

Towards more robust and efficient methods for the calculation of Protein-Ligand binding affinities

Lingle Wang

Submitted in partial fulfillment of the
requirements for the degree
of Doctor of Philosophy
in the Graduate School of Arts and Sciences

COLUMBIA UNIVERSITY

2012

©2012
Lingle Wang
All Rights Reserved

ABSTRACT

Towards more robust and efficient methods for the calculation of Protein-Ligand binding affinities

Lingle Wang

Biological processes often depend on protein-ligand binding events, so that accurate prediction of protein-ligand binding affinities is of central importance in structural based drug design. Although many techniques exist for calculating protein-ligand binding affinities, ranging from techniques that should be accurate in principle, such as free energy perturbation (FEP) theory, to relatively simple approximations based on empirically derived scoring functions, the counterbalancing demands of speed and accuracy have left us with no completely satisfactory solution thus far. This thesis will be focused on the methodology development towards more robust and reliable Protein-Ligand binding affinity calculation.

In Part I, we will present the WaterMap method, which will bridge the gap between the efficiency of empirical scoring functions and the accuracy of rigorous FEP methods. Unlike most other methods with the main focus on the direct interaction between the protein and the ligand, the WaterMap method we developed considers the explicit driving force from the solvent, in which several individual water molecules in the binding pocket play an active role in the binding process. We demonstrate that protein may adopt active site geometries that will destabilize the water molecules in the binding pocket through hydrophobic enclosure and/or correlated hydrogen bonds, and displacement of these water molecules by ligand groups complementary to protein surface will provide the driving force for ligand binding. In some extreme cases, the interactions are so unfavorable for water molecules that a void is formed in the binding pocket of protein. Our method also considers the contribution from occupation of ligand atoms in the dry regions of binding pocket, which in some cases provides the driving force for ligand binding.

FEP provides an in-principle rigorous method to calculate protein-ligand binding affini-

ties within the limitations of the potential energy model and it may have a potentially large impact on structure based drug design projects especially during late stage lead optimization when productive decisions about compound modification are made . However, converging explicit solvent simulations to the desired precision is far from trivial, especially when there are large structural reorganizations in the protein or in the ligand upon the formation of the binding complex or upon the alchemical transformation from one ligand to another. In these cases, there can be large energy barriers separating the different conformations and the ligand or the protein may remain kinetically trapped in the starting configuration for a very long time during brute-force FEP/MD simulations. The incomplete sampling of the configuration space results in the computed binding free energies being dependent on the starting protein or ligand configurations, thus giving rise to the well known quasi-nonergodicity problem in FEP.

In Part II, we will present a new protocol called FEP/REST, which combines the recently developed enhanced sampling technique REST (Replica Exchange with Solute Tempering) into normal FEP to solve the sampling problem in brute force FEP calculation. The computational cost of this method is comparable with normal FEP, and it can be very easily generalized to more complicated systems of pharmaceutical interest. We apply this method to two modifications of protein-ligand complexes which lead to significant conformational changes, the first in the protein and the second in the ligand. The new approach is shown to facilitate sampling in these challenging cases where high free energy barriers separate the initial and final conformations, and leads to superior convergence of the free energy as demonstrated both by consistency of the results (independence from the starting conformation) and agreement with experimental binding affinity data.

Part III focus on two topics towards the foundational understanding of hydrophobic interactions and electrostatic interactions. To be specific, the nonadditivity effect of hydrophobic interactions in model enclosures is studied in Chapter 9, and the competition between hydrophobic interaction and electrostatic interaction between a hydrophobe and model enclosure is studied in Chapter 10. The approximations in popular implicit solvent models, like the surface area model in hydrophobic interaction, and the quadratic dependence of electrostatic interaction on the magnitude of charge are investigated.

Six of the Chapters (Chapter 2-4, Chapter 6, and Chapter 9-10) have been published, and the other one (Chapter 7) has been accepted for publication and currently is in press. Each Part begins with its own introduction. Each chapter also contains its own abstract and introduction, and focus on one specific topic. They all share the common theme, that is to develop more robust and reliable methods to calculate protein-ligand binding affinities. The conclusions and discussions about future research directions are presented in Part IV.

Table of Contents

I	Development of WaterMap method	1
1	Introduction of WaterMap method	2
2	Thermodynamic properties of liquid water: an application of a nonparametric approach to computing the entropy of a neat fluid	6
2.1	Introduction	7
2.2	Methods	9
2.2.1	The Entropy expression of a neat fluid	9
2.2.2	Factorization of the orientational pair correlation function using generalized Kirkwood superposition approximation	11
2.2.3	The k'th nearest-neighbor method	13
2.2.4	Error analysis of the k'th nearest neighbor method	14
2.2.5	Calculation of the excess energy, enthalpy, and free energy	16
2.2.6	The finite-difference method of entropy calculation	16
2.2.7	Details of the simulation	17
2.3	Results and discussion	17
2.3.1	The Shannon entropies	17
2.3.2	Convergence properties	18
2.3.3	Error analysis	19
2.3.4	The radial dependence of orientational Shannon entropy	19
2.3.5	Inclusion of $g(\theta_1, \chi_1)$ in the factorization	19
2.3.6	Comparison of free energy results	20

2.3.7	Entropy calculation from FD method	21
2.4	Conclusion	21
3	A displaced-solvent functional analysis of model hydrophobic enclosures	42
3.1	Introduction	43
3.2	Methods	44
3.2.1	Derivation of the displaced solvent functional approach to computing protein ligand binding free energies	44
3.2.2	Simulation details	54
3.3	Results	58
3.4	Conclusion	60
4	Protein-Ligand binding: Contributions from wet and dry regions of the binding pocket	69
4.1	Introduction	70
4.2	Results and Discussions	71
4.3	Conclusion	75
4.4	Systems and Simulations	76
4.5	Methods	76
4.5.1	WaterMap calculation	77
4.5.2	Cavity calculation	77
4.5.3	Protein-ligand binding affinity analysis	78
II	Development of FEP/REST	89
5	Introduction of the FEP/REST method	90
6	Replica Exchange with Solute Scaling: A more efficient version of Replica Exchange with Solute Tempering	93
6.1	Introduction	94
6.2	Methodology	95
6.3	Results and Discussion	98

6.4	Conclusion	102
7	On achieving high accuracy and reliability in the calculation of relative protein-ligand binding affinities	110
7.1	Introduction	111
7.2	Results	113
7.3	Discussions and Conclusions	119
7.4	Methods	121
7.4.1	FEP/REST	121
7.4.2	Details of the simulations	122
III	Investigations about hydrophobic interaction and electrostatic interaction	134
8	Introduction of hydrophobic interactions and electrostatic interactions	135
9	Hydrophobic interactions in model enclosures from small to large length scales: nonadditivity in explicit and implicit solvent models.	137
9.1	Introduction	138
9.2	Definition of cooperative and anti-cooperative effects	140
9.3	Simulation details	141
9.3.1	Implicit solvent model calculations	143
9.4	Results and discussion	144
9.4.1	Comparison for different implicit solvent models in predicting the binding affinity	144
9.4.2	Comparison of MSA and SASA model predictions of the non-additivity effect	145
9.4.3	Non-additivity effect at wetting-dewetting transition	149
9.4.4	Higher order multi-body interactions	151
9.5	Conclusion	151

10 Competition between electrostatic and hydrophobic interactions	165
10.1 Introduction	166
10.2 Details of simulation	168
10.3 Results and discussion	170
10.3.1 Binding affinity results	170
10.3.2 Dependence of the binding affinity (Solvation free energy) on the mag- nitude of charge	171
10.4 Theoretical derivation for electrostatic contribution to the solvation free energy	173
10.5 Further evidence to validate the theory	175
10.6 conclusions	177
IV Conclusions	185
11 Conclusions and future research directions	186
V Bibliography	192
Bibliography	193
VI Appendices	214
A Error analysis in NN method and Constant pressure correction	215
A.1 $Var[\ln f(x)]$ for some special cases	215
A.1.1 Gaussian distribution	215
A.1.2 exponential distribution	216
A.1.3 exponential distribution in a finite range	216
A.1.4 linear distribution in a finite range	216
A.2 Determination of most proper weights	217
A.3 Constant pressure correction to ΔG_{sim} for the FD entropy	218

B Comparison between REST and TREM	221
B.1 Can TREM be more efficient than REST1?	221
B.2 Why REST2 is more efficient than TREM.	222
C Comparison between OPLS 2005 and OPLS 2.0 force fields for ligands	
CDA and CDB	223
C.1 Charge distributions on the Pyridine ring	223
C.2 The distribution of dihedral angle involved in the flipping of the P1 pyridine ring	225
D How to Treat bonded interactions in FEP simulation	227
E Structures of ligands studied in Chapter 4	231

List of Figures

2.1	NN estimate of first shell $S^{[t_1, t_2]}$ for TIP3P	26
2.2	NN estimate of second shell $S^{[t_1, t_2]}$ for TIP3P	27
2.3	NN estimate of third shell $S^{[t_1, t_2]}$ for TIP3P	28
2.4	Histogram method estimate of first shell $S^{[t_1, t_2]}$ for TIP3P	29
2.5	Histogram method estimate of second shell $S^{[t_1, t_2]}$ for TIP3P	30
2.6	Histogram method estimate of third shell $S^{[t_1, t_2]}$ for TIP3P	31
2.7	Total orientational excess entropy for TIP3P	32
2.8	Total orientational excess entropy for SPC	33
2.9	Total orientational excess entropy for SPC/E	34
2.10	Total orientational excess entropy for TIP4P	35
2.11	Total orientational excess entropy for TIP4P-Ew	36
2.12	Error analysis of first shell $S^{[t_1, t_2]}$ for TIP3P	37
2.13	Error analysis of second shell $S^{[t_1, t_2]}$ for TIP3P	38
2.14	r dependence of orientational Shannon entropy	39
2.15	Distribution function $g(\theta_1, \chi_1)$ by assuming independence of two angles	40
2.16	Distribution function $g(\theta_1, \chi_1)$	41
3.1	Protein ligand direct interaction and ligand charging free energy.	62
3.2	Surface and volume terms in IST	63
3.3	The effective volume of methane in different enclosures	64
3.4	Geometries of model enclosures	65
3.5	The coordinate system to define the orientation of water inside enclosure	66
3.6	Performance of DSF (WaterMap)	67

