table 4.5 the calculated and measured values for the base catalyzed reaction are in good agreement. Furthermore, this simplified model is in agreement with both other experimental data[141, 142] and the simulations in the previous section utilizing a more complicated model. This validation gives rise to considerable confidence in producing a simplified model of the protein enzyme environment provided by RNase A. Indeed, a minimal model including only modest parts of the protein environment gives rise a KIE signature in good agreement with the KIEs for the RNase A catalyzed reaction (“RNase A (QM)” in Table 4.5).

Comparison of the two structural models used in the KIE computations above offer a description and interpretation of why the experimental values differ. In particular, much like the notional interpretation of the acid versus base catalyzed reactions, the leaving group KIE of the RNase A catalyzed reaction is smaller than that for the base catalyzed reaction and corresponds to not only a shorter leaving group/phosphorous bond (1.95 versus 2.32 Å, respectively), but a shorter nucleophile/phosphorous bond

as well (Figure 4.7). These bonds are likely shorter due to the presence of hydrogen bonds with the simplified enzyme environment and provide a clear indication that the addition of these groups preferentially stabilizes the transition state with respect to that in solution. Clearly, it is expected that ribozymes make use of similar strategies when performing catalysis and KIEs and computational models can be useful in identifying the chemical groups involved.

Chapter 5

Computational Investigations of Phosphoryl Transfer in Ribozymes

The history, importance, and potential applications of ribozymes have already been discussed in Section 2.1. In this section several computational studies of ribozyme mechanism, particular that of the hepatitis delta virus ribozyme (HDVr) will be presented. The HDVr is one of several well known “small” ribozymes originally found in viruses. Such ribozymes are considerably easier to study, both experimentally and computationally, compared to their larger counterparts found in eukaryotic biochemistry. This is because they are more amenable to mutation and modification (due to their size) and frequently physically robust in laboratory conditions (perhaps due to the harsh conditions viruses frequently experience). Furthermore, it is quite clear that these ribozymes span a broad mechanistic range and are thus valuable archetypes for the mechanism of larger ribozymes[198, 199, 200, 201, 202, 203]. Indeed, several mechanistic insights regarding catalytic strategies such as site specific shifts of nucleobase pK_a values and recruitment of divalent metal ions to stabilize electrostatically strained structures have emerged from both experimental and computational analyses of these ribozymes[204, 205, 206, 207, 208, 209]. This trend is continued here by a recent computational study of the mechanism of the HDVr[7]. In Section 5.2 this is reproduced

with permission from the Journal of the American Chemical Society (submitted, the supporting information for this is also included as Appendix B).

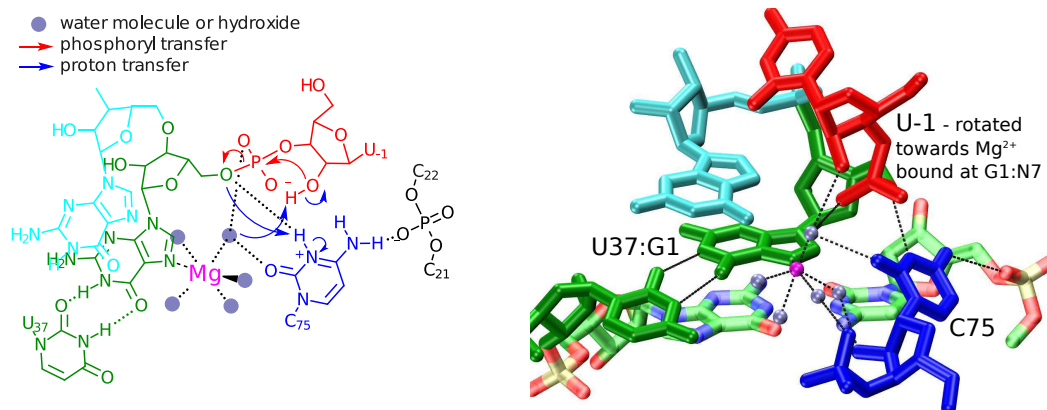
5.1 The Hepatitis Delta Virus Ribozyme

5.1.1 Background

The hepatitis delta virus ribozyme is a small, self-cleaving ribozyme found as a satellite of the hepatitis B virus[144, 210]. Of the few known “small,” viral ribozymes, it is the only one found in an animal virus[211] and is of particular interest due to the existence of similar sequences in both the human genome[18] and the genome of many other eukaryotes[14, 212]. Like most of the other known ribozymes, there has been considerable progress in characterizing the structure and mechanism of the HDVr. However, while the two main catalytic strategies in the HDVr appear to be analogous to those found in other ribozymes (*e.g.* the hammerhead and hairpin ribozymes), they are also apparently coupled in an unexpected fashion; this may be the origin of the HDVr’s superior catalytic rate.

Both the genomic (γ) and anti-genomic (α) strands of the HDVr are catalytically competent[144, 210, 213] and it is generally agreed that, in both cases, a minimum sequence of ~ 85 nt is necessary to support self-cleavage, with only 1 nt 5’ of the cleavage site between U-1/G1 (C-1/G1 in the α strand)[210, 214]. Numerous biochemical kinetic characterizations indicate that an optimal rate (a first order rate constant of ~ 1 - 10 min^{-1}) and efficiency are obtained at modest pH (~ 8)[215, 216] and high concentrations of Mg^{2+} (~ 50 - 100 mM)[217, 218, 219]. Autocatalysis is also known to occur in the absence of Mg^{2+} with molar amounts of NaCl[220, 221] or with other divalent metal ions[222, 223]. Mutation (or deletion) of C75 (C76 in the α strand) to any other standard residue greatly disrupts catalytic activity[219, 162]. Several non-standard uracil and cytosine analogue substitutions also diminish activity to varying degrees[224, 225, 226]. Inactive C75 Δ and C75U mutants have been rescued by imidazole titration[227] and C75U and C75n⁶C (6-azacytosine) mutants by site specific substitution of G1:O5’ with sulfur, an enhanced leaving group[226]. Interestingly, low concentrations of $\text{Co}(\text{NH}_3)_6^{3+}$ inhibit wild-type activity and abolish imidazole rescue of C75U mutants[228, 229]. These results strongly suggest that C75 acts as a general acid, protonating the G1:O5’ alkoxide

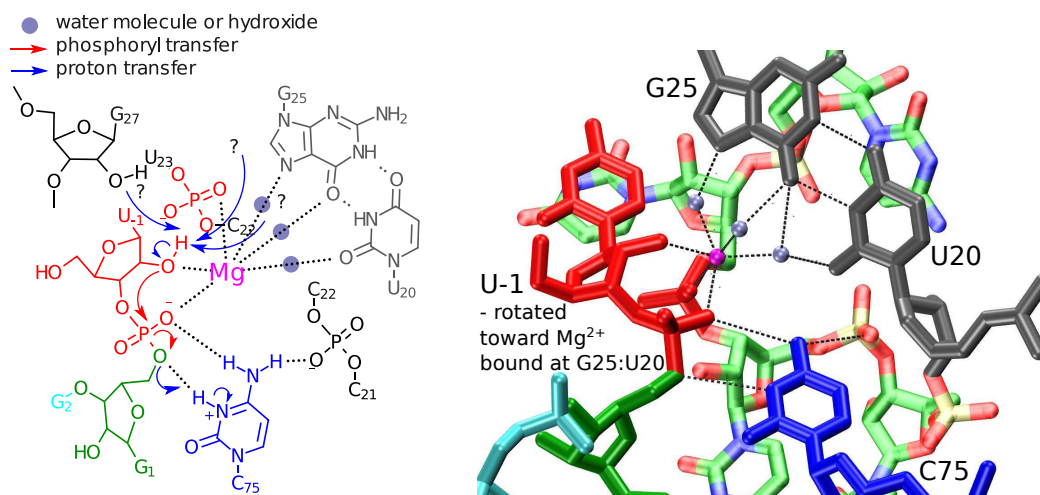
Figure 5.1: Proposed mechanism and schematic model (reproduced independently) from Ref. 229. In this model, a water or hydroxide molecule bound to a Mg^{2+} ion (magenta) coordinated to G1:N7 (part of a G1:U37 standard wobble, green) acts as a general base, activating the U-1:O2' nucleophile (red). The O5' primary alkoxide formed after phosphoryl transfer then accepts a proton from (protonated) C75:N3 (blue).



formed by 2'-*O*-transesterification. This hypothesis is further supported by Raman crystallographic measurements, which, by associating a resonance with C75, show that C75 has a microscopic pK_a in line with biochemical results[230]. Unfortunately, the role of Mg^{2+} , while clearly important, is not clarified by these results alone.

The catalytic (as opposed to structural) importance of Mg^{2+} has been supported by several studies[216, 220, 226, 228, 229]. A recent crystallographic study resolved an ion near the active site of an inhibited dU-1 mutant (U-1 itself could not be unambiguously resolved)[231]. This ion was present at a different location in the original crystal structure of a cleaved ribozyme[232]. However, other structures of C75U mutants crystallized with different divalent metal ions have resolved them in similar positions[228, 233, 234]. Interestingly, both ion positions are near GU wobble pairs. One pair (G1:U37, G1:U39 in the α strand) is a *cis* Watson-Crick/Watson-Crick pair (*c*WW, also called a “wobble” or “standard wobble” pair), and appears to be conserved as a purine:pyrimidine pair across many HDV-like ribozymes[212]. The other (G25:U20, G28:U23 in the α strand) has been resolved in both *cis* Hoogsteen/Watson-Crick (*c*HW)[228, 233] and *trans* Watson-Crick/Watson-Crick configurations (*t*WW, also called a “reverse wobble”)[231]. Moreover, this pair is rigorously conserved, along

Figure 5.2: Proposed mechanism and schematic model from Ref. 231. In this model, a Mg^{2+} ion (magenta) is coordinated to U-1:O2', U-1:*pro*-R_P, and U23:*pro*-S_P (red) and has three water mediated contacts to a G25:U20 reverse wobble (gray). This interaction is suggested to activate a general base, possibly, but not exclusively, a water molecule or G27:O2'. The O5' primary alkoxide formed after phosphoryl transfer then accepts a proton from (protonated) C75:N3 (blue).



with an adjacently stacked cytosine and other parts of the L3 strand[14, 212]. The *cHW* form has been found exclusively in the crystals of inactive C75U mutants, while the *tWW* form has been observed in crystals of dU-1 mutants (a similar configuration with quite elongated hydrogen bonds was also observed in the cleaved structure[232]). The significance of the reverse wobble configuration has been supported by both experiments on and simulations of a G25A:U20C double mutant[235, 236, 237] as well as all other isosteric pairs[238]. Recent experiments with chemical probes have also suggested the possibility for interconversion between the two configurations before and after cleavage[238], a scenario in line with the observation that the L3 strand is quite flexible[239] and potentially prone to misfolding[240, 241, 242]. Finally, Raman crystallography measurements have also suggested that the ~ 2 unit shift in the pK_a of C75 may correlate negatively with the presence of Mg^{2+} [230]. Unfortunately, none of these experiments directly suggest which Mg^{2+} ion (or ions) gives rise to this effect.

Several mechanisms have been proposed that attempt to integrate the experimental results just described, all building on the notion that the HDVr possesses two distinct

catalytic strategies[234]: 1.) use of an environmentally induced shift in the pK_a of C75 (likely acting as a general acid) and 2.) recruitment of one or more Mg^{2+} ions for electrostatic stabilization and/or activation of proton donors/acceptors (likely an acceptor acting as a general base). On the basis of Raman crystallography and competition experiments with $Co(NH_3)_6^{3+}$ [243], as well as site specific mutation of two guanine residues, Chen, *et al.* suggested a model in which a Mg^{2+} ion is inner sphere bound at G1:N7 and acts a general base, activating the U-1:O2' nucleophile; C75 would then act as a general acid, stabilizing the O5' alkoxide (Figure 5.1)[229]. This model was based on the superposition of a pre- and post-cleavage crystal structure. However, when a new pre-cleavage crystal was reported with a Mg^{2+} ion bound near U23 and the G25:U20 reverse wobble, another mechanistic model was proposed in which *this* ion facilitated the general base step (Figure 5.2)[231]. Importantly, this model was based on analogy to a crystal structure of the hammerhead ribozyme, as U-1 is inherently disordered in the electron density maps from pre-cleavage crystals lacking a C75U mutation[228]. Molecular dynamics (MD) simulations departing from this second model have characterized it as exceedingly stable[235, 236, 237]. However, they have been extremely short (~ 25 ns). By comparison, we have recently reported comparatively long MD trajectories (~ 350 ns) departing from an appropriately mutated C75U crystal structure[244]. These simulations show significant flexibility when either one or both Mg^{2+} ion binding sites are occupied. On the basis of these and subsequent simulations, we have suggested an alternate mechanistic model which is intermediate to those in the literature[?]. Namely, we suggest a transition between the two conformations, in line with the observed disorder in the crystal structures and in agreement with all of the experimental evidence used to support the existing models.

The existing body of experimental work has probed the mechanism of the HDVr via *indirect* means, namely perturbation by varying the pH or ion concentrations or by mutations. An extremely powerful technique for *direct* probing is the measurement of kinetic isotope effects (KIEs). In collaboration with Michael Harris and co-workers, we have shown how an integrated experimental/theoretical approach can be extremely valuable in analyzing the mechanisms of both model compounds[142] and a paradigmatic ribonuclease[6]. This approach has never before been applied to a ribozyme system and would represent a significant advance in the analysis of ribozyme mechanisms. Our new

simulation models and approaches offer a powerful platform with which to interpret these results and provide predictive mechanistic insights.

5.2 A Framework for Assessment of Metal-Assisted Nucleophile Activation in the Hepatitis Delta Virus Ribozyme

The hepatitis delta virus ribozyme is an efficient catalyst of RNA 2'-*O*-transphosphorylation and has emerged as a key experimental system for identifying and characterizing fundamental features of RNA catalysis. Recent structural and biochemical data have led to a proposed mechanistic model whereby an active site Mg^{2+} ion facilitates deprotonation of the O2' nucleophile and a protonated cytosine residue (C75) acts as an acid to donate a proton to the O5' leaving group (*Biochemistry*, **2010**, *49*, 6508-6518). This model assumes that the active site Mg^{2+} ion forms an inner-sphere coordination with the O2' nucleophile and a non-bridging oxygen of the scissile phosphate. These contacts, however, are not fully resolved in the crystal structure, and biochemical data are not able to unambiguously exclude other mechanistic models. In order to explore the feasibility of this model, we exhaustively mapped the free energy surfaces with different active site ion occupancies via quantum mechanical/molecular mechanical (QM/MM) simulations. Further, we incorporate a three-dimensional reference interaction site model for the solvated ion atmosphere that provides a realistic estimate of how the hypothetical ribozyme activated states are populated, and the degree to which active site Mg^{2+} ion binding would shift the nucleophile $\text{p}K_{\text{a}}$. The QM/MM results are in alignment with the available experimental data and suggest that, under the assumption of the metal ion binding mode described above, a pathway involving metal-assisted nucleophile activation is feasible and favorable over one in which the metal is absent. The simulation results are analyzed in the context of the existing experimental and computational data and key mechanistic predictions regarding transition state bonding are highlighted.

5.2.1 Introduction

The hepatitis delta virus ribozyme (HDVr) is a small, self-cleaving ribozyme found in the genome of a satellite of the hepatitis B virus[144, 210]. Of the few known viral ribozymes, it is the only one found in an animal virus[211] and is of particular interest due to the existence of similar sequences in both the human genome[18] and the genome of many other eukaryotes[14, 212]. It is now generally accepted that ribozymes like the HDVr employ a variety of catalytic strategies, including site specific shifts of nucleobase pK_a s and/or recruitment of divalent metal ions to stabilize electrostatically strained structures[198, 199, 202, 203]. These motifs are well established from detailed experimental and theoretical analysis of the hairpin and hammerhead ribozymes[204, 205, 206, 207, 208, 209].

Indeed, experimental studies of the HDVr have identified a specific cytosine residue, C75, as being critical for catalysis[219, 232, 227, 216, 245] and a broad range of evidence supports a scenario in which the pK_a of C75 is shifted ~ 2 units towards neutrality compared to both a single nucleotide in solution[230] and the cleaved product state[246]. Furthermore, biochemical data support the supposition that this residue donates a proton to the leaving group, thus acting as an acid catalyst[226, 247].

Metal ions also contribute to HDVr catalysis, in addition to their role in RNA folding. The HDVr requires millimolar concentrations of Mg^{2+} ions[213, 217, 218] (or some other divalent ion[222]) in order to reach an optimal reaction rate under near-physiological conditions. Detailed biochemical and kinetic studies revealed multiple functional divalent metal binding sites and demonstrated that molar concentrations of monovalent ions alone can support catalysis[220, 248, 221]. Site-bound metal ion interactions have been identified and characterized via crystallography[228, 233, 231], spectroscopy[243], chemical probing experiments[229, 238, 237], and molecular simulation[244, 235, 236, 237]. Moreover, pH-rate profiles for the reaction in the absence of Mg^{2+} and for mutants designed to disrupt binding of the proposed active site ion are inverted relative to the reaction of the native HDVr in Mg^{2+} [221, 237]. Phosphorothioate interference studies also revealed sites of potential site-bound metal ion interactions via coordination to one or more non-bridging oxygens, including the *pro-R_P* position of the scissile phosphate[249, 223, 250, 251, 225, 226, 221, 252]. Thiophilic metal ion rescue experiments also support a catalytic metal ion interacting with this position. However, unlike

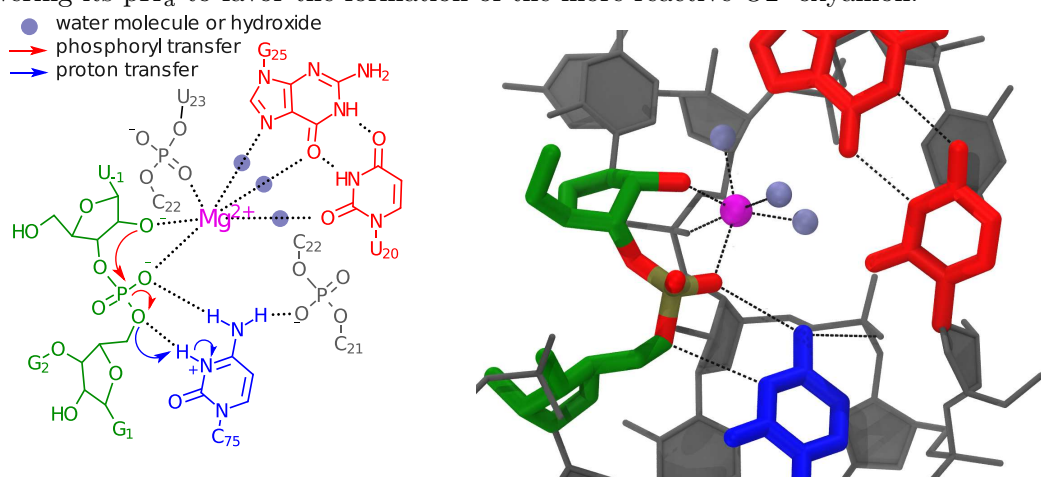
results for other metalloribozymes[253], the metal rescue is unconventional in that the substitution decreases the susceptibility of the ribozyme to inhibition by thiophilic metal ions rather than being activated by their presence[252]. Thus, evidence for the participation of an active site metal is very strong, but the details of its interactions in the transition state and its role in catalysis are not yet clear.

As with all mechanistic studies, a major difficulty lies in determining exactly what atomistic model best fits the data. Computational simulation allows for direct exploration of specific mechanistic pathways as well as identification of experimental observables that can potentially be used to discriminate between them. In the case of the HDVr, molecular simulations can be used to specifically position Mg^{2+} at locations thought to facilitate catalysis. The results of these simulations can then be used to evaluate: 1) the thermodynamics of ion association, 2) the catalytic competency of the resulting configurations, and 3) the structural and dynamical characteristics of mechanistic pathways and how they might affect experimental observables.

Recent *ab initio* quantum mechanical/molecular mechanical (QM/MM) studies, both via adiabatic calculations[254] and the string method[255], have focused primarily on analyzing the details of a specific set of active site metal interactions proposed by Chen, *et al.*[231] and Golden[234] (Figure 5.3). This mechanistic model was inferred from both crystallographic data and structural modeling via homology to the hammerhead ribozyme and is consistent with the available biochemical data[234]. The calculations have provided a detailed catalytic pathway consistent with a wide range of experimental data as well as a framework for evaluating predicted experimental outcomes. Significantly, these simulations predicted a change from a concerted mechanism when the active site Mg^{2+} is present, to a stepwise pathway via a protonated phosphorane intermediate in the presence of Na^+ (and absence of Mg^{2+}). However, the quantitative agreement of these simulations with predicted activation energies remains non-optimal, in part due to limitations arising from treatment of ion binding events based on NMR and kinetic measurements of the O2' $\text{p}K_a$ [255].