3.7	Performance of surface area based models	68
4.1	The hydration sites and dry region in the binding pocket of MUP	80
4.2	Predicted binding affinities for ligands binding to MUP	81
4.3	Comparison between WaterMap and MMGBSA	82
4.4	Relative binding affinity predictions using different methods	83
4.5	Comparison of Ligand HE4 and Ligand OC9	84
4.6	Combination of WaterMap with cavity contribution to rank-order congeneric ligands	85
4.7	Ligands LTL and TZL binding to the MUP receptor in the dry region. . . .	86
4.8	Molecular recognition motif between dry regions in the binding pocket and hydrophobic groups in the ligand	87
6.1	Temperature profile of four replicas using REST2 for trpcage	105
6.2	Protein heavy atom RMSD from native structure using REST2	106
6.3	Temperature profile and Protein heavy atom RMSD from native structure for the β -hairpin system using REST2	106
6.4	Protein heavy atom RMSD from native structure from a single REST2 replica at high temperature	107
6.5	Anti-correlation between the intra-molecular potential energy of the protein and the interaction energy between the protein and water	107
6.6	Distributions of intra-molecular potential energy of the protein and the in- teraction energy between protein and water	108
6.7	Comparison between TREM and REST2	109
7.1	The structures of T4L/L99A and the Thrombin systems	129
7.2	The Val111 side chain conformation sampled using different methods	130
7.3	The correct binding pose and the erroneous conformation sampled in simu- lation using OPLS 2005 force field	131
7.4	The conformations of the ligands in the Thrombin system sampled using different methods	132
7.5	Distributions of the ligand conformations with protein heavy atoms restrained	133

7.6	FEP/REST protocol	133
9.1	The geometries of model hydrophobic enclosures.	156
9.2	The geometries of model hydrophobic enclosures	157
9.3	Buried surface area/molecular mechanics prediction of methane-enclosure binding affinities	158
9.4	Buried surface area prediction of nonadditivity effects in model enclosures .	159
9.5	The buried MSA/SASA predictions of nonadditivity effect in methane trimer	160
9.6	Difference between MSA and SASA models for the prediction of the nonad- ditivity effects	161
9.7	The nonadditivity effect in the wetting and dewetting region of hydrophobic enclosures	162
9.8	The performance of GKSA as a function of the multi-body interactions included	163
9.9	The performance of GKSA as a function of the multi-body interactions included	164
10.1	Thermodynamic cycles connecting methane-plates binding affinities and charg- ing free energies	179
10.2	Dependence of the charging free energy on the magnitude of the charge . .	180
10.3	Transition from “hydrophobic like“ to “hydrophilic like“ as the magnitude of charge on the plates is increased	181
10.4	Configuration of water between charged plates	182
10.5	Dependence of charging free energy on the magnitude of charge for asym- metric plates	183
10.6	Asymmetry of the solvation free energies between cations and anions	184
A.1	Constant pressure corrections to free energy	220
C.1	The atom numbering on the P3 pyridine ring for ligand CDA	224
C.2	Distribution of Ligand conformations in complex using OPLS 2005	225
C.3	Distribution of ligand conformations in gas phase using OPLS 2005	226
D.1	The structure of the mixed molecule to define the mutation path	229
D.2	Instabilities of the bonded interactions with default setup of Desmond . . .	229

List of Tables

2.1	Orientalional Shannon entropies of the five water models	23
2.2	Comparison of entropy results from the NN method and cell theory	23
2.3	Results for the energy, enthalpy, and entropy of liquid water from various methods	24
2.4	Entropy results from FD method and comparison with other methods	25
3.1	Thermodynamics of methane enclosure binding	61
4.1	Decomposition of binding affinities into WaterMap and Cavity contribution	88
7.1	Predicted relative binding affinities for T4L/L99A system using different methods	125
7.2	Predicted relative binding affinities for Thrombin system using OPLS 2005 Force Field	125
7.3	Predicted relative binding affinities for Thrombin system using OPLS 2.0 Force Field	126
7.4	Lambda values, scaling factors and free energy difference between neighboring lambda windows for T4L/L99A system	127
7.5	Lambda values, scaling factors and free energy difference between neighboring lambda windows for Thrombin system	128
9.1	The binding thermodynamics of methane for the various model hydrophobic enclosures.	154
9.2	Multi-body potential of mean force in model hydrophobic enclosures	155

C.1 Charge distributions on atoms of the P3 pyridine ring for ligand CDA for two force fields.	224
---	-----

Acknowledgments

I would like to thank my advisors, Prof. Bruce J. Berne and Prof. Rich A. Friesner for their patience, support, and motivation. They are wonderful advisors, very insightful to choose research projects that are important but still tractable. Their advises are invaluable to my research. Their enthusiasm and inspiration about research is the most important thing I learned from them.

I would like to thank my thesis committee, Prof. Dave R. Reichman, Prof. Ruben Gonzalez, and Dr. Ruhong Zhou, for their helpful suggestions and comments about how to improve the dissertation. They also gave me a lot of advises during my PHD study on how to proceed the research projects.

I would like to especially thank Dr. Robert Abel. We cooperated on the WaterMap project. He brought me into the field of Protein-Ligand binding, and taught me a lot about how the basic principles in statistical mechanics are applied in real projects. I learned a lot from him about how to run the desmond and WaterMap program. Without his help, these projects couldn't be finished so quickly.

I would like to thank Prof. Tom Young. It was him, Dr. Robert Abel, Prof. B. J. Berne and Prof. R. A. Friesner who made the original WaterMap method working. We discussed a lot on the WaterMap project, and his suggestions and comments are very helpful to my research.

I would like to also thank Dr. Teng Lin, Dr. Byungchan Kim, Dr. Yujie Wu, Dr. Yuqing Deng from Schrodinger. They helped me to implement our FEP/REST algorithm into desmond. Dr. Ed Harder from Schrodinger is also acknowledged for making available to us the OPLS 2.0 force field. Without their help, these projects couldn't be finished so quickly.

All the members in the Berne and Friesner group are also acknowledged, including the

coordinator Betty Cusack and computing center technician Calman Lobel. We had a lot of discussions, and their suggestions are very helpful to my research. They have made my time at Columbia University more enjoyable and less stressful than graduate school typically is.

I would also like to thank all my friends, both at Columbia University, and all over the world. They have made my time in the US more enjoyable and helped me to get through some hard and depressed moments in life. I wish they all enjoy every day and achieve their goals soon.

Especially, I would like to thank my girlfriend, Hongyun Wang. I am indebted to her for her understanding, encouragement, quiet patience and unwavering support. I enjoy the time we spent together, and she made my life more interesting and much happier.

Most of all, I would like to thank my family, my parents and my sister. They helped me throughout this process with their support, their love and their enthusiasm. Born in a small town in China, I never imagined I would study at Columbia University for doctoral degree even in the wildest dream. It was my parents, with their forever love, with their never-ending source of support, and with their sacrifice of their own life, who made this dream come true. So I would like to express my deepest gratitude to my parents for their unflagging love and support throughout my life.

To my Mother and Father

Part I

Development of WaterMap method

Chapter 1

Introduction of WaterMap method

Water is unique among liquids for its biological significance and it plays an important role in the protein ligand binding process. It is widely believed that displacement of water molecules in the binding pocket of protein is a principal source of binding free energy. Water molecules in the binding pocket of protein are often entropically unfavorable due to the orientational and positional restraints imposed by the protein residues, or they are energetically unfavorable due to the breaking of hydrogen bonds when surrounded by hydrophobic groups of the protein. In both cases, they are free energetically unfavorable compared to those in bulk solution, and displacement of these free energetically unfavorable water molecules by ligand groups complementary to the protein surface will provide the driving force for ligand binding.

It has been demonstrated that there are two special regions in the binding pocket of the protein where the water molecules are extremely unfavorable, and empirical scoring functions will underestimate the ligand binding affinity when these water molecules are displaced. The first motif is called hydrophobic enclosure, where the water molecules are surrounded on multiple sides by hydrophobic protein side chains. In this case, water molecules will lose hydrogen bonds compared to bulk solution, and they are energetically very unfavorable. The second motif is called correlated hydrogen bonds, where water molecules in the binding pocket need to make several hydrogen bonds with the protein residues, and they are entropically very unfavorable compared to bulk solution.

The WaterMap method is designed to characterize the contribution of water displace-

ment to the protein ligand binding affinity. The WaterMap method for computing the free energy contribution of the ligand displacing the active site solvent molecules begins with the assumption the equilibrium properties of the hydration of the apo-receptor active site can be discerned from a converged explicitly solvated classical MD simulation of the protein. The positions of all water molecules that enter the protein active site during this dynamics simulation are recorded and clustered into high occupancy 1 Å radius spheres, which we denoted as the “hydration sites” of the active site cavity. Using methods borrowed from inhomogeneous solvation theory, the average system interaction energies, and excess entropies for the water in each hydration site are computed. The average system interaction energies of the water in the various hydration sites can be readily extracted from the dynamics simulation, and the excess entropies are calculated from a truncated expansion of the entropy in terms of solvent orientational and spatial correlation functions. Comparing the system interaction energy of water in a hydration site with that of water in the bulk fluid lets us estimate the enthalpic cost transferring the solvent in the hydration site from the protein active site to the bulk. The excess entropies calculated with this method may be used similarly. These calculations allow us to create a hydration thermodynamics map for a given receptor.

The thermodynamics of water in the binding pocket of protein can be used to estimate the free energy contribution of a ligand displacing water from the protein active site by noting that (1) if the ligand sterically overlaps with a given hydration site in its bound conformation, then it displaces water from that hydration site; and, (2) the higher the excess chemical potential of the solvent in a given hydration site, the more favorable its evacuation to the bulk fluid will be. With these assumptions in mind, a simple “displaced solvent functional” was formulated that attempts to correctly evaluate this contribution to the binding affinity by computing the transfer free energy of the solvent evacuated from the hydration sites by the ligand. The functional itself is

$$\Delta G_{bind} = \sum_{lig,hs} \Delta G_{hs} \left(1 - \frac{|\mathbf{r}_{lig} - \mathbf{r}_{hs}|}{R_{CO}} \right) \Theta (R_{CO} - |\mathbf{r}_{lig} - \mathbf{r}_{hs}|) \quad (1.1)$$

where ΔG_{bind} is the predicted binding free energy of the ligand evacuating the solvent from the active site, R_{co} is the distance cutoff for a ligand atom beginning to displace

a hydration site, ΔG_{hs} is the computed free energy of transferring the solvent in a given hydration site from the active site to the bulk fluid, and $\Theta(x)$ is the Heaviside step function. The contribution from each hydration site was capped, such that it would never contribute more than ΔG_{hs} to ΔG_{bind} no matter how many ligand atoms were in close proximity to it.

In Chapter 2 of this section, the inhomogeneous solvation theory (IST) is introduced, and an efficient method to calculate the entropy of water in different environments is introduced, which facilitates faster convergence and greater accuracy. In Chapter 3 of this section, the WaterMap method is applied to a number of model hydrophobic enclosures, and the WaterMap predicted binding affinities are compared with high accuracy free energy perturbation (FEP) results. The high correlation between the results using WaterMap and FEP indicates that WaterMap is a very useful model to characterize the contribution from water displacement to the binding affinity and a large amount of the binding affinity comes from the water displacement. In addition, the physical-chemical basis and the key approximation of WaterMap method are clarified, which facilitates an understanding of when the technique is expected to succeed and fail.

The hydrogen sites identified by WaterMap method usually correspond to the structured water molecules in the X-ray crystallography, and WaterMap only takes into consideration the contribution from the high solvent occupation regions in the binding pocket. In some extreme cases, a portion of the receptor active site is so unfavorable for water molecules that a void or a dry region is formed in the binding pocket. In Chapter 4 of this section, we demonstrate that the presence of dry regions has a nontrivial effect on ligand binding affinity if the ligand places atoms in these regions, and an additional term attributable to occupation of the dry regions by ligand atoms is introduced. The combination of these two terms has been shown to be more predictive in relative ligand binding affinity calculation for a set of congeneric ligands binding to MUP receptor which has a dry region in the binding pocket. This represents an important addition to the WaterMap method, and the combination of WaterMap with the cavity term will characterize the contribution from both wet and dry region in the binding pocket to ligand binding affinity. In addition, the molecular recognition between the dry region in the binding pocket and the hydrophobic

groups of the ligand occur on many different kinds of proteins, and we suggest that it may represent a general motif for molecular recognition.

Chapter 2

Thermodynamic properties of liquid water: an application of a nonparametric approach to computing the entropy of a neat fluid

Abstract

Due to its fundamental importance to molecular biology, great interest has continued to persist in developing novel techniques to efficiently characterize the thermodynamic and structural features of liquid water. A particularly fruitful approach, first applied to liquid water by Lazaridis and Karplus, is to use molecular dynamics or Monte Carlo simulations to collect the required statistics to integrate the inhomogeneous solvation theory equations for the solvation enthalpy and entropy. We here suggest several technical improvements to this approach, which may facilitate faster convergence and greater accuracy. In particular, we devise a nonparametric k 'th nearest neighbors (NN) based approach to estimate the water-water correlation entropy, and suggest an alternative factorization of the water-water

correlation function that appears to more robustly describe the correlation entropy of the neat fluid. It appears that the NN method offers several advantages over the more common histogram based approaches, including much faster convergence for a given amount of simulation data; an intuitive error bound that may be readily formulated without resorting to block averaging or bootstrapping; and the absence of empirically tuned parameters, which may bias the results in an uncontrolled fashion.

2.1 Introduction

Water is unique among liquids for its biological significance. It plays an active role in the formation of the structures of proteins, lipid bilayers, and nucleic acids in vivo, both through direct hydrogen bonding interactions with these biomolecules, and also through indirect interactions, where the unique hydrogen-bonded structure of liquid water is known to drive hydrophobic assembly [1]. It has been suggested that a robust characterization of the thermodynamic properties and structure of water solvating the active site of a protein is essential to rationalize the various binding affinities of small molecules that will displace that solvent to bind to the protein active site[2; 3].

As such, great interest has continued to persist in developing novel techniques to efficiently characterize the thermodynamic and structural features of liquid water in different environments. A particularly fruitful approach, first applied to liquid water by Lazaridis and Karplus[4; 5; 6], used molecular dynamics or Monte Carlo simulations to collect the required statistics to integrate the inhomogeneous solvation theory (IST) equations for the solvation enthalpy and entropy. In this theory, the solvation enthalpy is determined from an analysis of the change in the solute-solvent and solvent-solvent interaction energy terms, and the solvation entropy is computed from an expansion of the entropy in terms of increasingly higher order solute-solvent correlation functions[4]. This approach has been used to characterize the thermodynamics and structure of neat water[6], hydration of small hydrophobes[4], and the hydration of the active sites of proteins[7; 8]. Recently, it has also been extended to allow for the rapid computation of the relative binding affinities of a set of congeneric ligands with a given protein, via a semi-empirical

displaced-solvent functional[2].

Due to the increasing interest in applying this technique to water[9; 10; 11; 12] in various environments, we have chosen to reexamine the factorization and correlation function integration scheme originally suggested by Lazaridis and Karplus[6] for bulk water and later adopted by others[13]. We have found that several technical improvements in this scheme are possible, which may facilitate faster convergence and greater accuracy than the more typical expressions. In this paper, we (1) devise a nonparametric k 'th nearest neighbors (NN) [14] based approach to estimate the water-water correlation entropy, in lieu of the more common histogram based approaches; and (2) suggest an alternative factorization for the water-water correlation function that appears to more robustly describe the water-water correlation entropy of the neat fluid. To our knowledge, this is the first application of the NN method to compute the entropy of a neat fluid. It appears that the NN method offers several advantages over the more common histogram based approaches, including (1) much faster convergence for a given amount of simulation data, especially when the correlation function is highly structured; (2) an intuitive error bound may be readily formulated without resorting to block averaging or bootstrapping techniques, which may be problematic to apply to estimators of the entropy; and (3) the absence of empirically tuned parameters, such as the histogram bin width, which may bias the results in an unpredictable fashion. Our alternative factorization of the water-water correlation function explicitly includes correlations between the water-dipole-vector-intermolecular-axis angle with the angle of rotation of the water molecule about its dipole vector. This contribution, although neglected by others[6], has been found in our work to increase the agreement of results obtained by the entropy expansion with those obtained by less approximate methods, such as free energy perturbation theory. We also extensively compare the solvation entropies obtained from the truncated entropy expansion to those obtained from a finite difference analysis of free energy perturbation theory results. This comparison allows us to characterize the errors in both precision and accuracy associated with the NN method of integrating the entropy expansion presented here.

Our primary interest in developing this technique was to later adapt the method to study the solvation of solutes; thus, we were interested in determining realistic estimates of

the convergence of the technique when the isotropic symmetry of the fluid was not present. As such, when extracting the solvent configurations to compute the pair correlation function (PCF), we chose to use only the configurations of a distinguished solvent molecule with the rest of the system, instead of collecting statistics from all pairs of solvent molecules. Such a protocol allows for an interrogation of the relative convergence properties of the various methods that might be obscured by the additional statistics offered by taking advantage of the symmetry of the system.

2.2 Methods

2.2.1 The Entropy expression of a neat fluid

First derived by Green[15], and later by Raveché[16] and Wallace[17], the entropy of a fluid can be expressed as a sum of integrals over multi-particle correlation functions. For a molecular fluid[5], the expression is

$$\begin{aligned}
 s &= s^{id} + s_e = s^{id} - \frac{1}{2!} k \frac{\rho}{\Omega^2} \int [g^{(2)} \ln(g^{(2)}) - g^{(2)} + 1] d\mathbf{r} d\omega^2 \\
 &\quad - \frac{1}{3!} k \frac{\rho^2}{\Omega^3} \int [g^{(3)} \ln(\delta g^{(3)}) - g^{(3)} + 3g^{(2)}g^{(2)} - 3g^{(2)} + 1] d\mathbf{r}_1 d\mathbf{r}_2 d\omega^3 - \dots \quad (2.1)
 \end{aligned}$$

where, s^{id} is the entropy of an ideal gas with the same density and temperature as the fluid, s_e is the *excess* entropy of the fluid over that of the ideal gas, k is the Boltzmann's constant, and ρ is the number density, ω denotes the orientational variables of one molecule, Ω is the total volume of the orientational space (For nonlinear molecule like water, Ω is $8\pi^2$), $g^{(2)}$ is the pair correlation function, $g^{(3)}$ is the triplet correlation function, and $\delta g^{(3)}$ is the deviation of $g^{(3)}$ from the superposition approximation. In practice, it is very difficult or even impossible to converge the three-particle and higher order correlation terms. However, it has been established that, for most fluids, the largest contribution to the excess entropy comes from the two-particle correlation term[6], and the error induced by neglecting the higher order terms of the expansion may often be safely ignored.

Following the work of Lazaridis and Karplus[6], we evaluate the two-particle excess entropy of liquid water by separating the two-particle term into translational and orientational

components by factorization:

$$g(r, \omega^2) = g(r)g(\omega^2|r) \quad (2.2)$$

$$s_e^{(2)} = s_{trans}^{(2)} + s_{orient}^{(2)} \quad (2.3)$$

$$s_{trans}^{(2)} = -\frac{1}{2}k\rho \int [g(r) \ln g(r) - g(r) + 1] d\mathbf{r} \quad (2.4)$$

$$s_{orient}^{(2)} = \frac{1}{2}k\rho \int g(r) S^{orient}(r) d\mathbf{r} \quad (2.5)$$

$$S^{orient} = -\frac{1}{\Omega^2} \int J(\omega^2) g(\omega^2|r) \ln g(\omega^2|r) d\omega^2 \quad (2.6)$$

where r is the Oxygen-Oxygen distance of two water molecules, ω^2 are the angles that define the relative orientation of the two water molecules, $J(\omega^2)$ is the Jacobian of the angular variables, $g(r, \omega^2)$ is the pair correlation function, and $g(\omega^2|r)$ is the conditional-angular pair correlation function in the typical Bayesian notion. (Note that $g(r, \omega^2)$ is identical to $g^{(2)}$ as it appears in equation 2.1.) We denote the relative orientation of the two water molecules by the five angles[6] $[\theta_1, \theta_2, \phi, \chi_1, \chi_2]$, where θ_1, θ_2 are the angles between the intermolecular axis and the dipole vector of each molecule, ϕ describes the relative dihedral rotation of the dipole vector around the intermolecular axis, and χ_1, χ_2 describe the rotation of each molecule around its dipole vector. In the following discussion, we denote the entropy defined by formula 2.6 the orientational Shannon entropy[18], and denote the entropy defined by formula 2.5 the orientational excess entropy.

In line with prior work [6], we calculated the orientational Shannon entropy as defined by formula 2.6 for three different ranges of r : ($0 < r \leq 2.7$), ($2.7 < r \leq 3.3$), and ($3.3 < r \leq 5.6$), which correspond to the various peaks and troughs in the radial distribution function. In this way, the orientational excess entropy is related to Shannon entropy by :

$$s_{orient} = \frac{1}{2} N_i k S^{orient} \quad i=1,2,3 \quad (2.7)$$

where N_i is the average number of water molecules in the i -th shell.