The work presented here provides a detailed computational perspective on active site Mg^{2+} association and its contribution to catalysis by examining the distribution of states where divalent metal ion binding and nucleophile activation are thermodynamically connected within a molecular mechanics/three-dimensional reference interaction

Figure 5.3: Schematic diagram of the catalytic mechanism proposed by Chen, *et al.*[231, 234] Biochemical, structural and computational data together indicate that a protonated cytosine residue (blue) acts as an acid to transfer a proton to the O5' leaving group (blue arrows). Notably, the model shows that a hexacoordinated Mg^{2+} ion (magenta) with three water mediated contacts (gray spheres) to a GU reverse wobble (red) as well as a hydrogen bond between cytosine and the scissile phosphate aid in organizing the active site. The three remaining Mg^{2+} ligands are phosphate oxygens, including the O2' nucleophile. This interaction is proposed to play a role in activating the nucleophile by lowering its pK_a to favor the formation of the more reactive O2' oxyanion.



site model (MM/3D-RISM) framework. This integrated approach decomposes the pre-equilibrium and reactive steps so as to establish a common reference in which to compare free energy profiles (and barriers) calculated from extensive QM/MM simulations both with and without an active site-bound Mg^{2+} ion. The results are in close agreement with experimental rates and suggest that, under the assumption that the previously proposed metal ion binding mode is populated in the ground state[231, 234], a feasible mechanism involves nucleophile activation facilitated by metal ion coordination. This is then followed by cleavage of the leaving group bond facilitated by C75 acting as a general acid catalyst.

5.2.2 Computational Methods

Molecular Dynamics

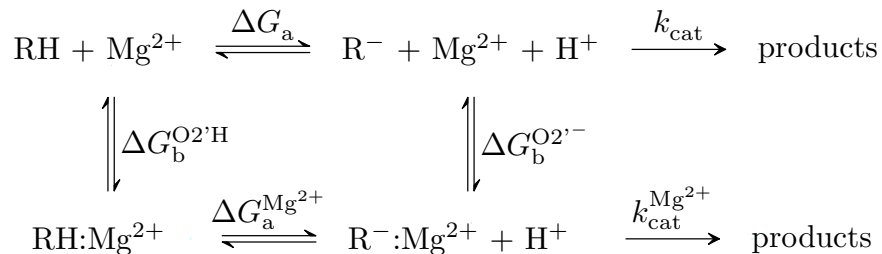
MM and QM/MM MD simulations were performed using the AMBER 14[136, 256] suite of programs. Atoms in the MM region were treated with the AMBER FF10 force field[109, 110, 111, 112] while those in the QM region were described by the AM1/d-PhoT semi-empirical Hamiltonian[116]. The solvent environment was modeled using the TIP4P-Ew[138] rigid water model and the associated alkali metal and halide ion parameters of Joung and Cheatham[169]; magnesium was modeled based on calculations by Mayaan, *et al.*[257]. Long range electrostatics were treated using periodic boundary conditions and the particle mesh Ewald method[119, 121]. Full details of the simulation protocol and system setup are given in the supporting information.

Long MD trajectories (>900 ns in total) were propagated from multiple initial states in which U-1:O2' was either neutral or deprotonated and a Mg^{2+} ion was or was not specifically bound at the active site. C75:N3 and C41:N3 were always protonated, as in previous works[235, 236, 244, 237]. In the case that Mg^{2+} was not bound, the ion was replaced with a Na^+ ion from the bulk and the system was re-equilibrated, providing a gradual transition towards the new ion environment. Trajectories were structurally analyzed for use in 3D-RISM and Non-Linear Poisson-Boltzmann/Surface Area (NLPB/SA) calculations and were also used as starting structures for QM/MM trajectories.

3D-RISM and NLPB Calculations

3D-RISM and NLPB/SA calculations were performed with AMBER[47] and its interface to the Adaptive Poisson-Boltzmann Solver[258, 259]. Multiple RISM closures, including the Kovalenko-Hirata (KH)[42] and n th order partial series expansion (PSE_n , $n=2,3$)[45] closures, were compared. The specific numerical details of both 3D-RISM and NLPB calculations are given in the supporting information. Calculations were performed on structures derived from long MD simulations by removing all solvent atoms (except bound Mg^{2+} , when appropriate) and, when applicable, changing the charge vector to that of the deprotonated nucleophile.

Figure 5.4: Reaction scheme for (de)protonation of the HDVr at U-1:O2' and Mg²⁺ binding events needed to attain catalytically active states with and without Mg²⁺ ($k_{\text{cat}}^{\text{Mg}^{2+}}$ and k_{cat} , respectively). RH indicates a (neutral) protonated reactant state (O2'H) while R⁻ indicates a deprotonated reactant state (O2'⁻). In all cases, C75 is assumed to be in a protonated state.



Relative free energies were derived from a thermodynamic cycle (Figure 5.4) describing protonation (ΔG_{a}) and binding of Mg²⁺ at U-1:O2' (ΔG_{b}). From these definitions, a p*K*_a shift, $\Delta\text{p}K_{\text{a}}$, for deprotonation of the O2' nucleophile between the Mg²⁺ bound and unbound active site configurations, can be defined in proportion to the difference in either the free energies of deprotonation or the free energies of Mg²⁺ binding.

$$\begin{aligned}
 \Delta\Delta G_{\text{a}} &\equiv \Delta G_{\text{a}}^{\text{Mg}^{2+}} - \Delta G_{\text{a}} \\
 &= \Delta G_{\text{b}}^{\text{O2'^{-}}} - \Delta G_{\text{b}}^{\text{O2'H}} \\
 \Delta\text{p}K_{\text{a}} &= \frac{\Delta\Delta G_{\text{a}}}{RT \ln 10}
 \end{aligned} \tag{5.1}$$

The components of the relative free energies are then estimated in the usual way from the solvation free energies[260]. Since both 3D-RISM and NLPB/SA provide robust estimates for the solvation free energy of Mg²⁺, they can directly calculate the binding free energies $\Delta G_{\text{b}}^{\text{O2'H}}$ and $\Delta G_{\text{b}}^{\text{O2'^{-}}$.

QM/MM Hamiltonian Replica Exchange

QM/MM Hamiltonian replica exchange umbrella sampling simulations were performed using a novel asynchronous protocol[2] with exchanges attempted at 5 ps intervals. Sampling in each replica averaged 60 - 100 ps (235 ns total). The replica states (>1300 per free energy surface) were defined by harmonic restraints, $U(\xi)$, on two atom transfer

coordinates:

$$U(\xi) = k(\xi - \xi_0)^2 \quad \xi \equiv r_{X-Y} - r_{Y-Z}. \quad (5.2)$$

Here r_{X-Y} indicates the distance between atoms X and Y. For phosphoryl transfer, ξ_{PhoT} , $X=\text{G1:O5}'$, $Y=\text{G1:P}$, $Z=\text{U-1:O2}'$, and $k = 50 \text{ kcal/mol-}\text{\AA}^2$ and for proton transfer of the general acid, $\xi_{\text{ProT}}^{\text{GA}}$, $X=\text{G:O5}'$, $Y=\text{C75:H3}$, $Z=\text{C75:N3}$, and $k = 60 \text{ kcal/mol-}\text{\AA}^2$ (C75:N3 is assumed to be protonated). The restraint locations, ξ_0 , were never more than 0.15 \AA apart. The resulting data was analyzed by the multistate Bennet acceptance ratio (MBAR)[71] in tandem with a Gaussian kernel density estimator (see Ref. 1 for details) as well as the recently developed variational free energy profile (vFEP) method[5]. The resulting free energy surfaces were then analyzed using routines from the DL-FIND library[184] in order to identify stationary points and minimum free energy paths.

5.2.3 Results

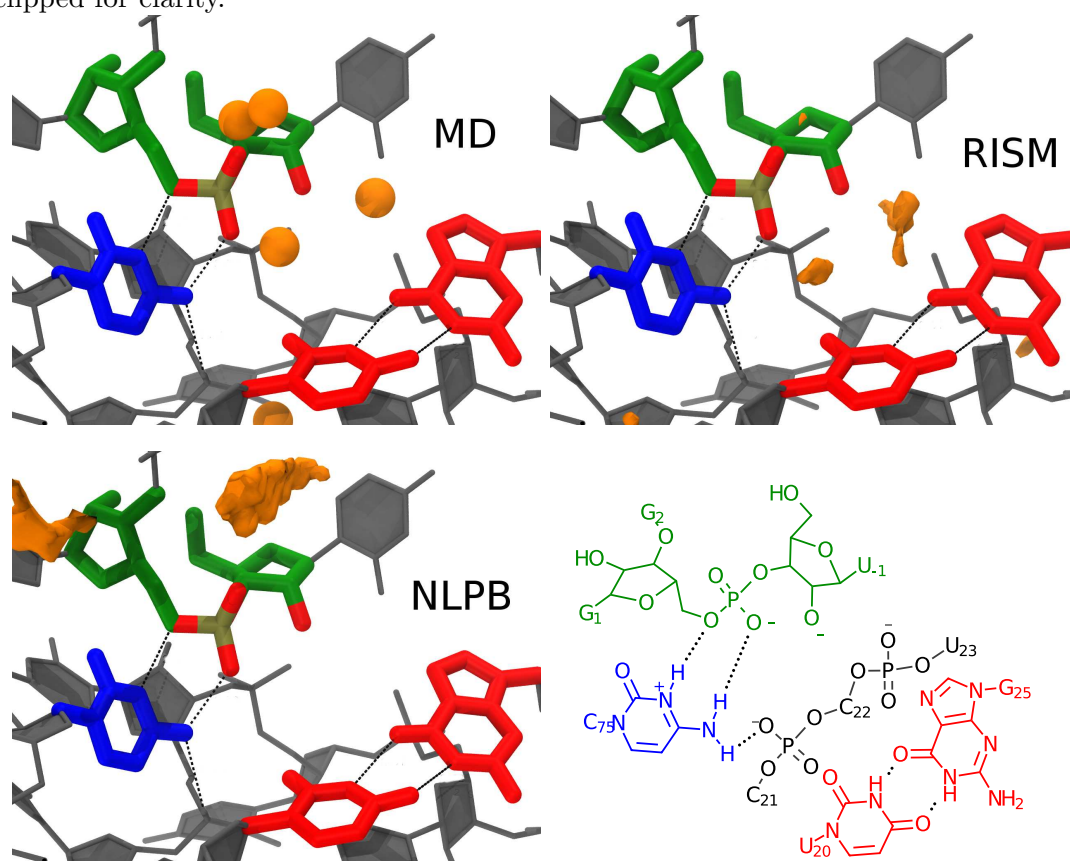
Simulations of the HDVr Ground States

First, multiple HDVr states (RH, R^- , RH:Mg^{2+} , and $\text{R}^-:\text{Mg}^{2+}$) were defined based on the protonation state of U-1:O2' and whether or not Mg^{2+} was specifically bound as in the crystallographic model[231, 234] (see Figure ??). As a baseline exploration of the conformational space available to each state, long-time MD trajectories were propagated for at least 100 ns of data collection. In all states, the trajectories displayed remarkable stability in an active, inline conformation of the U-1:O2' nucleophile. This was true regardless of whether or not a crystallographically resolved Mg^{2+} ion was artificially ejected from the active site (by swapping coordinates with a bulk Na^+ ion, as in the RH and R^- states) or the nucleophile was protonated or deprotonated (as in the R^- and $\text{R}^-:\text{Mg}^{2+}$ states). In order to evaluate these conformational searches, cluster analysis was performed and yielded single dominant clusters (>80% occupancy) with small fluctuations of the active site heavy atoms, although this was less true in the case of a deprotonated nucleophile (see supporting information).

3D-RISM and NLPB/SA Analysis

The ultimate goal of MM/NLPBSA and MM/3D-RISM type calculations is to assess energetics and solvation effects; however, these methods make significantly different

Figure 5.5: Na^+ pair distribution functions (orange isosurfaces) computed via NLPB (bottom left) and 3D-RISM-PSE3 (top right) compared with peak positions (orange spheres) from MD and volmap (top left). Isosurfaces correspond to a concentration of 300 times the bulk (140 mM). Density beyond 3 Å from the active site residues was clipped for clarity.



assumptions concerning the structure of the solvent environment (*e.g.* standard NLPB neglects ion-ion correlations, whereas 3D-RISM does not). Thus, as an initial test of the quality of 3D-RISM and NLPB calculations, the HDVr-Na⁺ pair distribution functions (PDFs) were calculated on a structure from a trajectory in which a Mg²⁺ ion was *not* bound (Figure 5.5). These PDFs describe, in an average sense, the three-dimensional distribution of Na⁺ ions around the HDVr solute without the effects of bound Mg²⁺. The ion density distribution predicted by 3D-RISM-PSE3 matches closely that from MD simulation which places a specific, buried Na⁺ ion near the nucleophile (Figure 5.5). NLPB, on the other hand, does not predict this density or even qualitatively agree with either the 3D-RISM or MD simulation ion density.

Table 5.1: MM/3D-RISM and MM/NLPB/SA results for binding free energies and p*K*_a shifts under varying background concentrations of NaCl (see Eqn. 5.1 and Figure 5.4 for definitions).

		[NaCl] (mM)		
method		140	200	1000
$\Delta G_b^{\text{O}2'\text{H}}$	RISM-KH	-11.7	-9.4	-1.2
	RISM-PSE2	-11.1	-9.0	-0.2
	NLPB/SA	-105.5	-105.0	-103.0
$\Delta G_b^{\text{O}2'^-}$	RISM-KH	-19.4	-17.0	-8.3
	RISM-PSE2	-16.4	-14.2	-4.7
	NLPB/SA	-183.1	-182.6	-180.5
ΔpK_a	RISM-KH	-5.6	-5.5	-5.1
	RISM-PSE2	-3.9	-3.8	-3.3
	NLPB/SA	-56.5	-56.5	-56.5

As a next step, Mg²⁺ binding free energies and p*K*_a shifts were calculated with varying background monovalent salt concentrations (Table 5.1). 3D-RISM calculations with the KH and PSE2 closures estimate binding free energies between -20 and -5 kcal/mol. NLPB/SA predicts very large binding free energies on the order of -100 kcal/mol. Omitting surface area terms decreased the magnitude of these by ~ 2 kcal/mol.

All methods predict lower binding affinities at higher background salt concentration, although the trend is much more pronounced with 3D-RISM than with NLPB/SA.

As described in Section 5.2.2 and Figure 5.4, the relative Mg^{2+} binding free energies in the O2' protonated and deprotonated states can be related to a $\text{p}K_{\text{a}}$ shift of the O2'-hydroxyl that occurs upon Mg^{2+} binding. The 3D-RISM results are much more consistent for this quantity, on average predicting a downward shift of 3.5-5.5 units. NLPB/SA predicts an extremely large shift of -56.5 units, regardless of the inclusion of surface area terms. Although 3D-RISM predicts slightly smaller shifts at higher salt concentrations, NLPB/SA predicts no such trend.

QM/MM Free Energy Surfaces

QM/MM free energy surfaces of the HDVr catalyzed reaction were calculated both with ($\text{R}^-:\text{Mg}^{2+}$) and without (R^-) Mg^{2+} present in the active site (see Figure 5.4). Figure 5.6 shows the free energy surfaces defined by axes corresponding to phosphoryl transfer (ξ_{PhoT}) and general acid proton transfer from C75 ($\xi_{\text{ProT}}^{\text{GA}}$). The results are largely indistinguishable, with only small differences in the location of their stationary points (Figure 5.6). The shapes of the reactant and product basins are also quite similar with nearly identical eigenvalues (Table 5.2). This analysis can be extended further by quantifying the reaction coordinate motions in each basin. To this end, the extent of coupling between the normal mode motions was analyzed by comparing the normal mode basis to the phosphoryl/proton transfer basis. We consider the modes to be completely coupled if the phosphoryl/proton transfer component magnitudes are equal (*i.e.* the normal mode basis is rotated 45° with respect to the axes). This coupling can be expressed on a scale of 0 – 1 by calculating the absolute value of the cosine of the angle of rotation (Table 5.2). For both free energy surfaces the normal modes indicate significant decoupling of phosphoryl and proton transfer in the product states. In the reactant states, however, there is stronger coupling, with slow oscillation primarily along the proton transfer coordinate (see supporting information). The (reactive) mode orthogonal to this describes mostly high frequency phosphoryl transfer motion. The transition states are likewise indicative of strongly coupled motions, indicating that motion along the general acid coordinate increases as the reaction progresses.

Figure 5.6: Free energy surfaces of the HDVr catalytic reaction starting from an activated (post base) state with (bottom) and without (top) Mg^{2+} bound at the position hypothesized by Chen, *et al.*[231]. Axis and abscissa correspond to atom transfer coordinates for general acid proton transfer and phosphoryl transfer ($\xi_{\text{ProT}}^{\text{GA}}$ and ξ_{PhoT} , respectively). Free energies are in kcal/mol relative to the reactant minima with 5 kcal/mol separation of contour lines. Minima and saddle points (black diamonds) and a minimum free energy path (white dots) are also shown.

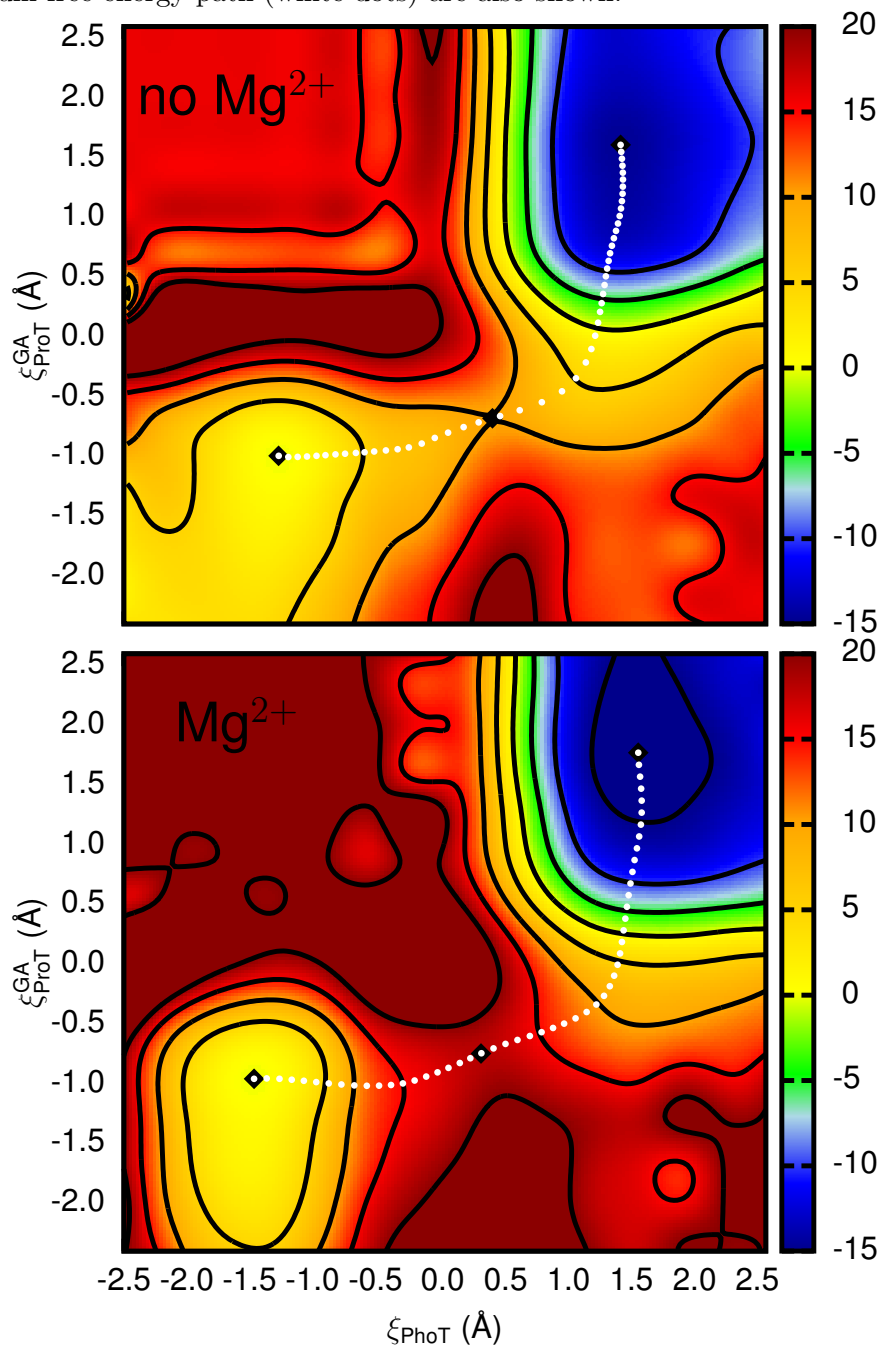
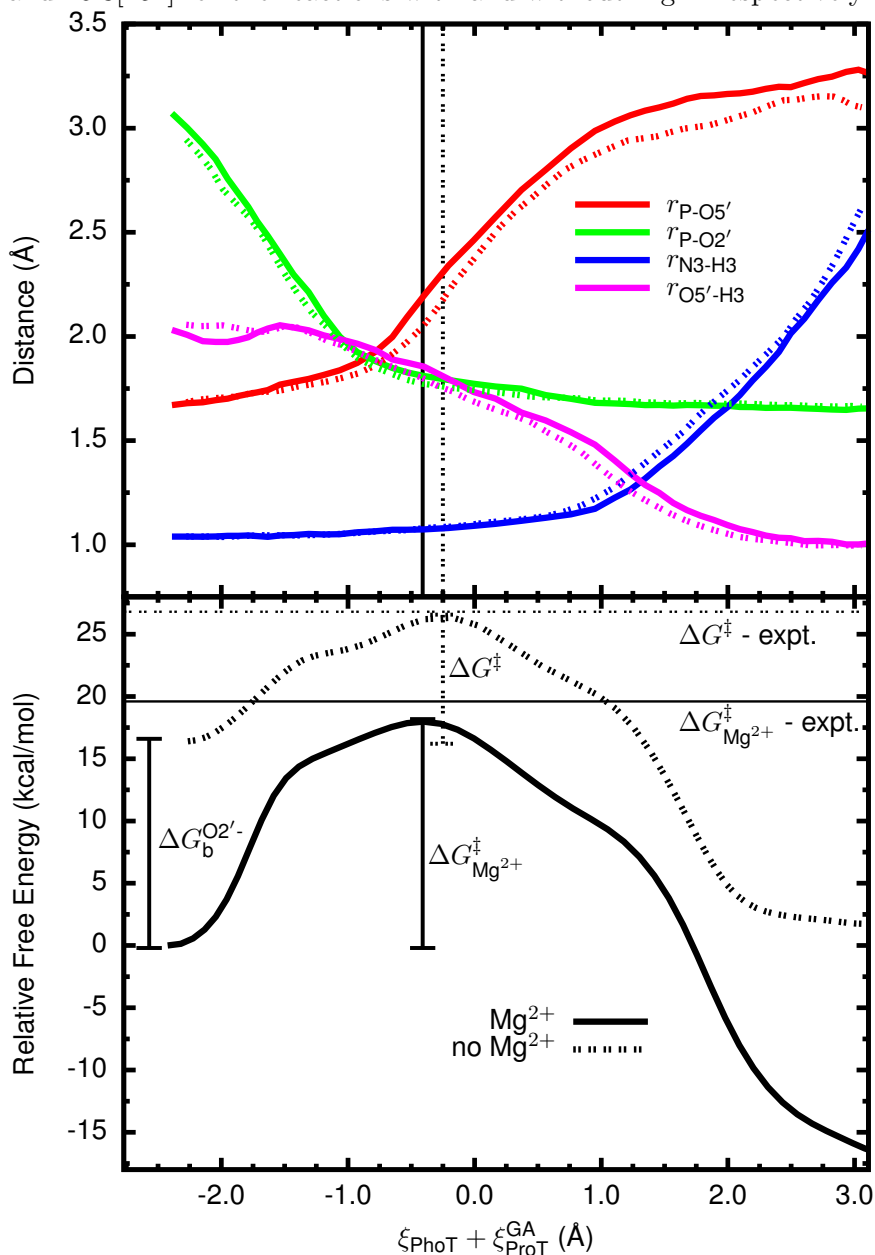


Table 5.2: Analysis of free energy profiles (FEPs) of the HDVr catalytic reaction (Figure 5.6) including normal mode analysis of stationary points and comparisons to previous studies. Relative free energies (ΔG) are in kcal/mol and frequencies (ν) are in ps^{-1} (see supporting information). Reaction coordinate (RC) coupling values represent a nominal measure of the amount of coupling between atom transfer modes (coupling = $|\cos \theta|$, where θ is the angle of rotation of the normal mode basis with respect to the phosphoryl transfer/proton transfer basis; 0 = completely decoupled; 1 = completely coupled). ^a Using k_{max} at 310 K, pH 7.0 from Ref. 237 ^b Using k_{obs} at 310 K, pH 6.0, 1 M NaCl from Ref. 252 ^b Using k_{obs} at 310 K, pH 6.1, 4 M NaCl and a genomic construct from Ref. 221 ^d Using string method with B3LYP/6-31G** from Ref. 255 ^e Not reported, estimated visually from graph.

	ΔG			RC		
	expt.	DFT ^d	FEP/corr.	coupling	ν_1	ν_2
reactant state						
Mg ²⁺	-	-	-	0.34	3.7	15.4
no Mg ²⁺	-	-	-	0.38	3.8	15.6
transition state						
Mg ²⁺	19.6 ^a	13	18.0	0.91	4.4 _i	17.9
no Mg ²⁺	26.8 ^b /24.2 ^c	2-4	10.0/26.4	0.96	3.8 _i	22.4
product state						
Mg ²⁺	-	\sim -6 ^e	-16.9	0.02	3.3	10.6
no Mg ²⁺	-	\sim -25 ^e	-14.6	0.11	3.3	11.0

Analysis of the relative free energies of the stationary points can be simplified by projecting the minimum free energy paths of each surface onto a sum of atom transfer

Figure 5.7: Reaction profiles (bottom) along the minimum free energy paths for the HDVr catalytic reaction starting from an activated (*i.e.* deprotonated) state with and without the presence of Mg^{2+} at the position hypothesized by Chen, *et al.*[231] (solid and dashed lines, respectively). The zeros of energy are set according to an estimated free energy of Mg^{2+} binding ($\Delta G_b^{\text{O}^{2-}} = 16.4$). Averages of selected bond lengths along the path are also shown (top). The calculated (including correction) and experimental barrier heights (from classical transition state theory) are 18.0 and 19.6-19.8[237, 252] and 26.4 and 26.8[252] for the reactions with and without Mg^{2+} respectively.



coordinates, $\xi_{\text{PhoT}} + \xi_{\text{ProT}}^{\text{GA}}$ (as described above in Section 5.2.2).¹ This leads to a simple, univariate function of the relative free energy linking the (activated) reactant and product states (Figure 5.7, bottom). Calculating these paths with either the MBAR or vFEP method gave indistinguishable results within statistical error. Although the energy differences between the reactant and transition states (neglecting Jacobian corrections) are quite different with and without Mg^{2+} (18.0 and 10.0 kcal/mol, respectively), several relevant average bond lengths are nearly identical along both paths (Figure 5.7, top). The only significant exception is the average length of the U-1:O5' to G1:P bond (Figure 5.7, top, red lines). Regardless, in both cases the start of the reaction is dominated by formation of the P-O2' bond (Figure 5.7, top, green lines) with an increasingly strong hydrogen bond between G1:O5' and C75:N3 as the transition state is crossed. This hydrogen bond concertedly changes to proton transfer as the P-O5' bond breaks (Figure 5.7, top, red and magenta lines). However, the covalent bond between the proton and its donor, C75:N3, does not appear to break until the reaction has progressed considerably (Figure 5.7, top, blue lines). In fact, this bond does not appear to form at all significantly until P-O5' bond breakage has nearly completed. This is similar to the observation that the normal modes display an increasing amount of coupling between the phosphoryl and proton transfer coordinates as the reaction progresses towards the transition state.

5.2.4 Discussion

The main purpose of the present work is to evaluate mechanistic scenarios departing from the suggested metal ion binding mode previously proposed in the literature[231, 234] by establishing the experimentally testable consequences and identifying key aspects deserving further consideration. The computational results reported here accomplish this goal by integrating several approaches aimed at different time and length scales (*i.e.* fast chemical events via QM/MM MD and slower/longer length RNA motions via MM MD) as well as solvent considerations (*i.e.* explicit MD and 3D-RISM). The resulting new information can be used to evaluate the specific atomistic reaction

¹ It is worth noting that this linear combination itself is probably *not* a good reaction coordinate, as its displacement vector does not overlap significantly with the active mode in two dimensions. This is also borne out in the relevant calculations in which a much lower barrier was obtained, as would be expected (data not shown).

pathways and allows structural/energetic analysis of the related ground states within a simple equilibrium scheme that relates Mg^{2+} binding and deprotonation of the O2' nucleophile (Figure 5.4).

Reaction Pathways and Ground States

In estimating relative activation energies from the QM/MM free energy surfaces calculated here, two major factors must be properly considered. First, the ground states (R^- and $\text{R}^-:\text{Mg}^{2+}$) of the free energy surfaces calculated here are not identical, as they correspond to two very different ionic bound states. Second, the ground states in these simulations are assumed to be pre-activated by equilibrium deprotonation of the O2' nucleophile (see Figure 5.4). For the enzymatic reaction at neutral or similar pH, however, the O2' is protonated in the ground state; deprotonation of the nucleophile is assumed to be an equilibrium process (specific base catalysis) that governs the population of the activated reactant state. Decreasing the population of the deprotonated O2' state in turn decreases the concentration of active ribozyme and therefore attenuates the intrinsic rate constant down to the experimentally observed value (*e.g.* the log-linear behavior observed in some pH-rate profiles). The population of active HDVr will depend on the correct protonation of both the O2' and C75.

Raman crystallographic measurements have indicated that the $\text{p}K_{\text{a}}$ of C75 may be anti-correlated with a Mg^{2+} binding event, but it is possible that it is anti-correlated with a different binding-site than the one investigated here[230]. Additionally, at pH values below the $\text{p}K_{\text{a}}$ of C75 (~ 6) this residue will be predominantly in its active form. As such, we proceed under the assumption that C75 is protonated in all states considered here. In the present work the deprotonation step is assumed to have already occurred. Instead, the pre-equilibrium assumptions can be used to estimate the fraction of active enzyme from the relative free energies of the various unreacted states (*i.e.* combinations of Mg^{2+} bound/unbound and O2' protonated/deprotonated). The combination of QM/MM and MM/3D-RISM provides a means of decomposing and dissecting the energetic contributions of processes that are not easily separated (or impossible to separate) experimentally (*i.e.* ion binding, deprotonation, and chemical reaction).

The free energy estimates from 3D-RISM generally agree in trend and magnitude (Table 5.1). The NLPB/SA estimates, however, do not appear to be physically realistic,

nor do the Na^+ pair distribution functions match well with those obtained from MD (Figure 5.5). As such, they are not considered further. The best estimate available here is from MM/3D-RISM-PSE2, as the KH closure is known to have difficulty in predicting small molecule hydration free energies[261] as well as excess chemical potentials for simple electrolytes[262] and excess ion distributions around DNA[263] (although still with greater fidelity than NLPB). It is difficult to assess the exact systematic/numerical errors of these estimates, but grid spacing effects likely will not cause errors greater than 0.5-1 kcal/mol[47] and systematic tests (data not shown) on the buffer size indicate errors less than 0.5 kcal/mol. Neglecting structural variations by clustering is probably the most substantial source of error, likely on the order of 1-2 kcal/mol.

Since it is well established that the HDVr is more catalytically proficient in the presence of divalent ions than monovalent ions[221, 252], it is perhaps initially surprising that the predicted free energy barrier from the activated precursor in the presence of a bound Mg^{2+} ion (18 kcal/mol) is almost double that in the Na^+ environment alone (10 kcal/mol, see Table 5.2 and Figure 5.7). However, as described, above, the reactant state on these surfaces is a *deprotonated* oxyanion and should clearly have strong, favorable electrostatic interactions with nearby cations. The negatively charged $\text{O}2'$ nucleophile near Mg^{2+} is thus better stabilized by the higher concentration of proximal positive charge. This charge preferentially stabilizes the reactant state over the transition state and effectively raises the barrier. In a related observation, P- $\text{O}2'$ bond formation is quite advanced in both transition states (see Figure 5.7), indicating that the negative charge is likely re-distributed away from the nucleophile and towards the leaving group, giving rise to less favorable electrostatic interactions. The fact that the transition state structures with and without Mg^{2+} are so similar (Figure 5.7, top) is sensible within a Hammond-Leffler framework since this stabilization effect should hold equally well for the charged product (note that Table 5.2 also shows that the “ligation” barrier is much lower in the absence of Mg^{2+}).

In order to achieve a meaningful comparison of the reaction barriers in the presence and absence of the active site Mg^{2+} , the ground states for these two reactions (RH and $\text{RH}:\text{Mg}^{2+}$) must be thermodynamically connected. The QM/MM simulations, however, depart from activated precursor states (R^- and $\text{R}^-:\text{Mg}^{2+}$) that assume nucleophile activation has occurred in a pre-equilibrium processes. Examining Figure 5.4, the most

direct route between the activated states of the two surfaces is to bind Mg^{2+} to the deprotonated nucleophile, R^- , to form $\text{R}^-:\text{Mg}^{2+}$. The best estimate of the free energy shift between the free energy profiles is thus ~ 16.4 kcal/mol (Figure 5.7). Clearly, taking this shift into account, the free energy barriers become much more closely aligned with experiment (Table 5.2); the barrier in the absence of Mg^{2+} is now properly situated above that in the presence of Mg^{2+} .

With the present calculations it is possible to estimate the relative populations of all four states within modest assumptions. First, the $\text{p}K_{\text{a}}$ shift between the Mg^{2+} bound and unbound states can be estimated from MM/3D-RISM. Although many of the shifts calculated here are quite large (4-5.5 units compared to ~ 1.3 in the presence of Ca^{2+} determined via NMR[255]), the clear indication is that bound Mg^{2+} reduces the energetic penalty of activating the nucleophile. This is of course only relevant in conjunction with the Mg^{2+} bound state being significantly populated, and this population can be estimated by effectively closing the cycle by inserting a relative free energy for one of the deprotonation steps. A reasonable estimate for ΔG_{a} can be made by assuming a $\text{p}K_{\text{a}}$ of ~ 13.5 (as for a free nucleotide[264]) for the unbound state $(\text{RH})^2$, in which case $\Delta G_{\text{a}} \approx 8.9$ kcal/mol and $\Delta G_{\text{a}}^{\text{Mg}^{2+}} = 3.6$ kcal/mol. Combining this with an estimate of -11.1 kcal/mol for $\Delta G_{\text{b}}^{\text{O}2\text{H}}$ would indicate that the lowest energy state in Figure 5.4 is that in which Mg^{2+} is bound and $\text{O}2'$ is protonated ($\text{RH}:\text{Mg}^{2+}$). The only other state with any appreciable population ($> 0.1\%$) is then the activated nucleophile bound to Mg^{2+} ($\text{R}^-:\text{Mg}^{2+}$).

Model Exploration and Testable Details

The key structural aspect of the mechanistic model studied here is the high level of rigidity in all of the ion bound states, regardless of the identity of the ion (or ions). This rigidity is present not only in the ground states (as evidenced by the low structural variation between MM MD trajectories), but also in the transition states (as evidenced by the high similarity between QM/MM saddle point geometries as well as their normal modes, Table 5.2). The latter case is especially interesting, as the indication is that

² Due to the RNA fold and large number of phosphates near the active site, the actual $\text{p}K_{\text{a}}$ is likely greater than this, perhaps as high as 15. However, using larger values does not change the overall interpretation here given the other errors inherent in the calculation.

the motions most principally aligned with the chemical step are not visibly perturbed by different metal ions. Therefore, with respect to chemical mechanism and transition state charge distribution, the most important aspect of this mechanistic model is the location and coordination pattern of the ions rather than their identity. If correct, this observation should offer non-trivial consequences for other aspects of the HDVr mechanism.

First of all, either both Mg^{2+} and Na^+ must be capable of promoting and stabilizing the same (or a highly similar) RNA fold or else other aspects of the HDVr are the main forces in attaining this state. The main candidates for these intrinsic stabilizing interactions are the highly specific hydrogen bonding pattern between C75:N3/G1:O5' and C75:N4/G1:*pro*-R_P/C22:*pro*-R_P and the rigorously conserved G25-U20 wobble pair (Figure 5.3)[212, 231]. The importance of the former has been supported by the apparent sensitivity to thio substitution at C22:*pro*-R_P[249, 250, 225] as well as rate enhancement of G1:S5'/C75c³C (3-deazacytosine) over G1:S5'/C75U mutants[226]. A similar inference can be made from the results of inactivating C75Z (zebularine) mutations, although an unambiguous interpretation is difficult due to a significant difference in the solution $\text{p}K_{\text{a}}$ of the presumed general acid (2.5 versus 4.2) in addition to removal of the exocyclic amine at C4[225]. The latter motif has been deemed critical through the characterization of multiple single/double and non-standard G25 mutants[238] as well as analysis of the pH-rate profiles and Mg^{2+} titrations of a G25A-U20C double mutant[237]. It has also been suggested that the hydrogen bonding pattern of this motif changes after the cleavage reaction, perhaps giving rise to differences in the affinity or geometry of a bound Mg^{2+} ion[238]. Although such a conformational change would presumably be much slower than chemistry, all of the present simulations seem to indicate long term stability of the wobble pair once the RNA fold is formed and thus describe no specific effect on the catalytic step.

Another aspect of the current mechanistic model is the high similarity between the reaction pathways with and without Mg^{2+} present. If these pathways are in fact similar, then a reasonable, zeroth order approximation would suggest that their heavy atom (¹⁸O) kinetic isotope effect (KIE) signatures would be similar. The present calculations are not sufficiently detailed (nuclear motions are not quantized) to make quantitative

predictions as to what the magnitude of these effects would be. Nonetheless, the mechanism described here makes specific predictions regarding the details of transition state bonding and the qualitative impact of this on observed KIEs can be predicted. Likewise, alternative mechanisms recently described by Ganguly, *et al.* for the Mg^{2+} and Na^+ reactions also make implicit predictions regarding the likely magnitude of isotope substitution on the reaction rate constant[255]. A key feature of the mechanisms in both their work and ours is the asynchronicity (although to different extents) of P-O5' bond cleavage, which lags behind P-O2' bond formation and leads to an associative transition state. Such a transition state takes the form of a phosphorane with significant charge accumulation on the non-bridging oxygens. Based on observed values for solution[141] and ribonuclease[6] catalyzed RNA 2'-*O*-transphosphorylation the O2' KIE is likely to be inverse due to a large contribution from formation of the new P-O2' vibrational mode. Non-bridging KIE values of 0.98-0.99 are observed for the late transition states for displacement of poor leaving groups such as *m*-nitrobenzyl and ribose[197, 142]. Conversely, the O5' effect is likely to be small since P-O bond cleavage is not far advanced. For specific base catalysis, in which the O5' departs as an oxyanion, the leaving group effect is large at 1.034[141]. However, in the ribonuclease A active site general acid catalysis from His119 (analogous to the role of C75) reduces this value to 1.015[6]. Interestingly, protonation from C75:N3 appears to occur late along the reaction coordinate. This feature is readily apparent in Figure 5.7 by noting the position with respect to the transition state where the lengths of the forming and cleaving bonds cross. The scenario here predicts that there would be minimal contribution to the O5' KIE from the formation of the new O5'-H bond. In contrast, in Ref. 255 this point is nearly coincident with the transition state.

An additional key difference between the mechanism predicted by Ganguly, *et al.* for the reaction with no Mg^{2+} and the results presented here is the temporary transfer of a proton from C75:N4 (the exocyclic amine) to the *pro*-R_P oxygen stabilizes a monoanionic phosphorane resulting in a stepwise mechanism. Results from analysis of non-enzymatic RNA 2'-*O*-transphosphorylation would suggest large inverse secondary KIEs for the non-bridging oxygens due to protonation and changes in hybridization resulting from formation of the phosphorane. Values of ~ 0.990 are observed for the

stepwise mechanism of acid catalysis of RNA 2'-*O*-triphosphorylation and for the associative transition state for phosphotriester reactions[141]. In contrast, reaction via a dianionic phosphorane or a similar negatively charged transition state would be characterized by very small non-bridging oxygen KIEs as observed for specific base reactions of phosphodiester with both good and poor leaving groups[197, 142].