2.2.2 Factorization of the orientational pair correlation function using generalized Kirkwood superposition approximation

The orientational pair correlation function (PCF) of water is a function of five angles, which is very difficult to converge from currently accessible molecular dynamics simulation time scales. The idea of factorization is to approximate the higher dimensional probability density function by the product of its lower dimensional marginal probability density functions. The generalized Kirkwood superposition approximation (GKSA)[19; 20; 21], allows an m -dimensional distribution to be estimated using corresponding $m-1$ -dimensional distributions:

$$\rho(x_1, x_2, \dots, x_m) = \begin{cases} \frac{\prod_{c_{m-1}^m} \rho_{m-1} \cdots \prod_{c_2^m} \rho_2}{\prod_{c_{m-2}^m} \rho_{m-2} \cdots \prod_{c_1^m} \rho_1} & m \text{ is odd} \\ \frac{\prod_{c_{m-1}^m} \rho_{m-1} \cdots \prod_{c_1^m} \rho_1}{\prod_{c_{m-2}^m} \rho_{m-2} \cdots \prod_{c_2^m} \rho_2} & m \text{ is even} \end{cases} \quad (2.8)$$

where ρ_{m-k} represents a specific probability density function of $m-k$ dimensionality, and c_{m-k}^m indicates all possible combinations of $m-k$ groupings from the set of m total variables. Reiss[20] and Singer[21] have demonstrated that the GKSA is the optimal approximation of an n -particle distribution for $n \geq 3$ from a variational point of view, and it has been applied in numerous settings[22; 23].

From the results of our simulations, and as indicated by Lazaridis and Karplus[6], the distribution has no structure along angle ϕ , *i.e.* $g(\phi)$ is close to 1 over the range of ϕ , and has no correlation with other angles. Thus, we approximated the 5 dimensional PCF by:

$$g(\theta_1, \theta_2, \phi, \chi_1, \chi_2) = g(\theta_1, \theta_2, \chi_1, \chi_2)g(\phi) \quad (2.9)$$

Note that for any properly defined orientational PCF $g(x_1, x_2 \cdots x_N)$,

$$\frac{1}{\Omega[x_1, x_2 \cdots x_n]} \int J(x_1, x_2 \cdots x_n) g(x_1, x_2 \cdots x_n) dx_1 dx_2 \cdots dx_n = 1 \quad (2.10)$$

where

$$\Omega[x_1, x_2 \cdots x_n] = \int J(x_1, x_2 \cdots x_n) dx_1 dx_2 \cdots dx_n \quad (2.11)$$

i.e., $\Omega^{[x_1, x_2 \dots x_n]}$ is the integral of the Jacobian $J(x_1, x_2 \dots x_n)$ over angular variables $x_1, x_2 \dots x_n$. Therefore, $g(x_1, x_2 \dots x_n)$ is proportional to $\rho(x_1, x_2 \dots x_n)$ with proportional coefficient $\Omega^{[x_1, x_2 \dots x_n]}$. Via application of the GKSA (formula 2.8), it follows

$$g(\theta_1, \theta_2, \chi_1, \chi_2) = \frac{g(\theta_1, \theta_2)g(\theta_1, \chi_1)g(\theta_1, \chi_2)g(\theta_2, \chi_1)g(\theta_2, \chi_2)g(\chi_1, \chi_2)}{g^2(\theta_1)g^2(\theta_2)g^2(\chi_1)g^2(\chi_2)} \quad (2.12)$$

Note that this factorization differs from that introduced by Karplus and Lazaridis[6] by the explicit inclusion of $g(\theta_1, \chi_1)$ and $g(\theta_2, \chi_2)$ terms. Taking this approximation of $g(x_1, x_2 \dots x_n)$ into the argument of the logarithm of formula 2.6 we find

$$S^{orient} = -\frac{1}{\Omega^2} \int J(\omega^2)g(\omega^2|r) \ln g(\omega^2|r)d\omega^2 \quad (2.13)$$

$$= -\sum_{C_2^4} \frac{1}{\Omega^{[x_1, x_2]}} \int J(x_1, x_2)g(x_1, x_2) \ln g(x_1, x_2)dx_1dx_2$$

$$+ 2 \sum_{C_1^4} \frac{1}{\Omega^{[x]}} \int J(x)g(x) \ln g(x)dx \quad (2.14)$$

$$= \sum_{C_2^4} S^{[x_1, x_2]} - 2 \sum_{C_1^4} S^{[x]} \quad (2.15)$$

where x_1, x_2 is any combination of two variables from the $[\theta_1, \theta_2, \chi_1, \chi_2]$ set, x is any variable from the $[\theta_1, \theta_2, \chi_1, \chi_2]$ set, $J(x_1, x_2)$ is the Jacobian of the corresponding two variables, and $J(x)$ is the Jacobian corresponding to variable x , $\Omega^{[x_1, x_2]}$ is the total accessible angular volume of variables x_1, x_2 , and $\Omega^{[x]}$ is the total accessible angular volume of variable x , $S^{[x_1, x_2]}$ is the Shannon entropy of angular variables x_1 and x_2 , and $S^{[x]}$ is the Shannon entropy of angular variable x .

We note that an ambiguity seems to exist in the literature as to how to properly apply an approximation of the type suggested in equation 2.12 to equation 2.6. We have adopted here to apply the approximation only to the logarithm of equation 2.6 (as was done in the original derivation of equation 2.1), which allows result 2.15 to be interpreted through the language of information theory [24]. An alternate approach, that has been adopted by others, has been to apply approximation 2.12 to both occurrences of the PCF in equation 2.6, taking care to renormalize the factorization of the PCF introduced in equation 2.12 so that meaningful results will still be obtained. Interestingly, the results of these two approaches do not numerically agree, which may not be obvious from cursory inspection.

We leave this proof as an exercise for the reader, which can be readily shown for instance from a correlated multidimensional Gaussian distribution.

2.2.3 The k'th nearest-neighbor method

The NN method[14] gives an asymptotically unbiased estimate of an integral of the form:

$$I = - \int \rho(x_1, x_2, \dots, x_s) \ln \rho(x_1, x_2, \dots, x_s) dx_1 dx_2 \dots dx_s \quad (2.16)$$

where $\rho(x_1, x_2, \dots, x_s)$ is the probability density function. Given a reasonable estimation of probability density function $f(x^i)$, the value of integral can be approximated as

$$I \approx -\frac{1}{n} \sum_{i=1}^n \ln f(x^i) \quad (2.17)$$

which follows from x^i being sampled from the true distribution $\rho(x^i)$. The NN method of nonparametrically estimating $f(x^i)$ at a point $x^i = (x_1^i, x_2^i \dots, x_s^i)$ is [25]

$$f(x^i) = \frac{k}{n} \frac{1}{V_s(R_{i,k})} \quad (2.18)$$

$$V_s(R_{i,k}) = \frac{\pi^{s/2} R_{i,k}^s}{\Gamma(\frac{1}{2}s + 1)} \quad (2.19)$$

where n is the number of data points in the sample, $V_s(R_{i,k})$ is the volume of an s -dimensional sphere with radius $R_{i,k}$, and $R_{i,k}$ is the Euclidean distance between the point x^i and its k -th nearest neighbor in the sample. This approximation amounts to assuming that the distance between neighboring sampled points in configuration space will be small where the probability density function is large, and vice versa. So this integration may be estimated as

$$I \approx -\frac{1}{n} \sum_{i=1}^n \ln f(x^i) = \frac{1}{n} \sum_{i=1}^n \ln \frac{n\pi^{s/2} R_{i,k}^s}{k\Gamma(\frac{1}{2}s + 1)} \quad (2.20)$$

However, the estimate in equation 2.20 is systematically biased [14] and will deviate from the correct result in the limit of large n by $L_{k-1} - \ln k - \gamma$, where $L_j = \sum_{i=1}^j \frac{1}{i}$ and $\gamma = 0.5772\dots$ is Euler's constant. By subtracting the bias $L_{k-1} - \ln k - \gamma$, the modified unbiased estimate is formulated as

$$I \approx \frac{s}{n} \sum_{i=1}^n \ln R_{i,k} + \ln \frac{n\pi^{s/2}}{\Gamma(\frac{1}{2}s + 1)} - L_{k-1} + \gamma \quad (2.21)$$

Now our goal is to modify our expressions for the Shannon entropies into a form that is amenable to a k'th NN evaluation of the integral. The expression of the two-dimensional orientational Shannon entropy has the form of

$$S^{[x_1, x_2]} = -\frac{1}{\Omega^{[x_1, x_2]}} \int J(x_1, x_2) g(x_1, x_2) \ln g(x_1, x_2) dx_1 dx_2 \quad (2.22)$$

where $J(x_1, x_2)$ is the Jacobian associated with x_1 and x_2 . Here, for χ_1 and χ_2 the Jacobian is 1, but for θ_1 and θ_2 the Jacobian is $\sin \theta_1$ and $\sin \theta_2$. However, by a change of variables from θ to $t = \frac{\pi}{2}(\cos \theta + 1)$, the Jacobian for t becomes 1, and the total angular volume is π for one dimensional distribution and π^2 for two dimensional distributions. Then, $g(x_1, x_2)$ is proportional to $\rho(x_1, x_2)$ in equation 2.16, with proportional coefficient π^2 . Following the NN method, the statistically unbiased estimation of the one and two-dimensional orientational Shannon entropies may now be approximated as

$$H_k^{[x]}(n) = \frac{1}{n} \sum_{i=1}^n \ln R_{i,k} + \ln \frac{n\pi^{1/2}}{\Gamma(\frac{1}{2} + 1)\Omega^{[x]}} - L_{k-1} + \gamma \quad (2.23)$$

$$H_k^{[x_1, x_2]}(n) = \frac{2}{n} \sum_{i=1}^n \ln R_{i,k} + \ln \frac{n\pi^1}{\Gamma(\frac{1}{2} \times 2 + 1)\Omega^{[x_1, x_2]}} - L_{k-1} + \gamma \quad (2.24)$$

where $H_k^{[x]}(n)$ is the k'th NN estimate of the Shannon entropy of random variable x from a sampling of n data points and $H_k^{[x_1, x_2]}(n)$ is the k'th NN estimate of the joint Shannon entropy of random variables x_1, x_2 from a sampling of n data points. Thus, we are now equipped to apply the NN method of estimating the entropy to liquid state problems. We also note that to compute the NN distances, we made use of the ANN code[26], which utilizes the k-d tree algorithm[27] for obtaining the k-th NN distances $R_{i,k}$ between sample points as necessary.