Finally, the active site here contains two specific phosphate mediated interactions with Mg^{2+} , whereas in the Na^+ mechanism there are multiple binding sites (see Figure 5.5) that do not identically overlap with these interactions (although the overall backbone structure is apparently maintained). Recent work probed thio effects at the G1 non-bridge oxygens and established a biphasic reaction profile upon substitution at G1:*pro*-R_P[252]. A rescue effect was observed in the presence of Cd^{2+} , but not Mn^{2+} , in agreement with earlier experiments[249]. However, this is not necessarily confirmation of a metal-oxygen contact as the observed effect could also be due to perturbation of the hydrogen bond network with C22:*pro*-R_P (a possibility that was also described in Ref. 252). This latter interpretation is convincing, as a significant thio “interference,” on par with that at G1:*pro*-R_P, was also previously observed at C22:*pro*-R_P³. The similarly small thio effects for a G1:*pro*-R_P modification in Na^+ alone (1 M NaCl and 100 mM EDTA) and on the fast reacting species in the Mg^{2+} reaction would seem to suggest that the mechanistic pathway in the absence of Mg^{2+} remains largely unperturbed and presumably achieves the same fold, at least on a time scale faster than chemistry[252]. If this is true, a similar observation should be able to be made upon thio substitution at U23:*pro*-S_P, which was not possible in previous studies which employed modification-interference methods that can only monitor *pro*-R_P substitutions[249].

Further Aspects

The present work investigated one distinct Mg^{2+} binding mode suggested based on crystallographic data and this mode appears to be consistent with a large amount of structural, biochemical and computational data. However, it cannot be excluded that other binding modes may also be consistent with these data and give similar results within the present thermodynamic framework. Nonetheless, a key characteristic of the

³ It is worth noting that in Ref. 249 the phosphates are numbered with respect to the 3' positions. As such, the G1, G2, and C22 phosphates are referred to as the -1, 1, and 21 positions, respectively.

mechanism described here is that it involves multiple compensating effects due to the presence of Mg^{2+} . The protonated ground state provides a favorable binding position which gives rise to $\text{p}K_{\text{a}}$ shift of the nucleophile towards neutrality. This promotion of the activated precursor state is necessarily offset by stabilization of the reactant over the transition state, as the ion provides a +2 charge while the general acid provides only +1. This effect becomes obvious when the Mg^{2+} ion is replaced by more diffusely held Na^{+} ions and the barrier is lowered, but at the expense of a far less populated precursor state. If this reasoning is correct, then other chemical modifications to the system which lower the O2' $\text{p}K_{\text{a}}$ should enhance the cleavage rate in low concentrations of Mg^{2+} or even rescue the reaction in the presence of monovalent ions alone. Such modifications might likely include (but of course not be limited to) the fluoromethyl substitutions made by Ye, *et al.* at the C2' position of a chimeric RNA oligomer[265].

5.2.5 Conclusion

Molecular simulations provide a convenient and robust framework for analyzing mechanistic pathways, especially those pertaining to complex biocatalytic systems. The key strength of such methods is that they provide unambiguous atomistic detail that can be mapped to experimental observables. In the present work we have performed simulations of two hypothetical reaction channels in the HDVr in order to critically assess a recently proposed mechanistic model involving a site-bound Mg^{2+} . Once these results are properly contextualized within the required thermodynamic assumptions of metal ion binding and the resulting $\text{p}K_{\text{a}}$ shift of the nucleophile, they agree with much of the experimental data as well as simple chemical intuition. The results are consistent with a mechanistic model whereby an active site metal ion facilitates nucleophile activation, and C75 acts as a general acid catalyst. This study represents a distinct advance in that the results are integrated in a well-defined thermodynamic framework. The specific structural and dynamical details of this mechanism suggest several areas where additional experiments could probe this model further, especially the measurement of KIEs. In addition, a hitherto neglected thio substitution at U23 and chemical modification of the U-1:O2' $\text{p}K_{\text{a}}$ could provide evidence that would require amendment of this mechanistic model in order to maintain consistency.

References

- [1] Brian K. Radak, Michael E. Harris, and Darrin M. York. Molecular simulations of RNA 2'-*O*-transesterification reaction models in solution. *J. Phys. Chem. B*, 117:94–103, 2013.
- [2] Brian K. Radak, Melissa Romanus, Emilio Gallicchio, Tai-Sung Lee, Ole Weidner, Nan-Jie Deng, Peng He, Wei Dai, Darrin M. York, Ronald M. Levy, and Shantenu Jha. *A Framework for Flexible and Scalable Replica-Exchange on Production Distributed CI*, pages 26:1–26:8. XSEDE '13. 2013.
- [3] Brian K. Radak, Melissa Romanus, Tai-Sung Lee, Haoyuan Chen, Ming Huang, Antons Treikalis, Vivekanandan Balasubramanian, Shantenu Jha, and Darrin M. York. Characterization of the three-dimensional free energy manifold for the uracil ribonucleoside from asynchronous replica exchange simulations. **submitted**.
- [4] Tai-Sung Lee, Brian K. Radak, Anna Pabis, and Darrin M. York. A new maximum likelihood approach for free energy profile construction from molecular simulations. *J. Chem. Theory Comput.*, 9:153–164, 2013.
- [5] Tai-Sung Lee, Brian K. Radak, Ming Huang, Kin-Yiu Wong, and Darrin M. York. Roadmaps through free energy landscapes calculated using the multidimensional vFEP approach. *J. Chem. Theory Comput.*, 10:24–34, 2014.
- [6] Hong Gu, Shuming Zhang, Kin-Yiu Wong, Brian K. Radak, Thakshila Disanayake, Daniel L. Kellerman, Qing Dai, Masaru Miyagi, Vernon E. Anderson, Darrin M. York, Joseph A. Piccirilli, and Michael E. Harris. Experimental and computational analysis of the transition state for ribonuclease A-catalyzed RNA 2'-*O*-transphosphorylation. *Proc. Natl. Acad. Sci. USA*, 110:13002–13007, 2013.

- [7] Brian K. Radak, Tai-Sung Lee, Michael E. Harris, and Darrin M. York. A framework for assessment of metal-assisted nucleophile activation in the hepatitis delta virus ribozyme. **submitted**.
- [8] Jamal M. Buzayan, Wayne L. Gerlach, and George Bruening. Satellite tobacco ringspot virus RNA: A subset of the RNA sequence is sufficient for autolytic processing. *Proc. Natl. Acad. Sci. USA*, 83:8859–8862, 1986.
- [9] Gerry A. Prody, John T. Bakos, Jamal M. Buzayan, Irving R. Schneider, and George Bruening. Autolytic processing of dimeric plant virus satellite RNA. *Science*, 231:1577–1580, 1986.
- [10] Biliang Zhang and Thomas R. Cech. Peptidyl-transferase ribozymes: Trans reactions, structural characterization and ribosomal RNA-like features. *Chem. Biol.*, 5:539–554, 1998.
- [11] Thomas A. Steitz and Peter B. Moore. RNA, the first macromolecular catalyst: The ribosome is a ribozyme. *Trends Biochem. Sci.*, 28:411–418, 2003.
- [12] Saba Valadkhan. snRNAs as the catalysts of pre-mRNA splicing. *Curr. Opin. Chem. Biol.*, 9:603–608, 2005.
- [13] Saba Valadkhan. The spliceosome: A ribozyme at heart? *Biol. Chem.*, 388:693–697, 2007.
- [14] Chiu-Ho T. Webb, Nathan J. Riccitelli, Dana J. Ruminski, and Andrej Lupták. Widespread occurrence of self-cleaving ribozymes. *Science*, 326:953–, 2009.
- [15] Marcos de la Pena and Inmaculada Garcia-Robles. Ubiquitous presence of the hammerhead ribozyme motif along the tree of life. *RNA*, 16:1943–1950, 2010.
- [16] Adam Roth, Zasha Weinberg, Andy G. Y. Chen, Peter B. Kim, Tyler D. Ames, and Ronald R. Breaker. A widespread self-cleaving ribozyme class is revealed by bioinformatics. *Nature Chem. Biol.*, 10:56–62, 2014.
- [17] Randi M Jimenez, Eric Delwart, and Andrej Lupták. Structure-based search reveals hammerhead ribozymes in the human microbiome. *J. Biol. Chem.*, 286:7737–7743, 2011.

- [18] Kouros Salehi-Ashtiani, Andrej Lupták, Alexander Litovchick, and Jack W Szostak. A genomewide search for ribozymes reveals an HDV-like sequence in the human CPEB3 gene. *Science*, 313:1788–1792, 2006.
- [19] Walter Gilbert. The RNA world. *Nature*, 319:618, 1986.
- [20] Jennifer A. Doudna and Thomas R. Cech. The chemical repertoire of natural ribozymes. *Nature*, 418:222–228, 2002.
- [21] Ronald R Breaker. Riboswitches and the RNA world. *Cold Spring Harb. Perspect. Biol.*, 4:1–15, 2012.
- [22] Andres Jäschke. Artificial ribozymes and deoxyribozymes. *Curr. Opin. Struct. Biol.*, 11:321–326, 2001.
- [23] D. L. Nelson and M. M. Cox. *Lehninger Principles of Biochemistry*. W. H. Freeman; 4th Edition, New York, NY, 2005.
- [24] Samuel E. Butcher. Structure and function of the small ribozymes. *Curr. Opin. Struct. Biol.*, 11:315–320, 2001.
- [25] Nils G. Walter and Shiamalee Perumal. *The small ribozymes: Common and diverse features observed through the FRET lens*, volume 13 of *Springer Series in Biophysics*, pages 103–127. Springer Berlin Heidelberg, 2009.
- [26] R. A. More O’Ferrall. Relationships between *E2* and *E1cB* mechanisms of β -elimination. *J. Chem. Soc. B*, pages 274–277, 1970.
- [27] William P. Jencks. General acid-base catalysis of complex reactions in water. *Chem. Rev.*, 72:705–718, 1972.
- [28] Eric V. Anslyn and Dennis A. Dougherty. *Modern Physical Organic Chemistry*. University Science Books, Sausalito, CA, 2006.
- [29] Michael E. Harris and Adam G. Cassano. Experimental analyses of the chemical dynamics of ribozyme catalysis. *Curr. Opin. Chem. Biol.*, 12:626–639, 2008.
- [30] N. B. Leontis and E. Westhof. Geometric nomenclature and classification of RNA base pairs. *RNA*, 7:499–512, 2001.

- [31] Neocles B. Leontis, Jesse Stombaugh, and Eric Westhof. The non-Watson-Crick base pairs and their associated isostericity matrices. *Nucleic Acids Res.*, 30:3497–3531, 2002.
- [32] J. D. Watson and F. H. C. Crick. Molecular structure of nucleic acids. *Nature*, 171:737–738, 1953.
- [33] Donald A. McQuarrie. *Statistical Mechanics*. University Science Books, Mill Valley, CA, 1973.
- [34] Terrell L. Hill. *An Introduction to Statistical Thermodynamics*. Dover Publications, Mineola, NY, 1986.
- [35] Peter Atkins and Julio de Paula. *Physical Chemistry*. W. H. Freeman and Company, New York, 7th edition, 2002.
- [36] John G. Kirkwood. Statistical mechanics of fluid mixtures. *J. Chem. Phys.*, 3:300–313, 1935.
- [37] Ivar Stakgold. *Boundary Value Problems of Mathematical Physics*. Society for Industrial and Applied Mathematics, Philadelphia, PA, 2000.
- [38] J.-P. Hansen and I. R. McDonald. *Theory of Simple Liquids*. Academic Press, Amsterdam, The Netherlands, 2006.
- [39] David Chandler and Hans C. Andersen. Optimized cluster expansions for classical fluids. II. Theory of molecular liquids. *J. Chem. Phys.*, 57:1930–1937, 1972.
- [40] Dmitrii Beglov and Benoît Roux. An integral equation to describe the solvation of polar molecules in liquid water. *J. Phys. Chem. B*, 101:7821–7826, 1997.
- [41] Andriy Kovalenko and Fumio Hirata. Three-dimensional density profiles of water in contact with a solute of arbitrary shape: A RISM approach. *Chem. Phys. Lett.*, 290:237–244, 1998.
- [42] Andriy Kovalenko and Fumio Hirata. Self-consistent description of a metal-water interface by the Kohn-Sham density functional theory and the three-dimensional reference interaction site model. *J. Chem. Phys.*, 110:10095–10112, 1999.

- [43] Andriy Kovalenko and Fumio Hirata. Potentials of mean force of simple ions in ambient aqueous solution. I. Three-dimensional reference interaction site model approach. *J. Chem. Phys.*, 112:10391–10417, 2000.
- [44] Jerome K. Percus and George J. Yevick. Analysis of classical statistical mechanics by means of collective coordinates. *Phys. Rev.*, 110:1–13, 1958.
- [45] Stefan M. Kast and Thomas Kloss. Closed-form expressions of the chemical potential for integral equation closures with certain bridge functions. *J. Chem. Phys.*, 129:236101, 2008.
- [46] J. S. Perkyns and B. Montgomery Pettitt. A dielectrically consistent interaction site theory for solvent-electrolyte mixtures. *Chem. Phys. Lett.*, 190:626–630, 1992.
- [47] Tyler Luchko, Sergey Gusarov, Daniel R. Roe, Carlos Simmerling, David A. Case, Jack Tuszynski, and Andriy Kovalenko. Three-dimensional molecular theory of solvation coupled with molecular dynamics in AMBER. *J. Chem. Theory Comput.*, 6:607–624, 2010.
- [48] Nicholas Metropolis, Arianna W. Rosenbluth, Marshall N. Rosenbluth, and Augusta H. Teller. Equation of state calculations by fast computing machines. *J. Chem. Phys.*, 21:1087–1092, 1953.
- [49] H. C. Andersen. Molecular dynamics simulations at constant pressure and/or temperature. *J. Chem. Phys.*, 72:2384–2393, 1980.
- [50] Shuichi Nosé and M. L. Klein. Constant pressure molecular dynamics for molecular systems. *Mol. Phys.*, 50:1055–1076, 1983.
- [51] William G. Hoover. Canonical dynamics: Equilibrium phase-space distributions. *Phys. Rev. A*, 31:1695–1697, 1985.
- [52] Glenn J. Martyna, Michael L. Klein, and Mark Tuckerman. Nosé-hoover chains: The canonical ensemble via continuous dynamics. *J. Chem. Phys.*, 97:2635–2643, 1992.
- [53] Robert Zwanzig. *Nonequilibrium Statistical Mechanics*. Oxford University Press, New York, New York, 2001.

- [54] Axel Brünger, Charles L. Brooks III, and Martin Karplus. Stochastic boundary conditions in molecular dynamics simulations of ST2 water. *Chem. Phys. Lett.*, 105:495–500, 1984.
- [55] John E. Mertz and B. Montgomery Pettitt. Molecular dynamics at a constant pH. *Int. J. Supercomput. Appl. High Perform. Comput.*, 8:47–53, 1994.
- [56] António M. Baptista, Paulo J. Martel, and Steffen B. Petersen. Simulation of protein conformational freedom as a function of pH: Constant-pH molecular dynamics using implicit titration. *Proteins*, 27:523–544, 1997.
- [57] Roland Bürigi, Peter A. Kollman, and Wilfred F. van Gunsteren. Simulating proteins at constant pH: An approach combining molecular dynamics and Monte Carlo simulation. *Proteins*, 47:469–480, 2002.
- [58] Maciej Dlugosz and Jan M. Antosiewicz. Constant-pH molecular dynamics simulations: A test case of succinic acid. *Chem. Phys.*, 302:161–170, 2004.
- [59] John Mongan, David A. Case, and J. Andrew McCammon. Constant pH molecular dynamics in generalized Born implicit solvent. *J. Comput. Chem.*, 25:2038–2048, 2004.
- [60] Michael S. Lee, Freddie R. Salsbury, Jr., and Charles L. Brooks III. Constant-pH molecular dynamics using continuous titration coordinates. *Proteins*, 56:738–752, 2004.
- [61] Harry A. Stern. Molecular simulation with variable protonation states at constant pH. *J. Chem. Phys.*, 126:164112, 2007.
- [62] Robert W. Zwanzig. High-temperature equation of state by a perturbation method. I. nonpolar gases. *J. Chem. Phys.*, 22:1420–1426, 1954.
- [63] M. P. Allen and D. J. Tildesley. *Computer Simulation of Liquids*. Oxford University Press, Oxford, 1987.
- [64] Daan Frenkel and Berend Smit. *Understanding Molecular Simulation*. Academic Press, San Diego, CA, 2002.

- [65] G. M. Torrie and J. P. Valleau. Nonphysical sampling distributions in Monte Carlo free-energy estimation: Umbrella sampling. *J. Comput. Phys.*, 23:187–199, 1977.
- [66] B. Widom. Some topics in the theory of fluids. *J. Chem. Phys.*, 39:2808–2812, 1963.
- [67] Charles H. Bennett. Efficient estimation of free energy differences from Monte Carlo data. *J. Comput. Phys.*, 22:245–268, 1976.
- [68] Shankar Kumar, Djamal Bouzida, Robert H. Swendsen, Peter A. Kollman, and John M. Rosenberg. The weighted histogram analysis method for free-energy calculations on biomolecules. I. The method. *J. Comput. Chem.*, 13:1011–1021, 1992.
- [69] Christian Bartels. Analyzing biased Monte Carlo and molecular dynamics simulations. *Chem. Phys. Lett.*, 331:446–454, 2000.
- [70] Marc Souaille and Benoît Roux. Extension to the weighted histogram analysis method: Combining umbrella sampling with free energy calculations. *Comput. Phys. Commun.*, 135:40–57, 2001.
- [71] Michael R. Shirts and John D. Chodera. Statistically optimal analysis of samples from multiple equilibrium states. *J. Chem. Phys.*, 129:124105, 2008.
- [72] Zhiqiang Tan, Emilio Gallicchio, Mauro Lapelosa, and Ronald M. Levy. Theory of binless multi-state free energy estimation with applications to protein-ligand binding. *J. Chem. Phys.*, 136:144102, 2012.
- [73] Emilio Gallicchio, Michael Andrec, Anthony K. Felts, and Ronald M. Levy. Temperature weighted histogram analysis method, replica exchange, and transition paths. *J. Phys. Chem. B*, 109:6722–6731, 2005.
- [74] John D. Chodera, William C. Swope, Jed W. Pitera, Chaok Seok, and Ken A. Dill. Use of the weighted histogram analysis method for the analysis of simulated and parallel tempering simulations. *J. Chem. Theory Comput.*, 3:26–41, 2007.

- [75] Michael R. Shirts, Eric Bair, Giles Hooker, and Vijay S. Pande. Equilibrium free energies from nonequilibrium measurements using maximum-likelihood methods. *Phys. Rev. Lett.*, 91:140601, 2003.
- [76] Michael R. Shirts and Vijay S. Pande. Comparison of efficiency and bias of free energies computed by exponential averaging, the Bennett acceptance ratio, and thermodynamic integration. *J. Chem. Phys.*, 122:144107, 2005.
- [77] Christophe Chipot and Andrew Pohorille, editors. *Free Energy Calculations: Theory and Applications in Chemistry and Biology*, volume 86 of *Springer Series in Chemical Physics*. Springer, New York, 2007.
- [78] Johannes Kästner. Umbrella sampling. *WIREs Comput. Mol. Sci.*, 1:932–942, 2011.
- [79] Donald Hamelberg, John Mongan, and J. Andrew McCammon. Accelerated molecular dynamics: A promising and efficient simulation method for biomolecules. *J. Chem. Phys.*, 120:11919–11929, 2004.
- [80] Christophe Chipot, Peter A. Kollman, and David A. Pearlman. Alternative approaches to potential of mean force calculations: Free energy perturbation versus thermodynamic integration. case study of some representative nonpolar interactions. *J. Comput. Chem.*, 17:1112–1131, 1996.
- [81] Johannes Kästner and Walter Thiel. Bridging the gap between thermodynamic integration and umbrella sampling provides a novel analysis method: “Umbrella integration”. *J. Chem. Phys.*, 123:144104, 2005.
- [82] Johannes Kästner and Walter Thiel. Analysis of the statistical error in umbrella sampling simulations by umbrella integration. *J. Chem. Phys.*, 124:234106, 2006.
- [83] Johannes Kästner. Umbrella integration in two or more reaction coordinates. *J. Chem. Phys.*, 131:034109, 2009.
- [84] Johannes Kästner. Umbrella integration with higher-order correction terms. *J. Chem. Phys.*, 136:234102, 2012.

- [85] Jodi E. Basner and Christopher Jarzynski. Binless estimation of the potential of mean force. *J. Phys. Chem. B*, 112:12722–12729, 2008.
- [86] Eric Darve and Andrew Pohorille. Calculating free energies using average force. *J. Chem. Phys.*, 115:9169–9183, 2001.
- [87] J. P. Valleau and D. N. Card. Monte Carlo estimation of the free energy by multistage sampling. *J. Chem. Phys.*, 57:5457–5462, 1972.
- [88] Hongxing Lei and Yong Duan. Improved sampling methods for molecular simulation. *Curr. Opin. Struct. Biol.*, 17:187–191, 2007.
- [89] Ayori Mitsutake, Yoshiharu Mori, and Yuko Okamoto. Enhanced sampling algorithms. *Methods Mol. Biol.*, 924:153–195, 2013.
- [90] Emilio Gallicchio and Ronald M. Levy. Advances in all atom sampling methods for modeling protein-ligand binding affinities. *Curr. Opin. Struct. Biol.*, 21:161–166, 2011.
- [91] Xiongwu Wu, Milan Hodoscek, and Bernard R. Brooks. Replica exchanging self-guided Langevin dynamics for efficient and accurate conformational sampling. *J. Chem. Phys.*, 137:044106, 2012.
- [92] Mikolai Fajer, Robert V. Swift, and J. Andrew McCammon. Using multistate free energy techniques to improve the efficiency of replica exchange accelerated molecular dynamics. *J. Comput. Chem.*, 30:1719–1725, 2009.
- [93] Wei Jiang, Milan Hodoscek, and Benoît Roux. Computation of absolute hydration and binding free energy with free energy perturbation distributed replica-exchange molecular dynamics. *J. Chem. Theory Comput.*, 5:2583–2588, 2009.
- [94] Yilin Meng, Danial Sabri Dashti, and Adrian E. Roitberg. Computing alchemical free energy differences with Hamiltonian replica exchange molecular dynamics (H-REMD) simulations. *J. Chem. Theory Comput.*, 7:2721–2727, 2011.
- [95] Mehrnoosh Arrar, Cesar Augusto F. de Oliveira, Mikolai Fajer, William Sinko,

- and J. Andrew McCammon. w-REXAMD: A Hamiltonian replica exchange approach to improve free energy calculations for systems with kinetically trapped conformations. *J. Chem. Theory Comput.*, 9:18–23, 2013.
- [96] Jason A. Wallace and Jana K. Shen. Continuous constant pH molecular dynamics in explicit solvent with pH-based replica exchange. *J. Chem. Theory Comput.*, 7:2617–2629, 2011.
- [97] Satoru G. Itoh, Ana Damjanović, and Bernard R. Brooks. pH replica-exchange method based on discrete protonation states. *Proteins*, 79:3420–3436, 2011.
- [98] Danial Sabri Dashti and Adrian E. Roitberg. pH-replica exchange molecular dynamics in proteins using a discrete protonation method. *J. Phys. Chem. B*, 116:8805–8811, 2012.
- [99] Wei Jiang and Benoît Roux. Free energy perturbation Hamiltonian replica-exchange molecular dynamics (FEP/H-REMD) for absolute ligand binding free energy calculations. *J. Chem. Theory Comput.*, 6:2559–2565, 2010.
- [100] Wei Jiang, Yun Luo, Luca Maragliano, and Benoît Roux. Calculation of free energy landscape in multi-dimensions with Hamiltonian-exchange umbrella sampling on petascale supercomputer. *J. Chem. Theory Comput.*, 8:4672–4680, 2012.
- [101] Christina Bergonzo, Niel M. Henriksen, Daniel R. Roe, Jason M. Swails, Adrian E. Roitberg, and Thomas E. Cheatham III. Multidimensional replica exchange molecular dynamics yields a converged ensemble of an RNA tetranucleotide. *J. Chem. Theory Comput.*, 10:492–499, 2014.
- [102] Juyong Lee, Benjamin T. Miller, Ana Damjanović, and Bernard R. Brooks. Constant pH molecular dynamics in explicit solvent with enveloping distribution sampling and Hamiltonian exchange. *J. Chem. Theory Comput.*, 10:2738–2750, 2014.
- [103] Yuji Sugita and Yuko Okamoto. Replica-exchange molecular dynamics method for protein folding. *Chem. Phys. Lett.*, 314:141–151, 1999.
- [104] Yuji Sugita, Akio Kitao, and Yuko Okamoto. Multidimensional replica-exchange method for free-energy calculations. *J. Chem. Phys.*, 113:6042–6051, 2000.

- [105] A. P. Lyubartsev, A. A. Martsinovski, S. V. Shevkunov, and P. N. Vorontsov-Velyaminov. New approach to Monte Carlo calculation of the free energy: Method of expanded ensembles. *J. Chem. Phys.*, 96:1776, 1992.
- [106] John D. Chodera and Michael R. Shirts. Replica exchange and expanded ensemble simulations as Gibbs sampling: Simple improvements for enhanced mixing. *J. Chem. Phys.*, 135:194110, 2011.
- [107] Satoru G. Itoh and Hisashi Okumura. Replica-permutation method with the Suwa-Todo algorithm beyond the replica-exchange method. *J. Chem. Theory Comput.*, 9:570–581, 2013.
- [108] John D. Chodera, William C. Swope, Frank Noe, Jan-Hendrik Prinz, Michael R. Shirts, and Vijay S. Pande. Dynamical reweighting: Improved estimates of dynamical properties from simulations at multiple temperatures. *J. Chem. Phys.*, 134:244107, 2011.
- [109] W. D. Cornell, P. Cieplak, C. I. Bayly, I. R. Gould, Kenneth M. Merz, Jr., D. M. Ferguson, D. C. Spellmeyer, T. Fox, J. W. Caldwell, and P. A. Kollman. A second generation force field for the simulation of proteins, nucleic acids and organic molecules. *J. Am. Chem. Soc.*, 117:5179–5197, 1995.
- [110] Junmei Wang, Piotr Cieplak, and Peter A. Kollman. How well does a restrained electrostatic potential (RESP) model perform in calculating conformational energies of organic biological molecules? *J. Comput. Chem.*, 21:1049–1074, 2000.
- [111] Alberto Pérez, Iván Marchán, David Svozil, Jiri Sponer, Thomas E. Cheatham III, Charles A. Loughton, and Modesto Orozco. Refinement of the AMBER force field for nucleic acids: Improving the description of α/γ conformers. *Biophys. J.*, 92:3817–3829, 2007.
- [112] Marie Zgarbová, Michal Otyepka, Jiří Šponer, Arnošt Mládek, Pavel Banáš, Thomas E. Cheatham III, and Petr Jurečka. Refinement of the Cornell et al. nucleic acids force field based on reference quantum chemical calculations of glycosidic torsion profiles. *J. Chem. Theory Comput.*, 7:2886–2902, 2011.

- [113] Nicolas Foloppe and Alexander D. MacKerell, Jr. All-atom empirical force field for nucleic acids: I. Parameter optimization based on small molecule and condensed phase macromolecular target data. *J. Comput. Chem.*, 21:86–104, 2000.
- [114] Alexander D. MacKerell, Jr. and Nilesh K. Banavali. All-atom empirical force field for nucleic acids: II. Application to molecular dynamics simulations of DNA and RNA in solution.
- [115] Evelyn Mayaan, Adam Moser, Alexander D. MacKerell, Jr., and Darrin M. York. CHARMM force field parameters for simulation of reactive intermediates in native and thio-substituted ribozymes. *J. Comput. Chem.*, 28:495–507, 2007.
- [116] Kwangho Nam, Qiang Cui, Jiali Gao, and Darrin M. York. Specific reaction parametrization of the AM1/d Hamiltonian for phosphoryl transfer reactions: H, O, and P atoms. *J. Chem. Theory Comput.*, 3:486–504, 2007.
- [117] Michael Gaus, Xiya Lu, Marcus Elstner, and Qiang Cui. Parameterization of DFTB3/3OB for sulfur and phosphorus for chemical and biological applications. *J. Chem. Theory Comput.*, 10:1518–1537, 2014.
- [118] Tom Darden, Darrin York, and Lee Pedersen. Particle mesh Ewald: An $N \log(N)$ method for Ewald sums in large systems. *J. Chem. Phys.*, 98:10089–10092, 1993.
- [119] Ulrich Essmann, Lalith Perera, Max L. Berkowitz, Tom Darden, Lee Hsing, and Lee G. Pedersen. A smooth particle mesh Ewald method. *J. Chem. Phys.*, 103:8577–8593, 1995.
- [120] Kwangho Nam, Jiali Gao, and Darrin M. York. An efficient linear-scaling Ewald method for long-range electrostatic interactions in combined QM/MM calculations. *J. Chem. Theory Comput.*, 1:2–13, 2005.
- [121] Ross C. Walker, Michael F. Crowley, and David A. Case. The implementation of a fast and accurate QM/MM potential method in Amber. *J. Comput. Chem.*, 29:1019–1031, 2008.
- [122] Timothy J. Giese, Haoyuan Chen, Thakshila Dissanayake, George M. Giambasu,

- Hugh Heldenbrand, Ming Huang, Erich R. Kuechler, Tai-Sung Lee, Maria T. Panteva, Brian K. Radak, and Darrin M. York. A variational linear-scaling framework to build practical, efficient next-generation orbital-based quantum force fields. *J. Chem. Theory Comput.*, 9:1417–1427, 2013.
- [123] Daniel Trzesniak, Anna-Pitschna E. Kunz, and Wilfred F. van Gunsteren. A comparison of methods to compute the potential of mean force. *Chem. Phys. Chem.*, 8:162–169, 2007.
- [124] Michel A. Cuendet and Mark E. Tuckerman. Free energy reconstruction from metadynamics or adiabatic free energy dynamics simulations. *J. Chem. Theory Comput.*, 10:2975–2986, 2014.
- [125] B.W. Silverman. *Density Estimation for Statistics and Data Analysis*. Chapman and Hall, New York, 1986.
- [126] M. P. Wand and M. C. Jones. *Kernel Smoothing*. Chapman and Hall, Boca Raton, FL, 1995.
- [127] Simon J. Sheather. Density Estimation. *Statist. Sci.*, 19:588–597, 2004.
- [128] Murray Rosenblatt. Remarks on some nonparametric estimates of a density function. *Annal. Math. Statist.*, 27:832–837, 1956.
- [129] Emanuel Parzen. On estimation of a probability density function and mode. *Annal. Math. Statist.*, 33:1065–1076, 1962.
- [130] R. A. Fisher. On the mathematical foundations of theoretical statistics. *Phil. Trans. R. Soc. Lond. A*, 222:309–368, 1922.
- [131] A.W.F. Edwards. *Likelihood*. Cambridge University Press, Cambridge, 1972.
- [132] Alan Grossfield. *WHAM: the weighted histogram analysis method* [Online], version 2.0.4; <http://membrane.urmc.rochester.edu/content/wham> (accessed April 2012).
- [133] John D. Chodera, Michael R. Shirts, and Kyle A. Beauchamp. *pymbar* [Online], version 2.1.0; <https://github.com/choderalab/pymbar> (accessed July 2014).

- [134] Christian Bartels and Martin Karplus. Multidimensional adaptive umbrella sampling: Applications to main chain and side chain peptide conformations. *J. Comput. Chem.*, 18:1450–1462, 1997.
- [135] John D. Chodera, Nina Singhal, Vijay S. Pande, Ken A. Dill, and William C. Swope. Automatic discovery of metastable states for the construction of Markov models of macromolecular conformational dynamics. *J. Chem. Phys.*, 126:155101, 2007.
- [136] D. A. Case, V. Babin, J. T. Berryman, R. M. Betz, Q. Cai, D.S Cerutti, T. E. Cheatham III, T. A. Darden, R. E. Duke, H. Gohlke, A. W. Götz, S. Gusarov, N. Homeyer, P. Janowski, J. Kaus, I. Kolossváry, A. Kovalenko, Tai-Sung Lee, S. Le Grand, T. Luchko, R. Luo, B. Madej, K. M. Merz, F. Paesani, D. R. Roe, A. Roitberg, C. Sagui, R. Salomon-Ferrer, G. Seabra, C. L. Simmerling, W. Smith, J. Swails, R. C. Walker, J. Wang, R. M. Wolf, X. Wu, and P. A. Kollman. *AMBER 14*. University of California, San Francisco, San Francisco, CA, 2014.
- [137] Ming Huang, Timothy J. Giese, Tai-Sung Lee, and Darrin M. York. Improvement of DNA and RNA sugar pucker profiles from semiempirical quantum methods. *J. Chem. Theory Comput.*, 10:1538–1545, 2014.
- [138] Hans W. Horn, William C. Swope, Jed W. Pitera, Jeffrey D. Madura, Thomas J. Dick, Greg L. Hura, and Teresa Head-Gordon. Development of an improved four-site water model for biomolecular simulations: TIP4P-Ew. *J. Chem. Phys.*, 120:9665–9678, 2004.
- [139] Abhinav Thota, Andre Luckow, and Shantenu Jha. Efficient large-scale replica-exchange simulations on production infrastructure. *Phil. Trans. R. Soc. A*, 369:3318–3335, 2011.
- [140] *RepEx* [Online]; <http://radical-cybertools.github.io/RepEx> (accessed August 2014).

- [141] Michael .E. Harris, Qing Dai, Hong Gu, Daniel L. Kellerman, Joseph A. Piccirilli, and Vernon E. Anderson. Kinetic isotope effects for RNA cleavage by 2'-*O*-transphosphorylation: Nucleophilic activation by specific base. *J. Am. Chem. Soc.*, 132:11613–11621, 2010.
- [142] Kin-Yiu Wong, Hong Gu, Shuming Zhang, Joseph A. Piccirilli, Michael E. Harris, and Darrin M. York. Characterization of the reaction path and transition states for RNA transphosphorylation models from theory and experiment. *Angew. Chem. Int. Ed.*, 51:647–651, 2012.
- [143] Qing Dai, John K. Frederiksen, Vernon E. Anderson, Michael E. Harris, and Joseph A. Piccirilli. Efficient synthesis of [2'¹⁸O]uridine and its incorporation into oligonucleotides: A new tool for mechanistic study of nucleotidyl transfer reactions by isotope effect analysis. *J. Org. Chem.*, 73:309–311, 2008.
- [144] L. Sharmeen, M. Y. Kuo, G. Dinter-Gottlieb, and J. Taylor. Antigenomic RNA of human hepatitis delta virus can undergo self-cleavage. *J. Virol.*, 62:2674–2679, 1988.
- [145] Barry J. Saville and Richard A. Collins. A site-specific self-cleavage reaction performed by a novel RNA in neurospora mitochondria. *Cell*, 61:685–696, 1990.
- [146] Wade C. Winkler, Ali Nahvi, Adam Roth, Jennifer A. Collins, and Ronald R. Breaker. Control of gene expression by a natural metabolite-responsive ribozyme. *Nature*, 428:281–286, 2004.
- [147] William G. Scott. Ribozymes. *Curr. Opin. Struct. Biol.*, 17:280–286, 2007.
- [148] Adrian R. Ferré-D'Amaré and William G. Scott. Small self-cleaving ribozymes. *Cold Spring Harb. Perspect Biol.*, 2:a003574, 2010.
- [149] Tatsuro Koike and Yasuo Inoue. Structure and reactivity of oligonucleotides. Part I. kinetics of the non-enzymatic transphorylation of adenylyl-(3'-5')-adenosine 3'-phosphate and other dinucleotides. *Chem. Lett.*, 1:569–572, 1972.

- [150] Eric Anslyn and Ronald Breslow. On the mechanism of catalysis by ribonuclease: Cleavage and isomerization of the dinucleotide UpU catalyzed by imidazole buffers. *J. Am. Chem. Soc.*, 111:4473–4482, 1989.
- [151] Pia Järvinen, Mikko Oivanen, and Harri Lönnberg. Interconversion and phosphoester hydrolysis of 2',5'- and 3',5'-dinucleoside monophosphates: Kinetics and mechanisms. *J. Org. Chem.*, 56:5396–5401, 1991.
- [152] Mikko Oivanen, Satu Kuusela, and Harri Lönnberg. Kinetics and mechanisms for the cleavage and isomerization of the phosphodiester bonds of RNA by Brønsted acids and bases. *Chem. Rev.*, 98:961–990, 1998.
- [153] Yingfu Li and Ronald R. Breaker. Kinetics of RNA degradation by specific base catalysis of transesterification involving the 2'-hydroxyl group.
- [154] D. M. Brown and D. A. Usher. Hydrolysis of hydroxyalkyl phosphate esters: Effect of changing ester group. *J. Chem. Soc.*, 87:6558–6564, 1965.
- [155] D. A. Usher, D. I. Richardson, Jr., and D. G. Oakenfull. Models of ribonuclease action. II. Specific acid, specific base, and neutral pathways for hydrolysis of a nucleotide diester analog. *J. Am. Chem. Soc.*, 92:4699–4712, 1970.
- [156] Andrew M. Davis, Adrian D. Hall, and Andrew Williams. Charge description of base-catalyzed alcoholysis of aryl phosphodiester: A ribonuclease model. *J. Am. Chem. Soc.*, 110:5105–5108, 1988.
- [157] Alvan C. Hengge, Karol S. Bruzik, Aleksandra E. Tobin, W. W. Cleland, and Ming-Daw Tsai. Kinetic isotope effects and stereochemical studies on a ribonuclease model: Hydrolysis reactions of uridine 3'-nitrophenyl phosphate. *Bioorg. Chem.*, 28:119–133, 2000.
- [158] Mikko Oivanen, Sergey N. Mikhailov, Nelly Sh. Padyukova, and Harri Lönnberg. Kinetics of mutual isomerization of the phosphonate analogs of dinucleoside 2',5'- and 3',5'-monophosphates in aqueous solution. *J. Org. Chem.*, 58:1617–1619, 1993.
- [159] Kwangho Nam, Jiali Gao, and Darrin M. York. Quantum mechanical/molecular

- mechanical simulation study of the mechanism of hairpin ribozyme catalysis. *J. Am. Chem. Soc.*, 130:4680–4691, 2008.
- [160] Kin-Yiu Wong, Tai-Sung Lee, and Darrin M. York. Active participation of the Mg^{2+} ion in the reaction coordinate of RNA self-cleavage catalyzed by the hammerhead ribozyme. *J. Chem. Theory Comput.*, 7:1–3, 2011.
- [161] Tai-Sung Lee and Darrin M York. Computational mutagenesis studies of hammerhead ribozyme catalysis. *J. Am. Chem. Soc.*, 132:13505–13518, 2010.
- [162] N. Kyle Tanner, Sophie Schaff, Gilbert Thill, Elisabeth Petit-Koskas, Anne-Marie Crain-Denoyelle, and Eric Westhof. A three-dimensional model of hepatitis delta virus ribozyme based on biochemical and mutational analyses. *Curr. Biol.*, 4:488–498, 1994.
- [163] L. Vitagliano, A. Merlino, A. Zagari, and L. Mazzarella. Productive and nonproductive binding to ribonuclease A: X-ray structure of two complexes with uridylyl (2',5') guanosine. *Protein Sci.*, 9:1217–1225, 2000.
- [164] D. A. Case, T. A. Darden, T. E. Cheatham III, C. L. Simmerling, J. Wang, R. E. Duke, R. Luo, R. C. Walker, W. Zhang, K. M. Merz, B. Roberts, S. Hayik, A. Roitberg, G. Seabra, J. Swails, A. W. Götz, I. Kolossváry, K. F. Wong, F. Paesani, J. Vanicek, R. M. Wolf, J. Liu, X. Wu, S. R. Brozell, T. Steinbrecher, H. Gohlke, Q. Cai, X. Ye, J. Wang, M.-J. Hsieh, G. Cui, D. R. Roe, D. H. Mathews, M. G. Seetin, R. Salomon-Ferrer, C. Sagui, V. Babin, T. Luchko, S. Gusarov, A. Kovalenko, and P. A. Kollman. *AMBER 12*. University of California, San Francisco, San Francisco, CA, 2012.
- [165] J. P. Ryckaert, G. Ciccotti, and H. J. C. Berendsen. Numerical integration of the Cartesian equations of motion of a system with constraints: Molecular dynamics of n-alkanes. *J. Comput. Phys.*, 23:327–341, 1977.
- [166] Shuichi Miyamoto and Peter A. Kollman. SETTLE: An analytic version of the SHAKE and RATTLE algorithms for rigid water models. *J. Comput. Chem.*, 13:952–962, 1992.