2.2.4 Error analysis of the k'th nearest neighbor method

It has been shown through an analysis of the limiting distribution[14] that the variance of the k-th NN estimate of the entropy $H_k(n)$ is

$$Var[H_k(n)] = \frac{Q_k + Var[\ln f(x)]}{n} \quad (2.25)$$

where $f(x)$ is the probability density function and $Q_k = \sum_{j=k}^{\infty} \frac{1}{j^2}$. Formally, this result follows from using the Poisson approximation of the binomial distribution to characterize

the fluctuations of $H_k(n)$ in the large n limit (please see ref. [14] for details). Since $H_k(n)$ is asymptotically unbiased[14], the asymptotic mean square error of the estimate is of the order given by equation 2.25. Typically, the true value $H(n)$ will be estimated by computing $H_k(n)$ for several values of k , typically 1 to 5. Since the analytical form of the variance is known, we may combine these estimates by a weighted averaging procedure, *i.e.* $H(n) = \sum w_k H_k(n)$. For independent variables with the same average, the weight which minimizes the variance of the estimate of the average is a weight proportional to the inverse of the variance of the variable (see the appendix A for details), *i.e.*,

$$w_k = \frac{1/(Q_k + Var[\ln f(x)])}{\sum_{i=1}^m 1/(Q_k + Var[\ln f(x)])} \quad \text{for } k = 1, 2 \dots m \quad (2.26)$$

where w_k is the ideal weight of $H_k(n)$ when averaging $H(n)$. Such calculations may also be readily extended to compute the standard deviation of such an estimate (appendix A). Interestingly, two well defined limits exist here: (1) if $Var[\ln f(x)]$ is small, then the proper weighting will be

$$w_k = \frac{1/Q_k}{\sum_{k=1}^m 1/Q_k} \quad \text{for } k = 1, 2 \dots m \quad (2.27)$$

and, (2) if $Var[\ln f(x)]$ is large, then the proper weighting will be a flat function which will lead to a simple arithmetic average. Therefore, the best possible estimate of $H(n)$ from m estimates of $H_k(n)$ will always be bound by these two limiting averages. Further, if these two limiting averages converge in the given sampling, it is highly probable the estimate of $H(n)$ is also converged. We also note here that an intuitive sense of which regime best fits the given data can be discerned by inspecting the relative noise in plots of the m $H_k(n)$ estimates as a function of n (where n is the amount of simulation time in this application). If the $H_1(n)$ estimate noticeably suffers greater fluctuations than the other estimates, then the $Var[\ln f(x)]$ term must be small, since the Q_1 component is dominating relative variances of the estimates. However, if the m $H_k(n)$ estimates all appear graphically to have fluctuations of a similar magnitude, then the $Var[\ln f(x)]$ term must be large, and the simple arithmetic average is more appropriate. Such inspection of our data revealed $Var[\ln f(x)]$ to be small. As such, the weighted average determined by application of eqn 2.27 was taken in this work as our best possible estimate of $H(n)$.

2.2.5 Calculation of the excess energy, enthalpy, and free energy

The excess molar energy of a fluid is simply

$$\Delta E = \frac{1}{2} \frac{\rho}{\Omega^2} \int g(r, \omega^2) u(r, \omega^2) d\mathbf{r} d\omega^2 \quad (2.28)$$

where $u(r, \omega^2)$ is the interaction energy between two molecules with distance r and orientation determined by ω^2 . This quantity is straight forward to extract from the simulation, as it is merely one half of the interaction energy between the water molecule of interest with the rest of the system. The molar excess enthalpy can be obtained by approximating the $\Delta(PV)$ term. For the liquid phase, the PV term may be safely neglected, and for the gas phase, we may use the ideal gas equation of state $PV = NkT$ to derive an excellent approximation to the PV term analytically. Combined with the excess entropy, we find the excess free energy of the fluid may be expressed as

$$\Delta G = \Delta E + \Delta(PV) - Ts_e \quad (2.29)$$

as is typical.

2.2.6 The finite-difference method of entropy calculation

In order to generate reference data to examine the accuracy of the k'th NN method of evaluating the entropy expansion, we pursued a finite difference analysis of the solvation free energy, as computed from free energy perturbation theory (FEP). The finite-difference (FD) method of computing an entropy from FEP data proceeds by first noting that the entropy is the temperature derivative of the free energy, and then attempting to accurately estimate this slope [28], ie

$$-\Delta S(T) = \left\langle \frac{\partial \Delta G}{\partial T} \right\rangle_P = \frac{\Delta G(T + \Delta T) - \Delta G(T - \Delta T)}{2\Delta T} \quad (2.30)$$

This method relies on the assumption that the heat capacity of the system is independent of temperature in the range $[T - \Delta T, T + \Delta T]$ [29]. This assumption appears to be valid near room temperature with ΔT even as large as $50K$ [28]. Here, we use the Bennett acceptance ratio[30] method to calculate the excess free energy of liquid water at $T = 298 \pm 20K$, and then use FD to calculate the excess entropy at $T = 298K$. The details of this method are

included in the appendix A. This data allows for independent validation of the NN approach and the approximations therein.

2.2.7 Details of the simulation

Dynamics trajectories were generated using the Desmond molecular dynamics program [31]. A 25 Å cubic box of the TIP4P[32] water model was first equilibrated to 298K and 1 atm with Nose-Hoover[33; 34] temperature and Martyna-Tobias-Klein[35] pressure controls, followed by 30 ns NVT dynamics simulation with a Nose-Hoover[33; 34] temperature control. In order to integrate the equations of motion of the system, the RESPA[36] integrator was used, where the integration step was 2 fs for the bonded and the nonbonded-near interactions and 6 fs for the nonbonded-far interactions. Configurations were collected every 1.002 ps. The cut-off distance was 9 Å for the Van der Waals interaction, and the particle-mesh Ewald[37] method was used to model the electrostatic interactions. Similar simulations were performed for the SPC[38], SPC/E[39], TIP3P[32] and TIP4P-Ew[40] water models.

When extracting the solvent configurations to compute the PCF, we chose to only use the configurations of a distinguished solvent molecule with the rest of the system, instead of collecting statistics from all pairs of solvent molecules. Our primary interest in developing this technique was to later adapt the method to study the solvation of solutes; thus, we were interested in determining realistic estimates of the convergence of the technique when the isotropic symmetry of the fluid was not present. Such a protocol allows for an interrogation of the relative convergence properties of the various methods that might be obscured by the additional statistics offered by taking advantage of the symmetry of the system.

2.3 Results and discussion

2.3.1 The Shannon entropies

The NN estimates of the two dimensional orientational Shannon entropies $S^{[t_1, t_2]}$ of the TIP3P water model for the three shells are given in figures 2.1, 2.2, and 2.3. The results reported in these figures were generally representative of those results obtained for the other models. We see from the figures that the weighted average estimate of all the Shannon

entropies are converged over the course of the simulations. The results of all the one and two dimensional orientational Shannon entropies for each of the three shells for all the water models studied are given in Table 2.1. By application of formula 2.4 and 2.7, we computed the translational excess entropies and orientational excess entropies for all the water models studied. All the final results are shown in 2.2. From the table, we see that for the TIP4P model the excess entropy result from the NN method $-13.67e.u.$ is very close to experimental value $-14.1e.u.$ We also note excellent agreement between the excess entropies computed here and those derived from cell theory[41]. The agreement for the TIP3P and SPC models was slightly diminished compared with the other models, for reasons that will be explained later.

2.3.2 Convergence properties

We extensively compared the commonly employed histogram method to compute the orientational Shannon entropy to the NN method weighted average (2.4, 2.5, and 2.6). We see clearly that the NN method weighted average converges much faster than histogram method for shells 1 and 2. For shell 3, both methods give similar results. This is easily understood: for the first and second shells, the water molecules are highly correlated, and the histogram results will have a strong dependency on the bin size used to do the integration; however, for the third shell, there is little correlation, so the histogram method has similar convergence properties compared to the NN method.

Figures 2.7, 2.8, 2.9, 2.10, 2.11 depict the total orientational excess entropies as a function of simulation time from the various histogram estimates and the NN weighted average estimate. For all the models studied, the 10° histogram estimate (which is most commonly used currently [6; 10]) gave results closest to the NN estimate. However, for a bin size of 20° , the entropy result is biased away from the correct result, and for bin sizes of 5° and 2.5° , much longer simulation time would be needed to converge the results. Since ideal bin size is problem specific, it cannot be deduced unless other reference data is already known. Thus, the absence of such a parametric bias in the NN method is a notable advantage of the technique.

2.3.3 Error analysis

As described in the methods section, we calculated the variance associated with the weighted average of the NN estimates for each of the one and two dimensional Shannon entropies. Since the NN estimate is asymptotically unbiased, the error of the estimate is also given by the variance. We calculated the error based on the weighted average, which assumes $Var \ln f(x)$ is 0. However, even in the extreme cases where $Var \ln f(x)$ goes to infinity and the five NN estimates contribute equally to the average, the variance of the arithmetic average only differs slightly from weighted average, and they are within the error bar of each other, strongly indicating the convergence of these calculations (Figures 2.12 and 2.13).

2.3.4 The radial dependence of orientational Shannon entropy

We calculated the orientational Shannon entropies in three radial regions, assuming the orientational distribution would be independent of r in each sub-region. To validate this approximation, we calculated the orientational Shannon entropies at different intervals of r from 2.5 to 4.0 Å. Typical Shannon entropies $S^{[t_1, t_2]}$ at different value of r are shown in Figure 2.14.

We see from the figure that the Shannon entropy increases as the distance between the two water molecules r increases, and goes to zero when r is sufficiently large. Additionally, the change of the Shannon entropy with respect to r is smooth in the respective first and second hydration shells. Because of the slow variation of the orientational Shannon entropy with respect to r , the sum of the orientational excess entropy at each interval will differ from the sum of the orientational excess entropy of the three shells only by at most $0.5e.u.$, which is within statistical uncertainty of the calculation. Thus, this approximation was not a large source of error in these calculations.

2.3.5 Inclusion of $g(\theta_1, \chi_1)$ in the factorization

The factorization of the PCF used here differs from the more common formulation[6] by the explicit inclusion of $g(\theta_1, \chi_1)$ and $g(\theta_2, \chi_2)$. The distribution functions $g(\theta_1) * g(\chi_1)$ and $g(\theta_1, \chi_1)$ for the TIP4P model are shown on Figures 2.15 and 2.16. Careful inspection of these figures suggests that $g(\theta_1, \chi_1)$ differs from $g(\theta_1)g(\chi_1)$ quantitatively, which is sup-

ported by the two dimensional Shannon entropy $S^{[\theta_1, \chi_1]}$ differing significantly from the sum of $S^{[\theta_1]}$ and $S^{[\chi_1]}$. For example, for the TIP4P model the first shell Shannon entropy of $S^{[\theta_1, \chi_1]}$ is -1.21, while $S^{[\theta_1]}$ is -0.34 and $S^{[\chi_1]}$ is -0.29. This result indicated a non-negligible correlation between χ_1 and θ_1 , which suggested that the explicit inclusion of $g(\theta_1, \chi_1)$ and $g(\theta_2, \chi_2)$ in our factorization would lead to greater quantitative precision. This also explains why our excess entropy result for the TIP4P model ($-13.67e.u.$) is about $1.5e.u.$ more negative than the previously reported value ($-12.2e.u.$)[6], which is in better agreement with both the FD estimate of the entropy of the model and the experimental estimate of liquid water.