- [167] H. J. C. Berendsen, J. P. M. Postma, W. F. van Gunsteren, A. Dinola, and J. R. Haak. Molecular dynamics with coupling to an external bath. *J. Chem. Phys.*, 81:3684–3690, 1984.
- [168] W. L. Jorgensen, J. Chandrasekhar, J. D. Madura, R. W. Impey, and M. L. Klein. Comparison of simple potential functions for simulating liquid water. *J. Chem. Phys.*, 79:926–935, 1983.
- [169] In Suk Joung and Thomas E. Cheatham III. Determination of alkali and halide monovalent ion parameters for use in explicitly solvated biomolecular simulations. *J. Phys. Chem. B*, 112:9020–9041, 2008.
- [170] Violeta López-Canut, Sergio Martí, Juan Bertrán, Vicente Moliner, and Iñaki Tuñón. Theoretical modeling of the reaction mechanism of phosphate monoester hydrolysis in alkaline phosphatase. *J. Phys. Chem. B*, 113:7816–7824, 2009.
- [171] Kwangho Nam, Jiali Gao, and Darrin M. York. *New QM/MM models for multi-scale simulation of phosphoryl transfer reactions in solution*, pages 201–218. John Wiley & Sons, Inc., 2008.
- [172] Violeta López-Canut, Javier Ruiz-Pernía, Iñaki Tuñón, Silvia Ferrer, and Vicent Moliner. Theoretical modeling on the reaction mechanism of *p*-nitrophenylmethylphosphate alkaline hydrolysis and its kinetic isotope effects. *J. Chem. Theory Comput.*, 5:439–442, 2009.
- [173] Timothy J. Giese, Brent A. Gregersen, Yun Liu, Kwangho Nam, Evelyn Mayaan, Adam Moser, Kevin Range, Olalla Nieto Faza, Carlos Silva Lopez, Angel Rodriguez de Lera, Gijs Schaftenaar, Xabier Lopez, Tai-Sung Lee, George Karypis, and Darrin M. York. QCRNA 1.0: A database of quantum calculations of rna catalysis. *J. Mol. Graph. Model.*, 25:423–433, 2006.
- [174] Enrique Marcos, Josep M. Anglada, and Ramon Crehuet. Description of pentacoordinated phosphorus under an external electric field: Which basis sets and semi-empirical methods are needed? *Phys. Chem. Chem. Phys.*, 10:2442–2450, 2008.

- [175] Kyoyeon Park, Andreas W. Götz, Ross C. Walker, and Francesco Paesani. Application of adaptive QM/MM methods to molecular dynamics simulations of aqueous systems. *J. Chem. Theory Comput.*, 8:2868–2877, 2012.
- [176] G. Monard, M. I. Bernal-Uruchurtu, A. van der Vaart, K. M. Merz Jr., and M. F. Ruiz-López. Simulation of liquid water using semiempirical Hamiltonians and the divide and conquer approach.
- [177] Soohaeng Yoo, Xiao Cheng Zeng, and Sotiris S. Xantheas. On the phase diagram of water with density functional theory potentials: The melting temperature of ice I_h with the Perdew-Burke-Ernzerhof and Becke-Lee-Yang-Parr functionals. *J. Chem. Phys.*, 130:221102, 2009.
- [178] Edina Rosta, Marcin Nowotny, Wei Yang, and Gerhard Hummer. Catalytic mechanism of RNA backbone cleavage by ribonuclease H from quantum mechanics/molecular mechanics simulations. *J. Am. Chem. Soc.*, 133:8934–8941, 2011.
- [179] Edina Rosta, H. Lee Woodcock, Bernard R. Brooks, and Gerhard Hummer. Artificial reaction coordinate “tunneling” in free-energy calculations: The catalytic reaction of RNase H. *J. Comput. Chem.*, 30:1634–1641, 2009.
- [180] Paul Maragakis, Arjan van der Vaart, and Martin Karplus. Gaussian-mixture umbrella sampling. *J. Phys. Chem. B*, 113:4664–4673, 2009.
- [181] S. J. Sheather and M. C. Jones. A reliable data-based bandwidth selection method for kernel density estimation. *J. R. Statist. Soc. B*, 53:683–690, 1991.
- [182] William Humphrey, Andrew Dalke, and Klaus Schulten. VMD: Visual Molecular Dynamics. *J. Mol. Graphics*, 14:33–38, 1996.
- [183] M. J. Frisch, G. W. Trucks, H. B. Schlegel, G. E. Scuseria, M. A. Robb, J. R. Cheeseman, G. Scalmani, V. Barone, B. Mennucci, G. A. Petersson, H. Nakatsuji, M. Caricato, X. Li, H. P. Hratchian, A. F. Izmaylov, J. Bloino, G. Zheng, M. Sonnenberg, M. Hada, M. Ehara, K. Toyota, R. Fukuda, J. Hasegawa, M. Ishida, T. Nakajima, Y. Honda, O. Kitao, H. Nakai, T. Vreven, J. A. Montgomery, J., J. E. Peralta, F. Ogliaro, M. Bearpark, J. J. Heyd, E. Brothers, K. N. Kudin,

- V. N. Straverov, R. Kobayashi, J. Normand, K. Raghavachari, A. Rendell, J. C. Burant, S. S. Iyengar, J. Tomasi, M. Cossi, N. Rega, J. M. Millam, M. Klene, J. E. Knox, J. B. Cross, V. Bakken, C. Adamo, J. Jaramillo, R. Gomperts, R. E. Stratmann, O. Yazyev, A. J. Austin, R. Cammi, C. Pomelli, J. W. Ochterski, R. L. Martin, K. Morokuma, V. G. Zakrzewski, G. A. Voth, P. Salvador, J. J. Dannenberg, S. Dapprich, A. D. Daniels, O. Farkas, J. B. Foresman, J. V. Ortiz, J. Cioslowski, and D. J. Fox. *Gaussian 09, Revision A.02*. Gaussian, Inc., Wallingford, CT, 2009.
- [184] Johannes Kästner, Joanne M. Carr, Thomas W. Keal, Walter Thiel, Adrian Wander, and Paul Sherwood. DL-FIND: An open-source geometry optimizer for atomistic simulations. *J. Phys. Chem. A*, 113:11856–11865, 2009.
- [185] Lars Ridder and Adrian J. Mulholland. Modeling biotransformation reactions by combined quantum mechanical/molecular mechanical approaches: From structure to activity. *Curr. Top. Med. Chem.*, 3:1241–1256, 2003.
- [186] Pascal Auffinger, Thomas E. Cheatham III, and Andrea C. Vaiana. Spontaneous formation of KCl aggregates in biomolecular simulations: A force field issue? *J. Chem. Theory Comput.*, 3:1851–1859, 2007.
- [187] Alan A. Chen and Rohit V. Pappu. Parameters of monovalent ions in the AMBER-99 forcefield: Assessment of inaccuracies and proposed improvements. *J. Phys. Chem. B*, 111:11884–11887, 2007.
- [188] Paul Haake and Richard V. Prigodich. Method for determination of phosphate anion-cation association constants from ^{31}P chemical shifts. *Inorg. Chem.*, 23:457–462, 1984.
- [189] P. A. Bash, L. L. Ho, Alexander D. MacKerell, Jr., D. Levine, and P. Hallstrom. Progress toward chemical accuracy in computer simulation of condensed phase reactions. *Proc. Natl. Acad. Sci. USA*, 93:3698–3703, 1996.
- [190] Yaoquan Tu and Aatto Laaksonen. On the effect of Lennard-Jones parameters on the quantum mechanical and molecular mechanical coupling in a hybrid molecular dynamics simulation of liquid water. *J. Chem. Phys.*, 111:7519–7525, 1999.

- [191] Ulla Pentikäinen, Katherine E. Shaw, Kittusamy Senthilkumar, Christopher J. Woods, and Adrian J. Mulholland. Lennard-Jones parameters for B3LYP/CHARMM27 QM/MM modeling of nucleic acid bases. *J. Chem. Theory Comput.*, 5:396–410, 2009.
- [192] Katherine E. Shaw, Christopher J. Woods, and Adrian J. Mulholland. Compatibility of quantum chemical methods and empirical (MM) water models in quantum mechanics/molecular mechanics liquid water simulations. *J. Phys. Chem. Lett.*, 1:219–223, 2010.
- [193] Demian Riccardi, Guohui Li, and Qiang Cui. Importance of van der Waals interactions in QM/MM simulations. *J. Phys. Chem. B*, 108:6467–6478, 2004.
- [194] Timothy J. Giese and Darrin M. York. Charge-dependent model for many-body polarization, exchange, and dispersion interactions in hybrid quantum mechanical/molecular mechanical calculations. *J. Chem. Phys.*, 127:194101, 2007.
- [195] Timothy J. Giese and Darrin M. York. Density-functional expansion methods: Grand challenges. *Theor. Chem. Acc.*, 131:1145, 2012.
- [196] Adam Moser, Kevin Range, and Darrin M. York. Accurate proton affinity and gas-phase basicity values for molecules important in biocatalysis. *J. Phys. Chem. B*, 114:13911–13921, 2010.
- [197] Barbara Gerratana, Gwendolyn A. Sowa, and W. W. Cleland. Characterization of the transition-state structures and mechanisms for the isomerization and cleavage reactions of uridine 3'-*m*-nitrobenzyl phosphate. *J. Am. Chem. Soc.*, 122:12615–12621, 2000.
- [198] Martha J. Fedor. Comparative enzymology and structural biology of RNA self-cleavage. *Annu. Rev. Biophys.*, 38:271–299, 2009.
- [199] Joachim Schnabl and Roland K. O. Sigel. Controlling ribozyme activity by metal ions. *Curr. Opin. Chem. Biol.*, 14:269–275, 2010.
- [200] Gerard C. L. Wong and Lois Pollack. Electrostatics of strongly charged biological

- polymers: Ion-mediated interactions and self-organization in nucleic acids and proteins. *Annu. Rev. Phys. Chem.*, 61:171–189, 2010.
- [201] Roland K O Sigel, Pascal Auffinger, Shi-Jie Chen, Jrg S Hartig, Wade C Winkler, Nils G Walter, Joseph A Piccirilli, Samuel E Butcher, Norbert Polacek, Hiroaki Suga, Joseph E Wedekind, and Victoria DeRose. *Structural and Catalytic Roles of Metal Ions in RNA*. Metal Ions in Life Sciences. The Royal Society of Chemistry, 2011.
- [202] Barbara L. Golden, Sharon Hammes-Schiffer, Paul R. Carey, and Philip C. Bevilacqua. *An integrated picture of HDV ribozyme catalysis*, volume 3 of *Bio-physics for the Life Sciences*, chapter 8, pages 135–167. Springer, New York, 2013.
- [203] W Luke Ward, Kory Plakos, and Victoria J. DeRose. Nucleic acid catalysis: Metals, nucleobases, and other cofactors. *Chem. Rev.*, 114:4318–4342, 2014.
- [204] Joseph W. Cottrell, Lincoln G. Scott, and Martha J. Fedor. The pH dependence of hairpin ribozyme catalysis reflects ionization of an active site adenine. *J. Biol. Chem.*, 286:17658–17664, 2011.
- [205] Tai-Sung Lee, Carlos Silva-Lopez, Monika Martick, William G. Scott, and Darrin M. York. Insight into the role of Mg^{2+} in hammerhead ribozyme catalysis from x-ray crystallography and molecular dynamics simulation. *J. Chem. Theory Comput.*, 3:325–327, 2007.
- [206] Tai-Sung Lee, George M Giambaşu, Carlos P Sosa, Monika Martick, William G Scott, and Darrin M York. Threshold occupancy and specific cation binding modes in the hammerhead ribozyme active site are required for active conformation. *J. Mol. Biol.*, 388:195–206, 2009.
- [207] Ian T. Suydam, Stephen D. Levandoski, and Scott A. Strobel. Catalytic importance of a protonated adenosine in the hairpin ribozyme active site. *Biochemistry*, 49:3723–3732, 2010.
- [208] Fabrice Leclerc. Hammerhead ribozymes: True metal or nucleobase catalysis? where is the catalytic power from? *Molecules*, 15:5389–5407, 2010.

- [209] Stephanie Kath-Schorr, Timothy J. Wilson, Nan-Sheng Li, Jun Lu, Joseph A. Piccirilli, and David M. J. Lilley. General acid-base catalysis mediated by nucleobases in the hairpin ribozyme. *J. Am. Chem. Soc.*, 134:16717–16724, 2012.
- [210] M. Y. Kuo, L. Sharmeen, G. Dinter-Gottlieb, and J. Taylor. Characterization of self-cleaving RNA sequences on the genome and antigenome of human hepatitis delta virus. *J. Virol.*, 62:4439–4444, 1988.
- [211] Michael M. C. Lai. The molecular biology of hepatitis delta virus. *Annu. Rev. Biochem.*, 64:259–286, 1995.
- [212] Chiu-Ho T. Webb and Andrej Lupták. HDV-like self-cleaving ribozymes. *RNA Biol.*, 8:719–727, 2011.
- [213] Huey-Nan Wu, Yu-June Lin, Fu-Pang Lin, Shinji Makino, and Ming-Fu Chang. Human hepatitis δ virus RNA subfragments contain an autocleavage activity. *Proc. Natl. Acad. Sci. USA*, 86:1831–1835, 1989.
- [214] A. T. Perrotta and M. D. Been. The self-cleaving domain from the genomic RNA of hepatitis delta virus: Sequence requirements and the effects of denaturant. *Nucleic Acids Res.*, 18:6821–6827, 1990.
- [215] T. S. Wadkins, A. T. Perrotta, A. R. Ferré-D’Amaré, J. A. Doudna, and M. D. Been. A nested double pseudoknot is required for self-cleavage activity of both the genomic and antigenomic hepatitis delta virus ribozymes. *RNA*, 5:720–727, 1999.
- [216] Shu-ichi Nakano, Durga M. Chadalavada, and Philip C. Bevilacqua. General acid-base catalysis in the mechanism of a hepatitis delta virus ribozyme. *Science*, 287:1493–1497, 2000.
- [217] Sarah P. Rosenstein and Michael D. Been. Self-cleavage of hepatitis delta virus genomic strand RNA is enhanced under partially denaturing conditions. *Biochemistry*, 29:8011–8016, 1990.
- [218] Michael D. Been, Anne T. Perrotta, and Sarah P. Rosenstein. Secondary structure

- of the self-cleaving RNA of hepatitis delta virus: Applications to catalytic RNA design. *Biochemistry*, 31:11843–11852, 1992.
- [219] Gilbert Thill, Marc Vasseur, and N. Kyle Tanner. Structural and sequence elements required for the self-cleaving activity of the hepatitis delta virus ribozyme. *Biochemistry*, 32:4254–4262, 1993.
- [220] Shu-ichi Nakano, David J. Proctor, and Philip C. Bevilacqua. Mechanistic characterization of the HDV genomic ribozyme: Assessing the catalytic and structural contributions of divalent metal ions within a multichannel reaction mechanism. *Biochemistry*, 40:12022–12038, 2001.
- [221] Anne T. Perrotta and Michael D. Been. HDV ribozyme activity in monovalent cations. *Biochemistry*, 45:11357–11365, 2006.
- [222] Y. A. Suh, P. K. Kumar, K. Taira, and S. Nishikawa. Self-cleavage activity of the genomic HDV ribozyme in the presence of various divalent metal ions. *Nucleic Acids Res.*, 21:3277–3280, 1993.
- [223] Hamid Fauzi, Junji Kawakami, Fumiko Nishikawa, and Satoshi Nishikawa. Analysis of the cleavage reaction of a *trans*-acting human hepatitis delta virus ribozyme. *Nucleic Acids Res*, 25:3124–3130, 1997.
- [224] A. K. Oyelere and S. A. Strobel. Site specific incorporation of 6-azauridine into the genomic HDV ribozyme active site. *Nucleosides Nucleotides Nucleic Acids*, 20:1851–1858, 2001.
- [225] Adegboyega K. Oyelere, Julia R. Kardon, and Scott A. Strobel. pK_a perturbation in genomic hepatitis delta virus ribozyme catalysis evidenced by nucleotide analogue interference mapping. *Biochemistry*, 41:3667–3675, 2002.
- [226] Subha Das and Joseph Piccirilli. General acid catalysis by the hepatitis delta virus ribozyme. *Nature Chem. Biol.*, 1:45–52, 2005.
- [227] Anne T. Perrotta, I-hung Shih, and Michael D. Been. Imidazole rescue of a cytosine mutation in a self-cleaving ribozyme. *Science*, 286:123–126, 1999.

- [228] Ailong Ke, Kaihong Zhou, Fang Ding, Jamie H. D. Cate, and Jennifer A. Doudna. A conformational switch controls hepatitis delta virus ribozyme catalysis. *Nature*, 429:201–205, 2004.
- [229] Jui-Hui Chen, Bo Gong, Philip C. Bevilacqua, Paul R. Carey, and Barbara L. Golden. A catalytic metal ion interacts with the cleavage site GU wobble in the HDV ribozyme. *Biochemistry*, 48:1498–1507, 2009.
- [230] Bo Gong, Jui-Hui Chen, Elaine Chase, Durga M. Chadalavada, Rieko Yajima, Barbara L. Golden, Philip C. Bevilacqua, and Paul R. Carey. Direct measurement of a pK_a near neutrality for the catalytic cytosine in the genomic HDV ribozyme using Raman crystallography. *J. Am. Chem. Soc.*, 129:13335–13342, 2007.
- [231] Jui-Hui Chen, Rieko Yajima, Durga M. Chadalavada, Elaine Chase, Philip C. Bevilacqua, and Barbara L. Golden. A 1.9 Å crystal structure of the HDV ribozyme precleavage suggests both Lewis acid and general acid mechanisms contribute to phosphodiester cleavage. *Biochemistry*, 49:6508–6518, 2010.
- [232] Adrian R. Ferré-D’Amaré, Kaihong Zhou, and Jennifer A. Doudna. Crystal structure of a hepatitis delta virus ribozyme. *Nature*, 395:567–574, 1998.
- [233] Ailong Ke, Fang Ding, Joseph D. Batchelor, and Jennifer A. Doudna. Structural roles of monovalent cations in the HDV ribozyme. *Structure*, 15:281–287, 2007.
- [234] Barbara L. Golden. Two distinct catalytic strategies in the hepatitis delta virus ribozyme cleavage reaction. *Biochemistry*, 50:9424–9433, 2011.
- [235] Narayanan Veeraraghavan, Abir Ganguly, Jui-Hui Chen, Philip C. Bevilacqua, Sharon Hammes-Schiffer, and Barbara L. Golden. Metal binding motif in the active site of the HDV ribozyme binds divalent and monovalent ions. *Biochemistry*, 50:2672–2682, 2011.
- [236] Narayanan Veeraraghavan, Abir Ganguly, Barbara L. Golden, Philip C. Bevilacqua, and Sharon Hammes-Schiffer. Mechanistic strategies in the HDV ribozyme: Chelated and diffuse metal ion interactions and active site protonation. *J. Phys. Chem. B*, 115:8346–8357, 2011.

- [237] Ji Chen, Abir Ganguly, Zulaika Miswan, Sharon Hammes-Schiffer, Philip C. Bevilacqua, and Barbara L. Golden. Identification of the catalytic Mg^{2+} ion in the hepatitis delta virus ribozyme. *Biochemistry*, 52:557–567, 2013.
- [238] Dominique Lévesque, Cédric Reymond, and Jean-Pierre Perreault. Characterization of the trans Watson-Crick GU base pair located in the catalytic core of the antigenomic HDV ribozyme. *PLoS ONE*, 7:40309, 2012.
- [239] S. R. Lynch and I Tinoco Jr. The structure of the L3 loop from the hepatitis delta virus ribozyme: A *syn* cytidine. *Nucleic Acids Res.*, 26:980–987, 1998.
- [240] Durga M. Chadalavada and Scott M. Knudsen and Shu-ichi Nakano and Philip C. Bevilacqua. A role for upstream RNA structure in facilitating the catalytic fold of the genomic hepatitis delta virus ribozyme. *J. Mol. Biol.*, 301:349–367, 2000.
- [241] Susan E. Senchak Durga M. Chadalavada and Philip C. Bevilacqua. The folding pathway of the genomic hepatitis delta virus ribozyme is dominated by slow folding of the pseudoknots. *J. Mol. Biol.*, 317:559–575, 2002.
- [242] Trevor S. Brown, Durga M. Chadalavada, and Philip C. Bevilacqua. Design of a highly reactive HDV ribozyme sequence uncovers facilitation of RNA folding by alternative pairings and physiological ionic strength. *J. Mol. Biol.*, 341:695–712, 2004.
- [243] Bo Gong, Jui-Hui Chen, Philip C. Bevilacqua, Barbara L. Golden, and Paul Richard Carey. Competition between $\text{Co}(\text{NH}_3)_6^{3+}$ and inner sphere Mg^{2+} ions in the HDV ribozyme. *Biochemistry*, 48:11961–11970, 2009.
- [244] Tai-Sung Lee, George M. Giambaşu, Michael E. Harris, and Darrin M. York. Characterization of the structure and dynamics of the HDV ribozyme in different stages along the reaction path. *J. Phys. Chem. Lett.*, 2:2538–2543, 2011.
- [245] I-hung Shih and Michael D. Been. Involvement of a cytosine side chain in proton transfer in the rate-determining step of ribozyme self-cleavage. *Proc. Natl. Acad. Sci. USA*, 98:1489–1494, 2001.

- [246] Andrej Lupták, Adrian R. Ferré-D'Amaré, Kaihong Zhou, Kurt W. Zilm, and Jennifer A. Doudna. Direct pK_a measurement of the active-site cytosine in a genomic hepatitis delta virus ribozyme. *J. Am. Chem. Soc.*, 123:8447–8452, 2001.
- [247] Anne T. Perrotta, Timothy S. Wadkins, and Michael D. Been. Chemical rescue, multiple ionizable groups, and general acid-base catalysis in the HDV genomic ribozyme. *RNA*, 12:1282–1291, 2006.
- [248] Shu-ichi Nakano, Andrea L. Cerrone, and Philip C. Bevilacqua. Mechanistic characterization of the HDV genomic ribozyme: Classifying the catalytic and structural metal ion sites within a multichannel reaction mechanism. *Biochemistry*, 42:2982–2994, 2003.
- [249] Y. H. Jeoung, P. K. Kumar, Y. A. Suh, K. Taira, and S. Nishikawa. Identification of phosphate oxygens that are important for self-cleavage activity of the HDV ribozyme by phosphorothioate substitution interference analysis. *Nucleic Acids Res*, 22:3722–3727, 1994.
- [250] Nita S. Prabhu, Gail Dinter-Gottlieb, and Philip A. Gottlieb. Single substitutions of phosphorothioates in the HDV ribozyme G73 define regions necessary for optimal self-cleaving activity. *Nucleic Acids Res.*, 25:5119–5124, 1997.
- [251] Kristen Raines and Philip A. Gottlieb. Enzymatic incorporation of 2'-thio-CTP into the HDV ribozyme. *RNA*, 4:340–345, 1998.
- [252] Pallavi Thaplyal, Abir Ganguly, Barbara L. Golden, Sharon Hammes-Schiffer, and Philip C. Bevilacqua. Thio effects and an unconventional metal ion rescue in the genomic hepatitis delta virus ribozyme. *Biochemistry*, 52:6499–6514, 2013.
- [253] John K. Frederiksen and Joseph A. Piccirilli. Identification of catalytic metal ion ligands in ribozymes. *Methods*, 49:148–166, 2009.
- [254] Abir Ganguly, Philip C. Bevilacqua, and Sharon Hammes-Schiffer. Quantum mechanical/molecular mechanical study of the HDV ribozyme: Impact of the catalytic metal ion on the mechanism. *The Journal of Physical Chemistry Letters*, 2:2906–2911, 2011.

- [255] Abir Ganguly, Pallavi Thaplyal, Edina Rosta, Philip C. Bevilacqua, and Sharon Hammes-Schiffer. Quantum mechanical/molecular mechanical free energy simulations of the self-cleavage reaction in the hepatitis delta virus ribozyme. *J. Am. Chem. Soc.*, 136:1483–1496, 2014.
- [256] Romelia Salomon-Ferrer, Andreas W. Götz, Duncan Poole, Scott Le Grand, and Ross C. Walker. Routine microsecond molecular dynamics simulations with AMBER on GPUs. 2. Explicit solvent particle mesh Ewald. *J. Chem. Theory Comput.*, 9:3878–3888, 2013.
- [257] Evelyn Mayaan, Kevin Range, and Darrin M. York. Structure and binding of Mg(II) ions and di-metal bridge complexes with biological phosphates and phosphoranes. *J. Biol. Inorg. Chem.*, 9:807–817, 2004.
- [258] Nathan A. Baker, David Sept, Simpson Joseph, Michael J. Holst, and J. Andrew McCammon. Electrostatics of nanosystems: Application to microtubules and the ribosome. *Proc. Natl. Acad. Sci. USA*, 98:10037–10041, 2001.
- [259] Robert Konecny, Nathan A. Baker, and J. Andrew McCammon. iAPBS: A programming interface to the adaptive Poisson-Boltzmann solver. *Comput. Sci. Disc.*, 5:015005, 2012.
- [260] Samuel Genheden, Tyler Luchko, Sergey Gusarov, Andriy Kovalenko, and Ulf Ryde. An MM/3D-RISM approach for ligand binding affinities. *J. Phys. Chem. B*, 114:8505–8516, 2010.
- [261] Jean-François Truchon, B. Montgomery Pettitt, and Paul Labute. A cavity corrected 3D-RISM functional for accurate solvation free energies. *J. Chem. Theory Comput.*, 10:934–941, 2014.
- [262] In Suk Joung, Tyler Luchko, and David A. Case. Simple electrolyte solutions: Comparison of DRISM and molecular dynamics results for alkali halide solutions. *J. Chem. Phys.*, 138:044103, 2013.
- [263] George M. Giambaşu, Tyler Luchko, Daniel Herschlag, Darrin M. York, and David A. Case. Ion counting from explicit-solvent simulations and 3D-RISM. *Biophys. J.*, 106:883–894, 2014.

- [264] Irina Velikyan, Sandipta Acharya, Anna Trifonova, Andras Földesi, and Jyoti Chattopadhyaya. The pK_a 's of 2'-hydroxyl group in nucleosides and nucleotides. *J. Am. Chem. Soc.*, 123:2893–2894, 2001.
- [265] Jing-Dong Ye, Nan-Sheng Li, Qing Dai, and Joseph A. Piccirilli. The mechanism of RNA strand scission: An experimental measure of the Brønsted coefficient, β_{nuc} . *Angew. Chem. Int. Ed.*, 119:3788–3791, 2007.
- [266] Andrew Pohorille, Christopher Jarzynski, and Christophe Chipot. Good practices in free-energy calculations. *J. Phys. Chem. B*, 114:10235–10253, 2010.
- [267] Fangqiang Zhu and Gerhard Hummer. Convergence and error estimation in free energy calculations using the weighted histogram analysis method. *J. Comput. Chem.*, 33:453–465, 2012.
- [268] Mark N. Kobra. Systematic and statistical error in histogram-based free energy calculations. *J. Comput. Chem.*, 24:1437–1446, 2003.
- [269] Daniel R. Roe and Thomas E. Cheatham III. PTRAJ and CPPTRAJ: Software for processing and analysis of molecular dynamics trajectory data. *J. Chem. Theory Comput.*, 9:3084–3095, 2013.
- [270] Vinod K. Misra and David E. Draper. A thermodynamic framework for Mg^{2+} binding to RNA. *Proc. Natl. Acad. Sci. USA*, 98:12456–12461, 2001.
- [271] Vojtěch Mlýnský, Pavel Banáš, Jiří Šponer, Marc W. van der Kamp, Adrian J. Mulholland, and Michal Otyepka. Comparison of *ab initio*, DFT, and semiempirical QM/MM approaches for description of catalytic mechanism of hairpin ribozyme. *J. Chem. Theory Comput.*, 10:1608–1622, 2014.

Appendix A

Supporting Information for: Molecular Simulations of RNA 2'-*O*-Transesterification Reaction Models in Solution

A.1 Umbrella Sampling

A.1.1 Detailed Equilibration Protocol

Simulations were set up by first selecting two quantum mechanical gas phase minima of each of the solutes, one corresponding to an appropriate reactant state and the other to separated products. These structures were then solvated in a rhombic dodecahedron composed of 2640 water molecules. For simulations with ions, water molecules were randomly converted into sodium and chloride ions until the system was net neutral at 140 mM when considering only the water molecules at a density of 1.0 g/cm³. This is equivalent to assuming that the neutralized solute (*i.e.* as if it were a salt) was placed into a 140 mM bulk ionic environment with no additional changes to the system volume.

The “equilibration” of each system was performed in two stages: relatively long trajectories with only molecular mechanical (MM) interactions to account for slower ion motions and shorter quantum mechanical/molecular mechanical (QM/MM) trajectories

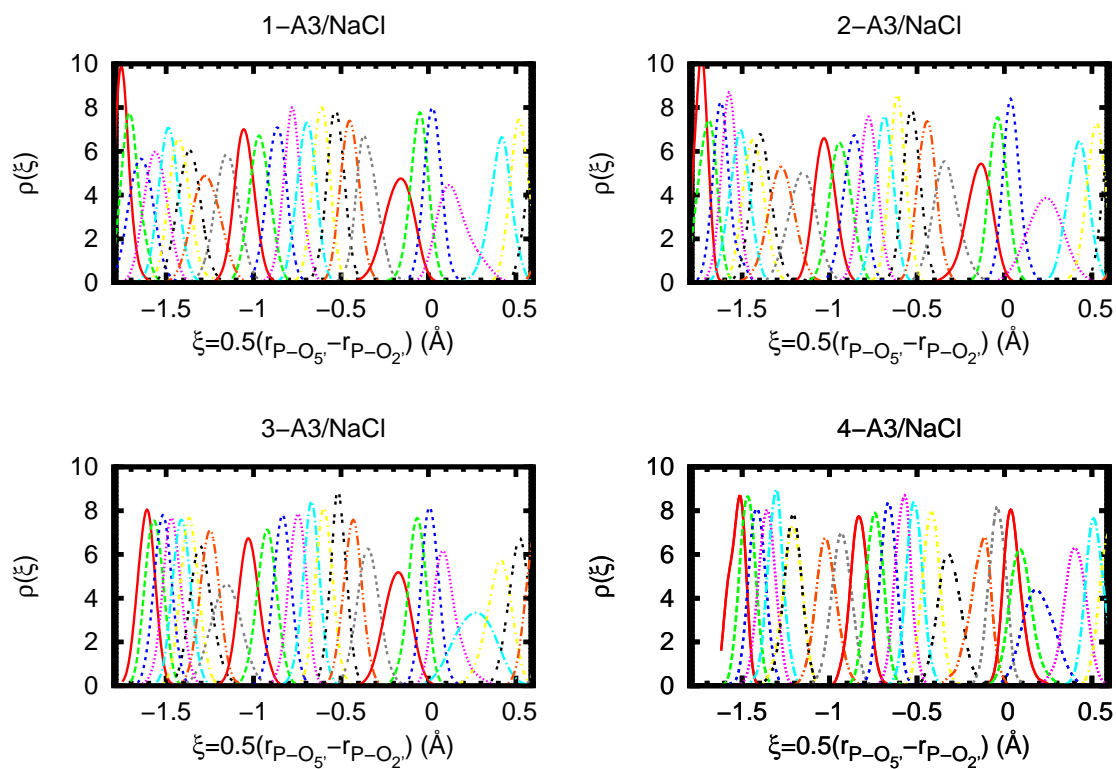
to relax the density and internal solute geometries. For the first stage, the solute atoms were assigned fixed charges corresponding to their gas phase Mulliken charges and harmonically restrained at their initial coordinates [$U_{\text{restraint}}(\mathbf{x}) = k(\mathbf{x} - \mathbf{x}_{\text{initial}})^2$, $k = 100$ kcal/mol-Å²]. All bonds to hydrogen, including those in the solute, were constrained at their equilibrium geometries. The system was minimized, followed by 1 ns of dynamics with temperature and pressure coupling and 10 ns with only temperature coupling. In the second stage, the solute was briefly minimized to remove minor artifacts of the MM description, followed by 200 ps of dynamics with pressure and temperature coupling. The unit cell volumes from the last 50 ps of the corresponding reactant and product state simulations were then averaged, lattice vectors were reassigned, and another 400 ps of constant volume and temperature dynamics was performed on both systems. Throughout the QM/MM equilibration trajectories, weak harmonic restraints along key geometric quantities were applied to maintain in-line attack of the nucleophile in the reactant and to prevent the quantum region from becoming too large due to separation of the products.

During production, 24 or 32 “windows” were defined by a harmonic biasing potential, $U_{\text{bias}}(\xi) = k(\xi - \xi_0)^2$ ($k = 80$ or 120 kcal/mol-Å²), whose centers, ξ_0 , were placed no more than 0.1 Å apart. In order to minimize correlation between windows, starting structures were randomly assigned at regular intervals from the last 200 ps of the nearest endpoint equilibration (reactant or products), which were evenly divided.

A.1.2 Overlap of Reaction Coordinate Distributions

It is well known that overlap of the potential energy (or reaction coordinate) distributions is essential to accurate free energy calculations[266], especially in the calculation of free energy profiles[267]. Simple visual inspection of the distributions observed here beget significant confidence that the simulations are well converged with a high degree of overlap (see A.1).

Figure A.1: Representative plots of reaction coordinate distributions from umbrella sampling simulations of different reaction models (models n-A3/NaCl in the text, where $n = 1-4$).



A.1.3 Lennard-Jones Parameters

Sodium Parameters

During preliminary simulations it was observed that sodium ions would occasionally approach and bind, apparently irreversibly, to the dianionic quantum region. This was deemed problematic for two reasons. First, the present work is not aimed at studying the direct participation of metal ions in phosphoryl transfer; such events would presumably be better modeled if the ions (and perhaps even some of the solvation shell) were treated quantum mechanically. Second, experiments show that the binding coefficients for sodium and phosphates are only of the order 10^{-1} and 10^1 for mono- and dianionic species respectively[188]. Binding of ions over the majority of the simulation is therefore indicative of either physical inaccuracy (the model predicts an incorrect binding free energy) or an improper phase space sampling (simulation does not properly sample bound and unbound configurations). Both of these issues can be (non-optimally) circumvented by effectively augmenting the phase space to exclude bound configurations. One way to accomplish this is to set the Lennard-Jones interactions between the solute and sodium ions such that: 1) at some distance, $r_{k_B T}$, near the solvation shell boundary the potential between atoms is close in energy to the thermal fluctuations of the system (*i.e.* $U_{LJ}(r_{k_B T}) = k_B T$) and 2) beyond this distance the potential is either negligible or only weakly attractive. The Lennard-Jones potential is thus considered:

$$U_{LJ}(r) = \epsilon \left[\left(\frac{R_{\min}}{r} \right)^{12} - 2 \left(\frac{R_{\min}}{r} \right)^6 \right] = \frac{A}{r^{12}} - \frac{B}{r^6}$$

Setting $U_{LJ}(r_{k_B T}) = k_B T$ and solving for positive solutions of R_{\min} gives

$$R_{\min} = r_{k_B T} \left[1 + \left(1 + \frac{k_B T}{\epsilon} \right)^{\frac{1}{2}} \right]^{\frac{1}{6}}$$

Since an infinite number of R_{\min} and ϵ pairs can meet the first condition, the second condition is arbitrarily satisfied by setting $\epsilon = 0.01$ kcal/mol. The values for $r_{k_B T}$ are then taken to correspond to 2.5 times the location of the first peak of the radial distribution function of sodium and the oxygens of water. During the parameterization of several alkali metal and halide ions Joung and Cheatham calculated these peaks, at

300K, to be 2.38 Å and 2.35 Å for TIP3P and TIP4P-Ew respectively (the difference in temperature is assumed to be negligible)[169]. Inserting the appropriate values into the equation for R_{\min} gives 8.5675 Å for TIP3P and 8.4596 Å for TIP4P-Ew systems ($\epsilon = 0.01$ kcal/mol for both). For all of the simulations reported here, these values are used to calculate the Lennard-Jones interaction between sodium ions and all carbon and oxygen atoms in the solute.

Quantum Region Parameters

For all systems, chemical environments were assigned to be consistent with the reactant state. Since both force fields use different atom typing schemes, there is not always a one to one parameter correspondence. For example oxygen atoms in hydroxyl groups and esters are different types in the AMBER force field, but are the same type in the CHARMM force field. Both the CHARMM and AMBER force fields do, however, use the same mixing rules for interactions involving atoms of different types. That is, for two atom types i and j , the combined interaction parameters are:

$$\epsilon_{ij} = (\epsilon_i \epsilon_j)^{\frac{1}{2}} \quad \text{and} \quad R_{\min,ij} = \frac{R_{\min,i} + R_{\min,j}}{2},$$

or equivalently:

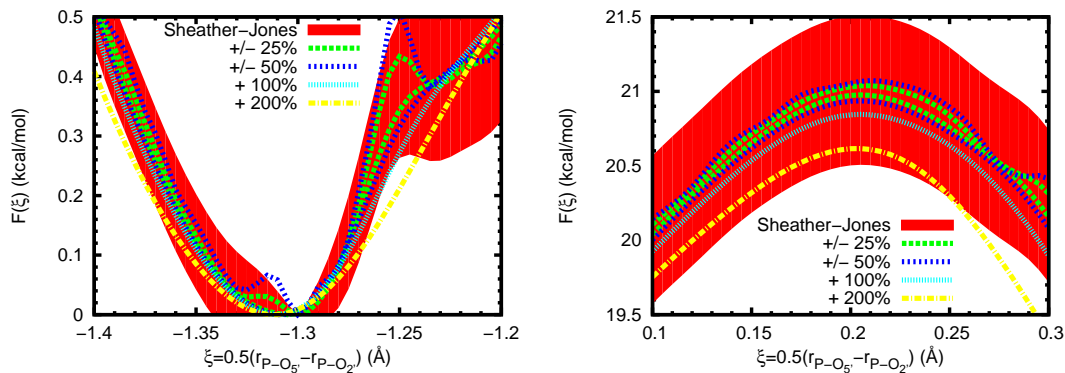
$$A_{ij} = \epsilon_{ij} R_{\min,ij}^{12} \quad \text{and} \quad B_{ij} = 2\epsilon_{ij} R_{\min,ij}^6.$$

The ϵ/R_{\min} representation is usually preferred in the definition of forcefields since mixed atom interactions are then readily formulated (*i.e.* there is no direct mixing rule for A and B parameters), while the A/B representation is more convenient for comparing repulsive and attractive components. It is also important to note that an infinite number of ϵ/R_{\min} pairs can give rise to a particular A/B pair, but the mixed atom interactions will be different in many of the pairs.

Table A.1: Lennard-Jones parameters for atoms in the quantum region. R_{\min} , ϵ , A , and B values are in \AA , kcal/mol, kcal- \AA^{12} /mol, and kcal- \AA^6 /mol respectively.

atom	environment	AMBER FF10				CHARMM27			
		$R_{\min}/2$	ϵ	$A \times 10^{-3}$	B	$R_{\min}/2$	ϵ	$A \times 10^{-3}$	B
O	hydroxyl	1.7210	0.2104	581.80	699.75	1.7700	0.1521	589.07	598.66
	esters	1.6837	0.1700	361.40	495.73				
	non-bridge	1.6612	0.2100	379.88	564.89				
P		2.1000	0.2000	6025.9	2195.6	2.1500	0.5850	23376.0	7396.0
C	ribose ring	1.9080	0.1094	1043.1	675.61	2.2750	0.0200	1574.6	354.92
	5' methylene					2.0100		997.47	472.69
H	hydroxyl	0.0000	0.0000	0.0	0.0000	0.2245	0.0460	3.1×10^{-9}	7.5×10^{-4}
	alkyl	1.4870	0.0157	7.5161	21.726	1.3400	0.0240	3.2948	17.785
	vicinal oxygen	1.3870	0.0157	3.2597	14.308	1.3200	0.0220	2.5216	14.896
	5' methylene					1.3400	0.0280	3.8439	20.749

Figure A.2: Plots demonstrating the effect of increasing/decreasing the bandwidths on the calculated free energy profile (data correspond to model 4-A3/NaCl in the text). Extreme close ups of both the reactant (left) and transition state (right) show smoothing errors quite small compared to the statistical errors. The shaded region represents an estimated 95% confidence interval relative to the reactant state.