2.3.6 Comparison of free energy results

From these simulations, we computed the excess molar energies and excess free energies of the various water models. The results of these calculations for all models studied are listed in table 2.3 along side the relevant literature values. The excess free energies we have obtained here show excellent agreement (within 0.5 kcal/mol uniformly) with the high precision FEP results obtained by Shirts *et. al.*[42]. Interestingly, the TIP4P model gives results closest to the experimental quantities.

The SPC/E, TIP4P, and TIP4P-Ew models all give free energy results somewhat closer to the Shirts[42] results than the other models. This may not be accidental. In our calculations, the higher order multi-particle correlation entropies were ignored. There is some literature precedence expecting these higher order contributions to the excess entropy to vanish at the temperature of solid-liquid phase transition[43; 44]. Recently, Saija has shown that for the TIP4P model, the temperature of maximum density (TMD) coincides with the temperature where higher order contributions to the entropy should vanish[13]. Studies of temperature dependence of the densities of the different water models studied here[45] have shown that the TMD of the TIP4P model occurred at 258K, the TMD of the SPC/E model occurred at 235K[46], the TMD of the TIP4P-Ew model occurred at 272K[40], and the density of the SPC and TIP3P models increases monotonically as temperature decreases in the range [220, 370][45]. This indicates, for the TIP3P and SPC models, multi-particle correlation entropy may contribute more to the total entropy than for the other models,

which may be why our quantitative accuracy for them is somewhat diminished. However, the molecular detail afforded by this technique in yielding both a value of the entropy and a physical interpretation of its meaning, in terms of the fluid structure implied by the shape of the PCF, gives it a comparative advantage over techniques such as FEP, which will generally only yield a value of the entropy without any additional molecular understanding of the system.

2.3.7 Entropy calculation from FD method

We calculated the excess free energy of water at temperature $298 \pm 20K$ with the Bennett acceptance ratio[30] method, and obtained entropies at $298K$ by the FD formula. The results are presented in Table 2.4. The excess entropies computed from the FD method are consistently larger in magnitude than those computed from the NN method, consistent with us neglecting the contributions from the higher order terms of the expansion.

As in the proceeding section, the NN and FD excess entropies of the SPC/E water are in very close agreement; however, the agreement of the NN and FD entropies of the SPC and TIP3P models is much poorer. We again expect the reason for this discrepancy to be due to the TMD of the SPC/E model being close to the range of temperatures treated in this study, while the TMDs of the SPC and TIP3P models fall well outside this range. Thus, the higher order terms of the entropy expansion are expected to make larger contributions to the excess entropies for the SPC and TIP3P models versus the contribution made to the excess entropy of the SPC/E water.

2.4 Conclusion

Our results indicate that the NN method of computing entropies in the liquid state offers several compelling advantages over the more common histogram approaches, including (1) much faster convergence for a given amount of simulation data; (2) an intuitive error bound for the uncertainty of the calculation without resorting to block averaging or bootstrapping techniques, which may be problematic to apply to estimators of the entropy; and (3) not relying on empirically tuned parameters, such as the histogram bin width, which may

bias the results in an unpredictable fashion. We also found that inspection of the limiting behaviours of $Var \ln f(x)$ may be used to both analyze the convergence of the given calculation, and develop the best possible estimate of the entropy given a set of calculated $H_k(n)$. Although we also found that a judicious choice of the histogram bin width may mitigate these advantages, such a choice is difficult to make without prior knowledge of the properties of the limiting distribution, which may not be available when new problems are investigated.

Our alternative factorization of the water-water correlation function, which explicitly included correlations between the angle formed by the water dipole vector and the intermolecular axis with the angle of rotation of the water molecule about its dipole vector, was found to increase the agreement of results obtained by the entropy expansion with those obtained by less approximate methods, such as FEP and the FD benchmark calculations. This result suggests that this contribution should not be ignored in future studies of the excess entropy of liquid water and other fluids.

water models		$\mathcal{S}^{[t_1, t_2]}$	$\mathcal{S}^{[t_1, \chi_1]}$	$\mathcal{S}^{[t_1, \chi_2]}$	$\mathcal{S}^{[\chi_1, \chi_2]}$	$\mathcal{S}^{[t_1]}$	$\mathcal{S}^{[\chi_1]}$
Shell1	TIP4P	-1.33	-1.21	-1.15	-1.02	-0.34	-0.29
	SPC	-1.67	-1.28	-1.24	-0.89	-0.50	-0.27
	TIP3P	-1.65	-1.16	-1.14	-0.74	-0.47	-0.23
	SPC/E	-1.70	-1.32	-1.29	-0.94	-0.51	-0.29
	TIP4P-Ew	-1.44	-1.29	-1.23	-1.05	-0.39	-0.30
Shell2	TIP4P	-0.59	-0.44	-0.46	-0.38	-0.10	-0.10
	SPC	-0.69	-0.42	-0.46	-0.30	-0.11	-0.09
	TIP3P	-0.60	-0.29	-0.34	-0.18	-0.09	-0.06
	SPC/E	-0.71	-0.46	-0.50	-0.33	-0.13	-0.10
	TIP4P-Ew	-0.68	-0.51	-0.53	-0.38	-0.12	-0.12
Shell3	TIP4P	-0.010	-0.007	-0.002	-0.003	-0.001	-0.000
	SPC	-0.014	-0.007	-0.005	-0.001	-0.002	-0.000
	TIP3P	-0.015	-0.003	-0.003	-0.001	-0.002	-0.000
	SPC/E	-0.013	-0.007	-0.005	-0.003	-0.001	-0.000
	TIP4P-Ew	-0.012	-0.007	-0.004	-0.001	-0.001	-0.000

Note: $t = \frac{\pi}{2}(\cos(\theta) + 1)$, all these entropies are unitless.

Table 2.1: Orientational Shannon entropies of the five water models

	EXP	TIP4P	TIP3P	SPC	SPC/E	TIP4P-Ew
$s_{trans}^{(2)}$	–	-3.15(3.14 ^a)	-2.99	-2.99	-3.19	-3.33
$s_{orient}^{(2)}$	–	-10.52(9.10 ^a)	-8.58	-10.20	-11.53	-11.76
$s_{ex}^{(2)}$	–	-13.67(-12.2 ^a)	-11.57	-13.19	-14.72	-15.09
s_{ex}	-14.05 ^b	-14.32 ^c	-13.36 ^c	-14.01 ^c	-14.79 ^c	-14.99 ^c

Entropies in cal/mol K (e.u.).

^adata from Lazaridis [6]

^bdata from Wagner [47]

^cdata from Henschman by cell theory [41]

Table 2.2: Comparison of entropy results from the NN method and cell theory

water models	TIP4P	TIP3P	SPC	SPC/E	TIP4P-Ew
excess energy	-9.85	-9.49	-9.90	-11.08	-10.91
excess enthalpy	-10.43	-10.07	-10.48	-11.66(-10.48 ^a)	-11.49(-10.45 ^b)
excess enthalpy*	-10.41	-10.09	-10.47	-11.69(-10.51 ^a)	-11.61(-10.57 ^b)
excess entropy from NN	-13.67	-11.57	-13.19	-14.72	-15.09
excess entropy**	-14.43	-13.39	-14.46	-15.57	-15.53
excess free energy from NN	-6.36	-6.63	-6.55	-7.27(-6.09 ^a)	-7.00(-5.96 ^b)
excess free energy*	-6.11	-6.10	-6.16	-7.05(-5.87 ^a)	-6.98(-5.94 ^b)
excess free energy from exp	-6.33				
excess enthalpy from exp	-10.52				

Energies in kcal/mol, entropies in cal/mol K (e.u.)

* results from Shirts[42]

** results from Shirts[42] by extracting enthalpy from free energy

^aInclude polarization correction [39]

^bInclude polarization correction [40]

Table 2.3: Results for the energy, enthalpy, and entropy of liquid water from various methods

water models	TIP4P	TIP3P	SPC	SPC/E	TIP4P-Ew
excess free energy at 278K	-6.35**	-6.21(-6.24 ^a)	-6.36(-6.39 ^a)	-7.19(-7.23 ^a)	–
excess free energy at 298K	-6.03**	-5.95	-6.06	-6.89	–
excess free energy at 318K	-5.73**	-5.71(-5.69 ^a)	-5.80(-5.78 ^a)	-6.66(-6.62 ^a)	–
excess entropy from FD	-15.2**	-13.8(±0.8 ^b)	-15.2(±0.8 ^b)	-15.3(±0.8 ^b)	–
excess entropy from NN	-13.67	-11.57	-13.19	-14.72	-15.09
excess entropy from FEP*	-14.43	-13.39	-14.46	-15.57	-15.53

Energies in kcal/mol, entropies in cal/mol K (e.u.)