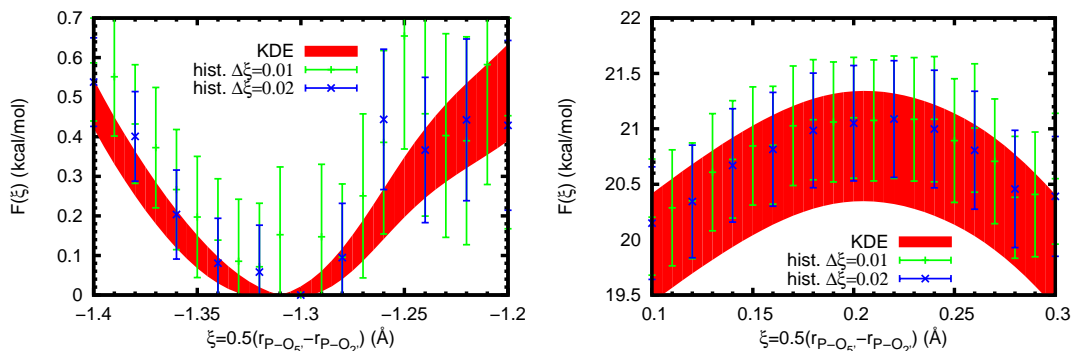
A.2 Kernel Density Estimation

Although kernel density estimation is a fairly mature field (Ref. 125 and 126 are both excellent overviews), its use in the chemical literature is perhaps somewhat novel. The astute reader will notice immediate analogy to histogramming, which is extremely common in the computation of free energy profiles[68, 70] and for which the bias/variance tradeoff is well established[268]. Here we briefly examine the effects of bandwidth selection, which plays the same role as bin width selection in histogramming, and briefly compare to the results obtained with histograms.

A.2.1 Effect of Bandwidth Choice

Many methods exist for optimal selection of the bandwidth[126], although extension to multiple dimensions still seems to be a challenge. In one dimension, the method of Sheather and Jones seems to be most consistently optimal amongst different types of probability distributions, but for simple, unimodal distributions the differences are often small[181]. Taking the Sheather-Jones approach as a base, we examine the effects of increasing (“over-smoothing”) and decreasing (“under-smoothing”) by a fixed percentage in each umbrella sampling simulation. The results show that the differences

Figure A.3: Plots comparing the use of a normal kernel density estimator versus a more traditional histogram estimator when calculating a free energy profile (data correspond to model 4-A3/NaCl in the text). Extreme close ups of both the reactant (left) and transition state (right) show that the kernel estimator has comparably smoother derivatives and more well defined minima and maxima. The shaded region and error bars represent an estimated 95% confidence interval relative to the reactant state.



in energy are extremely difficult to discern when considered on a scale comparable to the statistical error, even when the bandwidth is increased by as much as 100% (A.2, close attention to the axes is necessary). This would suggest that even crude, over-smoothing estimates of the bandwidths would still give energies with acceptable errors. Because less smooth estimates can have undesirable physical properties (*e.g.* spurious minima/maxima), this over-smoothing approach is recommended. Additional advantages may be found when, for example, attempting to optimize biasing forces, since in this case smoothing could also improve numerical stability. However, it should be mentioned that the systematic error in the progress coordinate values is of the order of the bandwidth used (0.02 - 0.05 in this work) and should be considered, for example, when drawing conclusions regarding the location of stationary points.

A.2.2 Comparison with Histograms

Due to the dominance of histogram estimators in the chemical literature, it would appear informative to directly compare them with kernel density estimators (KDEs), especially in the present context of calculating free energy profiles. The primary theoretical advantage of KDEs is that they provide smoother estimates with well defined derivatives (this is not formally the case for histograms). The estimate is also defined over the

full range of data, as opposed to a series of discrete “bins” and therefore the degree of smoothing is not linked to the resolution of the estimate. For example, in the left side of A.3 the spurious “bump” near -1.25 can be removed by doubling the bin width at the cost of halving the resolution. Presumably the lower resolution estimate introduces bias into the location of extrema (although in this case the minima are coincidentally located for both bin widths). Interestingly, the same “bump” is observed in the “under-smoothed” estimates in A.2, but in this case increasing the degree of smoothing does not as severely augment the result. Lastly, it is comparatively simple to automatically determine extrema when using KDEs, both because the KDE is able to remove more statistical noise and because the derivative can be directly determined as well, although this was not done in the present work. If a minimum/maximum is naively defined as having both adjacent points be higher/lower in free energy, then a pass through the raw data for the three profiles in A.3 discovers two minima and one maxima in the kernel estimate as opposed to seven minima and 6 maxima in the low resolution estimate (bin width of 0.02 \AA) and a befuddling 31 minima and 32 maxima at higher resolution (bin width of 0.01 \AA).

Appendix B

Supporting Information for: A Framework for Assessment of Metal-Assisted Nucleophile Activation in the Hepatitis Delta Virus Ribozyme

B.1 Molecular Dynamics

B.1.1 Detailed Solvent Equilibration Protocol

The initial model structure coordinates¹ were modified to a topology in which the nucleophile (U-1:O2') was deprotonated and C41 and C75 were protonated at N3. This structure was then neutralized by the addition of Na⁺ ions by placement in the areas of most negative electrostatic potential. The resulting arrangement was solvated in a truncated octahedral cell and additional water molecules were randomly converted to either Na⁺ or Cl⁻ ions until the NaCl concentration was 140 mM when considering only the water content at a density of 0.995 g/mL (the bulk density of the TIP4P-Ew model

¹ Barbara Golden, *personal communication*.

at 298 K[138]). After steepest descent minimization, the system was heated from 0 to 300 K in 300 ps of dynamics plus an additional 200 ps of dynamics. The system density was then relaxed for 500 ps with added pressure coupling. These trajectories utilized a 2 fs time step (all others used 1 fs) and harmonic Cartesian restraints on the RNA heavy atoms and Mg^{2+} ions [$U_{\text{restraint}}(\mathbf{x}) = k(\mathbf{x} - \mathbf{x}_{\text{initial}})^2$, $k = 2\text{-}5$ kcal/mol-Å²]. Unless stated otherwise, an Andersen thermostat[49] was used at 300 K with “massive” collisions every 2000 integration steps. Pressure coupling was performed with a Berendsen barostat[167] applied at 1 bar with a 5 ps coupling constant.

The density relaxed structure was modified as necessary by the addition or removal of protons. If a modification left the system with a net charge, then a bulk ion was randomly removed. After modification, the system was again minimized by steepest descent and heated from 0 to 300 K in 300 ps plus an additional 500 ps of dynamics. A sequence of “annealing” was then performed. An annealing step consisted of heating from 300 - 600 K in 300 ps, 500 ps of dynamics at 600 K, cooling from 600 - 300 K in 300 ps, and 2.5 ns of dynamics at 300 K. This was repeated twice, followed by 2 ns of dynamics and geometric removal of harmonic Cartesian restraints (same as above, except no restraints on Mg^{2+} ions) at 100 ps intervals (500 ps total). Finally, steepest descent minimization and the heating step were performed again before production dynamics.

In order to avoid excessive repetition of the above procedure, some trajectories were instead initialized from starting coordinates taken after 50 - 150 ns of production. After coordinate modification, the annealing steps were omitted, but all other steps remained the same. This protocol was reserved for instances in which the ionic environment was not expected to change greatly, namely local changes in Mg^{2+} binding either via steered MD or by swapping with a randomly chosen bulk Na^+ ion.

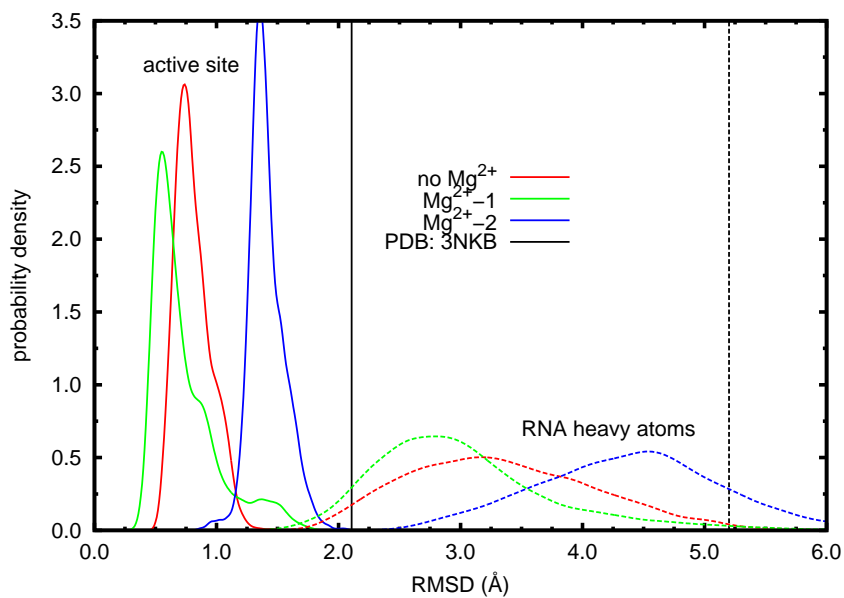
B.1.2 Cluster Analysis

Table B.1: Cluster analysis of ≥ 100 ns MD trajectories produced very few clusters with greater than 1% occupancy. A single dominant cluster ($>80\%$ occupancy) is obtained in all cases. Mg^{2+} -1 indicates the ion position suggested by Chen, *et al.*[231] while Mg^{2+} -2 indicates an alternate, but similar, binding mode that does not include direct coordination to U-1:O2' (but still U23:*pro*-S_P and U-1:*pro*-R_P).

	cluster	% occupation	avg. RMSD wrt centroid (Å)
no Mg^{2+} (100 ns)	1	91.4	1.16
	2	17.0	1.06
	3	16.0	1.19
	4	<0.1	1.23
Mg^{2+} -1 (195 ns)	1	91.6	1.30
	2	7.0	1.46
	3	1.1	1.42
	4	0.3	1.29
Mg^{2+} -2 (100 ns)	1	89.0	1.46
	2	8.8	1.37
	3	2.2	1.48
	4	<0.1	1.66

Structural analysis was performed via a hierarchical agglomerative clustering algorithm (as implemented in CPPTRAJ[269]) using the mass-weighted root mean square deviation (RMSD) of heavy atoms in the active site, defined as residues U-1, G1, U20, C21, C22, U23, C24, G25, and C75. Because the U-1 nucleobase is solvent exposed it is free to undergo syn/anti inversion which causes large changes in the RMSD while the rest of the active site remains constant. Accordingly, all nucleobase heavy atoms from U-1 were omitted from RMSD calculations. For each trajectory the cluster count (a required input parameter) was sequentially reduced from six until no more than one cluster had a fractional occupation below 1%. In all cases this procedure gave rise to

Figure B.1: RMSD distributions of long MD trajectories (≥ 100 ns) with respect to the centroid (active site atoms only) of the most populous cluster when an Mg^{2+} ion was bound at the position hypothesized by Chen, *et al.* (Mg^{2+} -1). An alternate binding mode (Mg^{2+} -2) was also explored, as well as a trajectory in which no Mg^{2+} ion was initially bound. Much narrower distributions are obtained when considering only the active site atoms (solid lines) as opposed to all RNA heavy atoms (dashed lines).



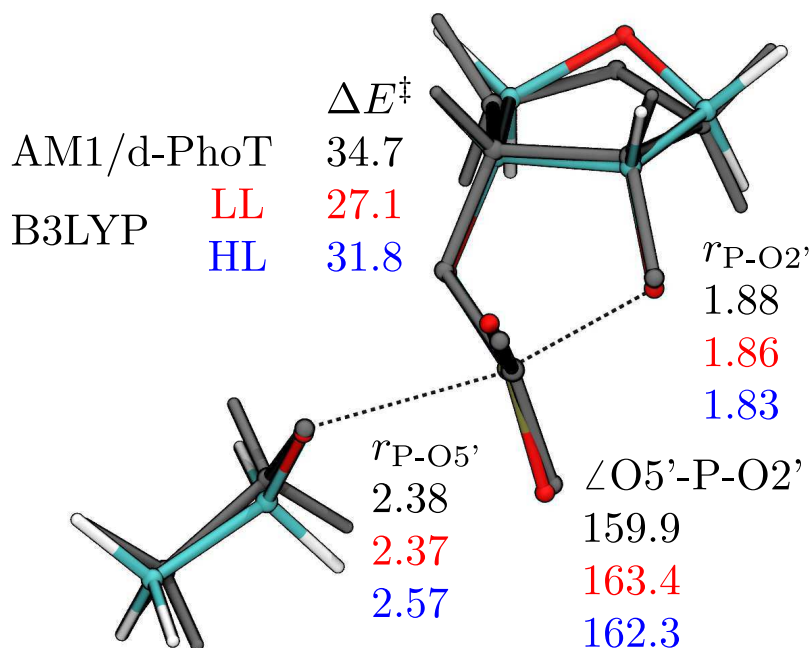
primary clusters with >80% occupancy (B.1).

B.2 RISM and NLPB Calculations

Reference interaction site model (RISM) calculations were performed using the Kovalenko-Hirata (KH)[42] and n th order partial series expansion (PSE n , $n=2,3$)[45] closures. Susceptibility inputs for 3D-RISM were obtained from 1D-RISM calculations using a dielectrically consistent formulation[46] on a grid of 32768 points and 5 MDIIS vectors (all other numerical parameters were set to the default). The temperature and dielectric constant were set to 300 K and 78.4461950541166, respectively. The solvent was modeled with the cSPC/E[47] water model and appropriate ion parameters[169]. A constant density approximation was assumed in which the density of water sites was lowered from the neat density of 55.428 M (≈ 0.0334 molecules/ \AA^3) as the salt concentration increased. 3D-RISM calculations used a 96 \AA buffer distance with a grid spacing of 0.5 \AA and no solute-solvent interaction cutoff (“-solvcut” $\gg 96$). The residual was solved to a tolerance of 10^{-6} using 5 MDIIS vectors and a step size of 0.7.

Because Non-Linear Poisson-Boltzmann (NLPB) calculations are so widespread in the literature, especially on RNA systems (*e.g.* Ref. 270, 235, 237), a slightly different protocol was used in order to make our results comparable to those previously published. Following Misra and Draper[270] dielectric constants of 80.0 and 2.0 were used for the solvent and solute interior respectively. Solute atom radii were taken from the default mbondi set in AmberTools 14 with a 1.4 \AA solvent probe radius. The only exception to this was explicit Mg^{2+} ions, which were given a radius of 1.45 \AA ; this value best reproduces the experimental solvation free energy of -433.3 kcal/mol[270]. Implicit Na^+ and Cl^- ions were modeled with an exclusion radius of 2.0 \AA . Fine grid lengths and spacings were chosen as close as possible to those from 3D-RISM while still satisfying the multigrid constraints imposed by the Adaptive Poisson-Boltzmann Solver[258, 259]. Coarse grid lengths used the same number of grid points but twice the box length (*i.e.* half the fine grid resolution).

Figure B.2: Comparison of quantum chemical models on a model system for RNA backbone cleavage. The displayed structures are overlays of optimized saddle point geometries in vacuum from B3LYP/6-311++G(3df,2p)//6-31++G** (high level/HL, colored) and AM1/d-PhoT (dark gray). Select quantities for both levels of theory, as well as B3LYP/6-31G** (low level/LL) are also shown. Energies are in kcal/mol, distances are in Å, and angles are in degrees.



B.3 Semiempirical Model Validation

B.3.1 Quantum Chemical Calculations

Density functional theory (DFT) and semiempirical quantum calculations were performed on a model system representative of the RNA backbone (B.2). Optimizations were performed with the Gaussian 09 program[183] either directly (for the B3LYP and M06-2X functionals) or with AMBER 14 as an external routine (for AM1/d-PhoT). Optimizations were done in both the gas phase and, when available, a polarizable continuum model with radii chosen as in Ref. 142.

B.3.2 Benchmark Calculations on a Model for RNA Cleavage

Fast, approximate semiempirical quantum methods are reported in the chemical literature for a wide range of purposes. Especially in the last several years, direct comparisons with more accurate DFT based MD simulations have been possible, albeit with sampling on a much shorter time scale (<10 ps trajectories) and more severe approximations for long range electrostatics[178, 255]. This has led to some doubts as to the qualitative and quantitative accuracy of semiempirical models[271]. Nonetheless, our work has previously shown that AM1/d-PhoT accurately reproduces high level DFT energies and geometries in vacuum. Furthermore, these tests indicated that the modest basis sets used for QM/MM MD, especially Pople-type sets neglecting diffuse functions, can lead to significant underestimation of reaction energies for highly charged phosphoryl transfer reactions[1].

As an extension of our previous tests, we use a dinucleotide-like model for specific base catalyzed cleavage of the RNA backbone. Previous studies found that variations on this model reasonably capture the energetics and bonding environment found in the transition state for RNA 2'-*O*-transphosphorylation by ribonuclease A[6]. Here it is only required that such a model provide a reasonable baseline for phosphoryl transfer chemistry, the effects of the enzyme environment being added later via QM/MM. As seen in B.2, AM1/d-PhoT compares favorably with the high level B3LYP reference in both geometry and energies, although with slight overestimation of the latter. Importantly, by about as much as AM1/d-PhoT *overestimates* the reference energy (~ 3 kcal/mol), neglecting diffuse functions on heavy atoms *underestimates* it. These errors are arguably insignificant, as they are on the order of differences between results from the B3LYP and M06-2X functionals, although these are reduced by an order of magnitude when using implicit solvent (see supporting information). The most important consideration here is that QM/MM simulations with DFT (usually with the B3LYP functional) often *do* neglect diffuse functions for the sake of efficiency[255], but under the assumption that this will not greatly affect the energy. However, the example shown here indicates that the loss in accuracy by doing so may in fact be comparable to abandoning DFT entirely and still has the distinct disadvantage of being several orders of magnitude less efficient than semiempirical methods.

Table B.2: Quantum chemical calculations at various levels of theory for phosphoryl transfer in a model system in both the gas phase and solvent (either with an implicit polarizable continuum model (PCM) or explicit MM model). Reactant (R) and transition state (\ddagger) geometries are designated by the value of an atom transfer coordinate ($\xi = r_{\text{P-O5}'} - r_{\text{P-O2}'}$, in Å) in the optimized structure. Barrier heights (ΔE^\ddagger) and reaction energies (ΔE) are given in kcal/mol using the classical energies without zero-point correction. ^a Actually 6-311++G(3df,2p)//6-31++G**. ^b From QM/MM umbrella sampling reported in Ref. 1. The system in that study also differed slightly by inclusion of a full ribose ring.

		ξ_{R}	ξ^\ddagger	ΔE^\ddagger	ΔE	
gas phase	B3LYP	6-31G*	-2.87	0.73	27.0	-48.3
		6-31G**	-2.87	0.74	27.1	-48.1
		6-31+G*	-2.88	0.55	30.4	-49.3
		6-31++G**	-2.87	0.51	30.0	-49.9
		6-311++G(3df,2p) ^a	-	-	31.8	-50.3
		6-31G*	-2.87	0.85	30.2	-41.6
		6-31G**	-2.86	0.85	30.3	-41.4
		6-31+G*	-2.91	0.00	0.0	-42.2
		6-31++G**	-2.91	0.70	33.5	-42.3
		6-311++G(3df,2p) ^a	-	-	34.7	-43.2
	AM1/d-PhoT	-2.62	0.50	34.7	-29.6	
	AM1/d-PhoT+SMAP	-2.87	0.49	30.2	-31.7	
PCM	B3LYP	6-31G*	-1.96	0.74	15.4	10.0
		6-31G**	-1.97	0.74	15.4	10.1
		6-31+G*	-2.42	0.61	19.8	5.9
		6-31++G**	-2.41	0.60	19.8	5.9
		6-311++G(3df,2p) ^a	-	-	22.1	5.8
		6-31G*	-2.68	0.81	16.0	15.0
	M06-2X					

6-31G**	-2.69	0.80	16.1	15.1
6-31+G*	-2.59	0.71	19.7	12.2
6-31++G**	-2.59	0.71	19.7	12.2
6-311++G(3df,2p) ^a	-	-	21.9	12.2
AM1/d-PhoT/MM ^b	-2.64	0.48	26.2	-

B.4 Free Energy Surface Stationary Point Analysis

A key advantage of the variational free energy profile method[5] is that it readily constructs smooth, differentiable free energy surfaces. It is thus straightforward to identify and assess stationary points via normal mode analysis provided that a mass can be meaningfully assigned to the relevant coordinates. Unfortunately, this can sometimes only be done approximately. In the present case the progress coordinates are taken to be atom transfer coordinates approximated as the asymmetric stretching modes of the appropriate collinear triatomic. For simplicity, the reduced masses were calculated using force constants from the AMBER force field (B.3). Large changes in the parameters (up to 100%) changed the masses by less than 1%. The mass matrix so obtained was then used to mass weight the Hessian at each stationary point and yielded the desired eigenvalues after appropriate unit conversions.

Table B.3: Parameters used for reduced mass, μ , determination of atom transfer coordinates. Since the reduced mass of each mode is dependent on the potential (assuming the system is at a stationary point), harmonic spring constants, k , were taken from the AMBER force field.

mode	μ (g/mol)	bond	k (kcal/mol-Å ²)
O-P-O	21.2050	P-O	230.0
O-H-N	1.0405	O-H	553.0
		N-H	434.0