** results from Franz Saija[13]

* results from Shirts[42] by extracting enthalpy from free energy

^a results in parentheses includes constant pressure correction(appendix A)

^b indicates the error associated with the entropy

Table 2.4: Entropy results from FD method and comparison with other methods

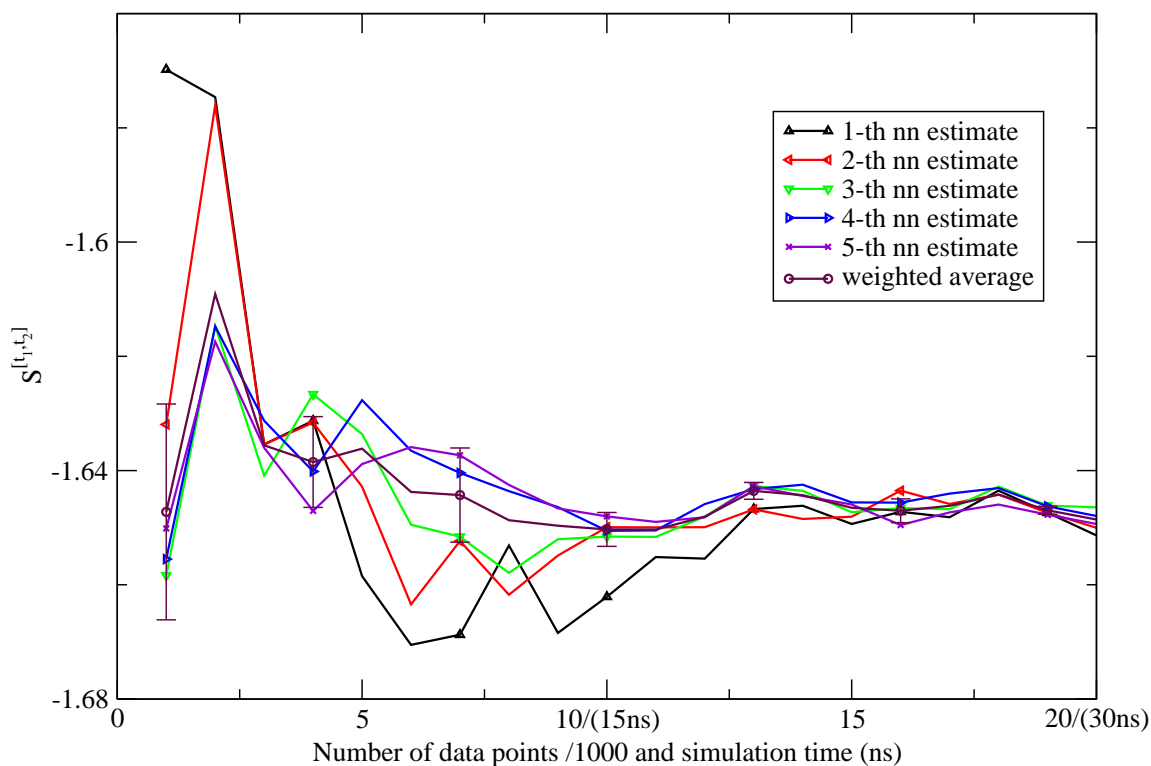


Figure 2.1: The first shell orientational Shannon entropy $S^{[t_1, t_2]}$ for the TIP3P model as a function of the number of data points (labeled on the horizontal axis in front of “/” in units of 1000) and the corresponding simulation time (labeled on the horizontal axis in parentheses) using the NN method. The weighted average estimate and the associated error bar were also depicted.

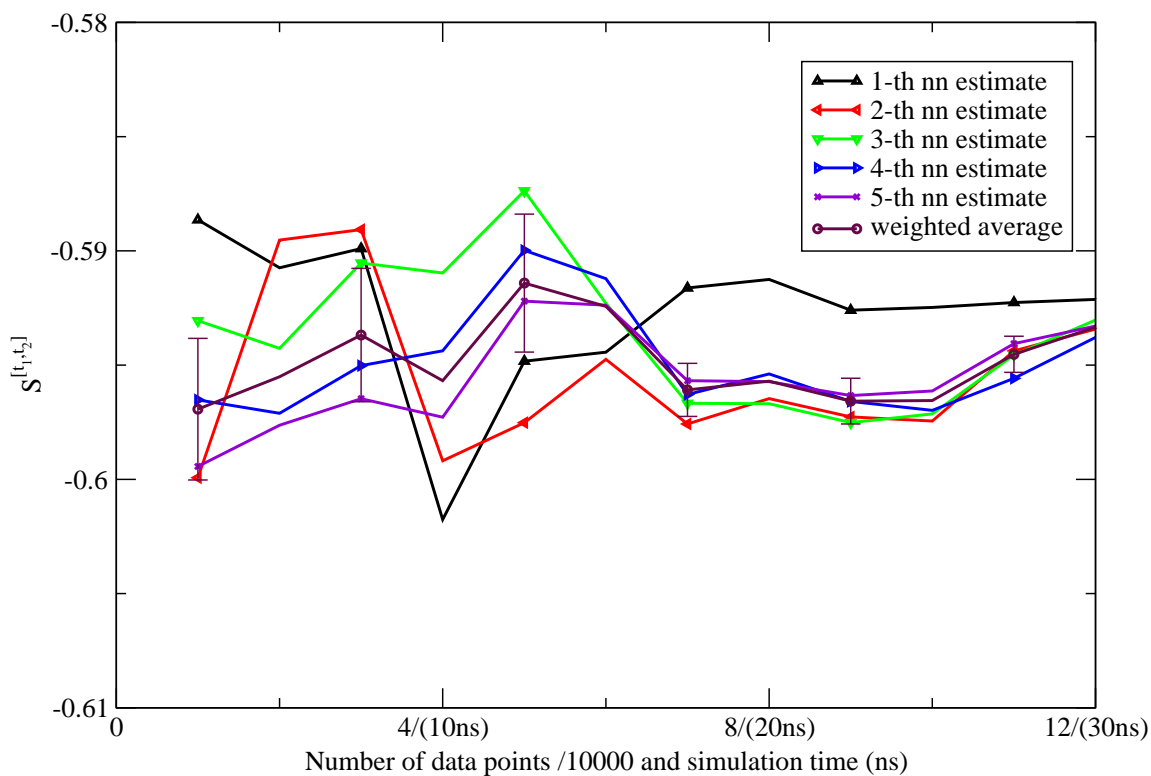


Figure 2.2: The second shell orientational Shannon entropy $S^{[t_1, t_2]}$ for the TIP3P model as a function of the number of data points (labeled on the horizontal axis in front of "/" in units of 10000) and the corresponding simulation time (labeled on the horizontal axis in parentheses) using the NN method. The weighted average estimate and the associated error bar were also depicted.

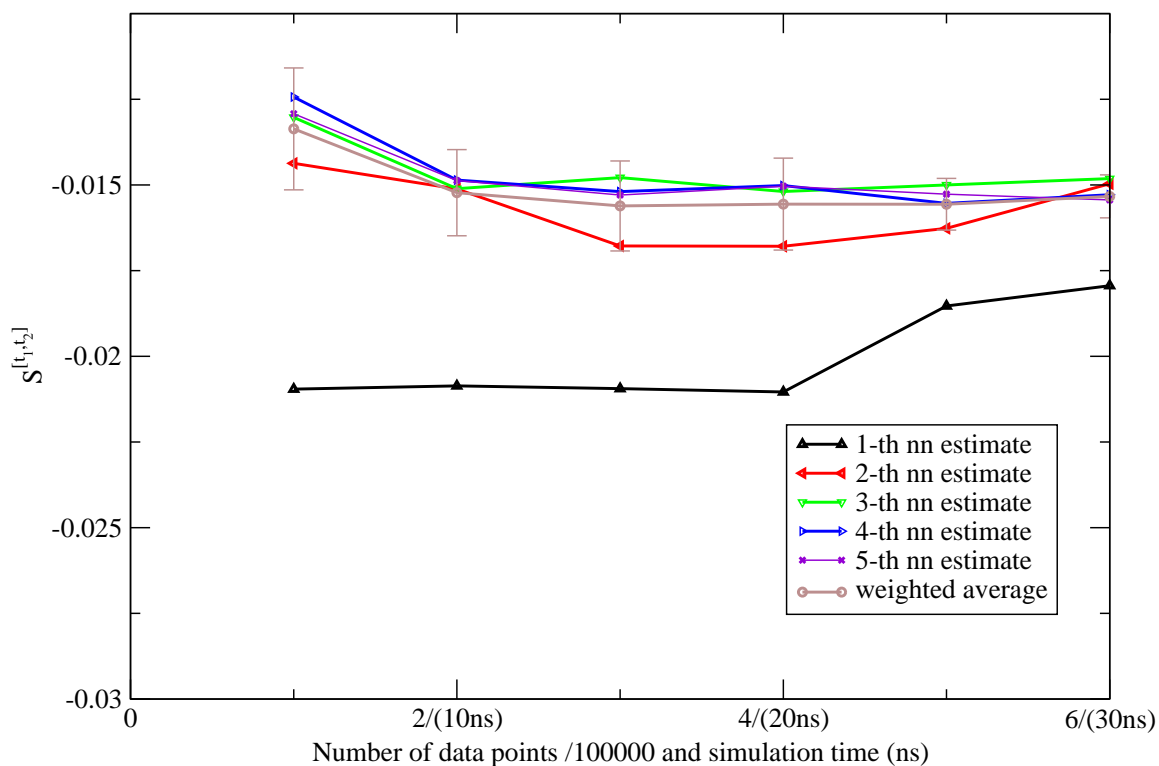


Figure 2.3: The third shell orientational Shannon entropy $S^{[t_1, t_2]}$ for the TIP3P model as a function of the number of data points (labeled on the horizontal axis in front of "/" in units of 100000) and the corresponding simulation time (labeled on the horizontal axis in parentheses) using the NN method. The weighted average estimate and the associated error bar were also depicted.

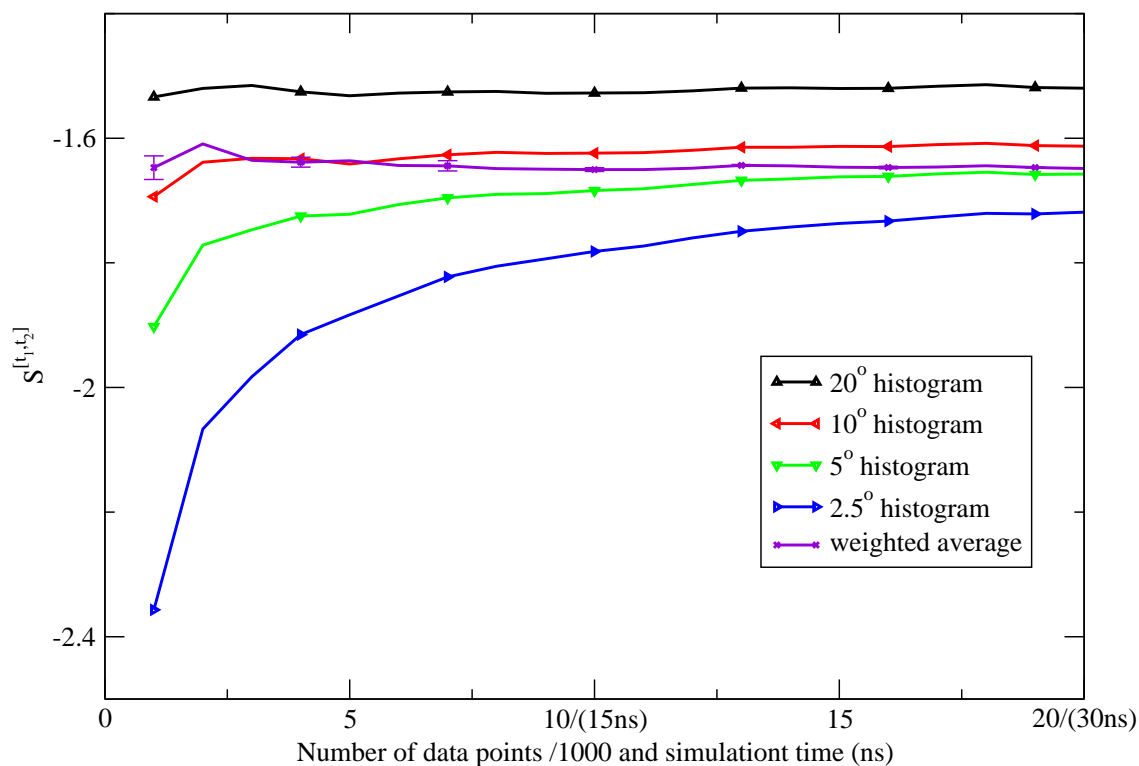


Figure 2.4: The first shell orientational Shannon entropy $S^{[t_1, t_2]}$ for the TIP3P model as a function of the number of data points (labeled on the horizontal axis in front of "/" in units of 1000) and the corresponding simulation time (labeled on the horizontal axis in parentheses) using histogram method. The weighted average of the NN estimates and the associated error bar were also depicted.

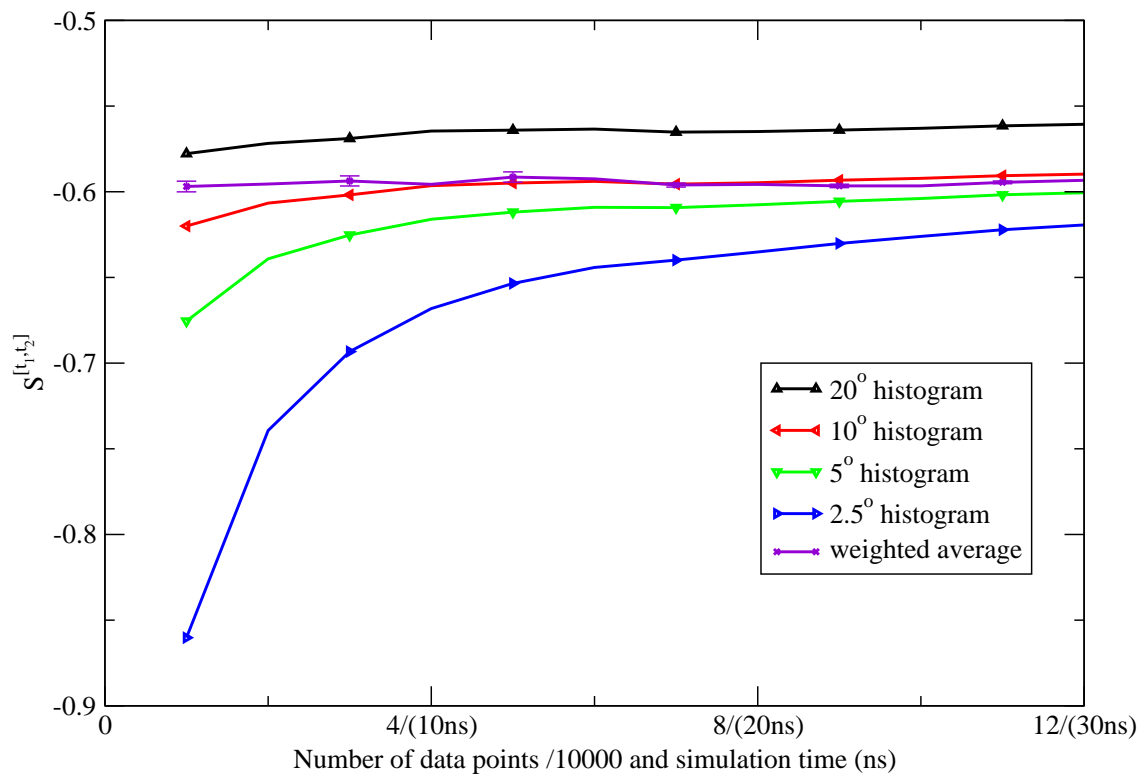


Figure 2.5: The second shell orientational Shannon entropy $S^{[t_1, t_2]}$ for the TIP3P model as a function of the number of data points (labeled on the horizontal axis in front of "/" in units of 10000) and the corresponding simulation time (labeled on the horizontal axis in parentheses) using histogram method. The weighted average of the NN estimates and the associated error bar were also depicted.

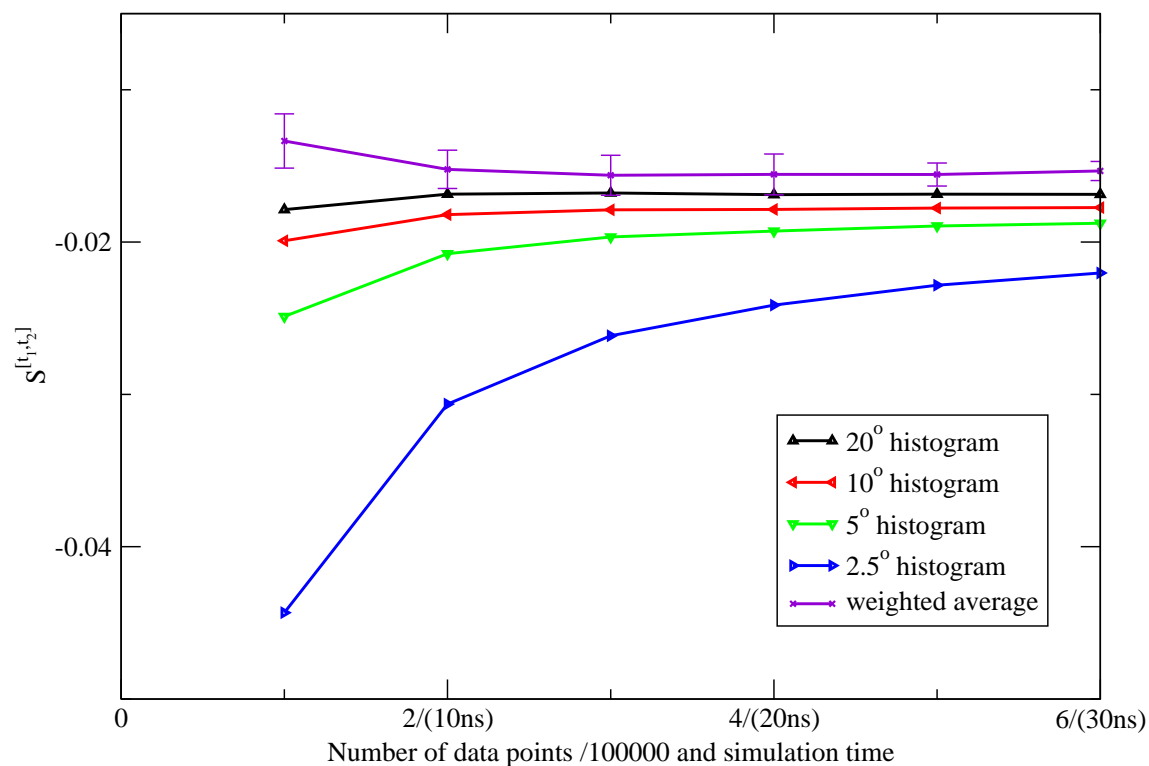


Figure 2.6: The third shell orientational Shannon entropy $S^{[t_1, t_2]}$ for the TIP3P model as a function of the number of data points (labeled on the horizontal axis in front of "/" in units of 100000) and the corresponding simulation time (labeled on the horizontal axis in parentheses) using histogram method. The weighted average of the NN estimates and the associated error bar were also depicted.

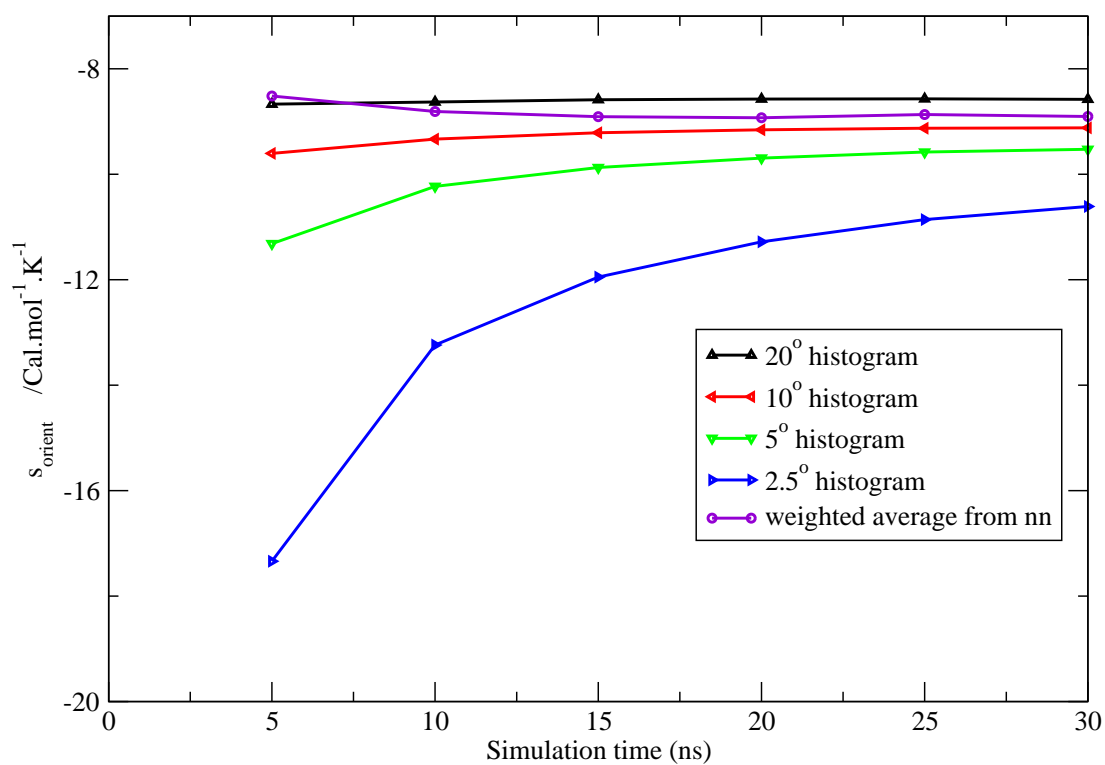


Figure 2.7: Total orientational excess entropy as a function of simulation time from the NN method and histogram method with different bin width for the TIP3P model.

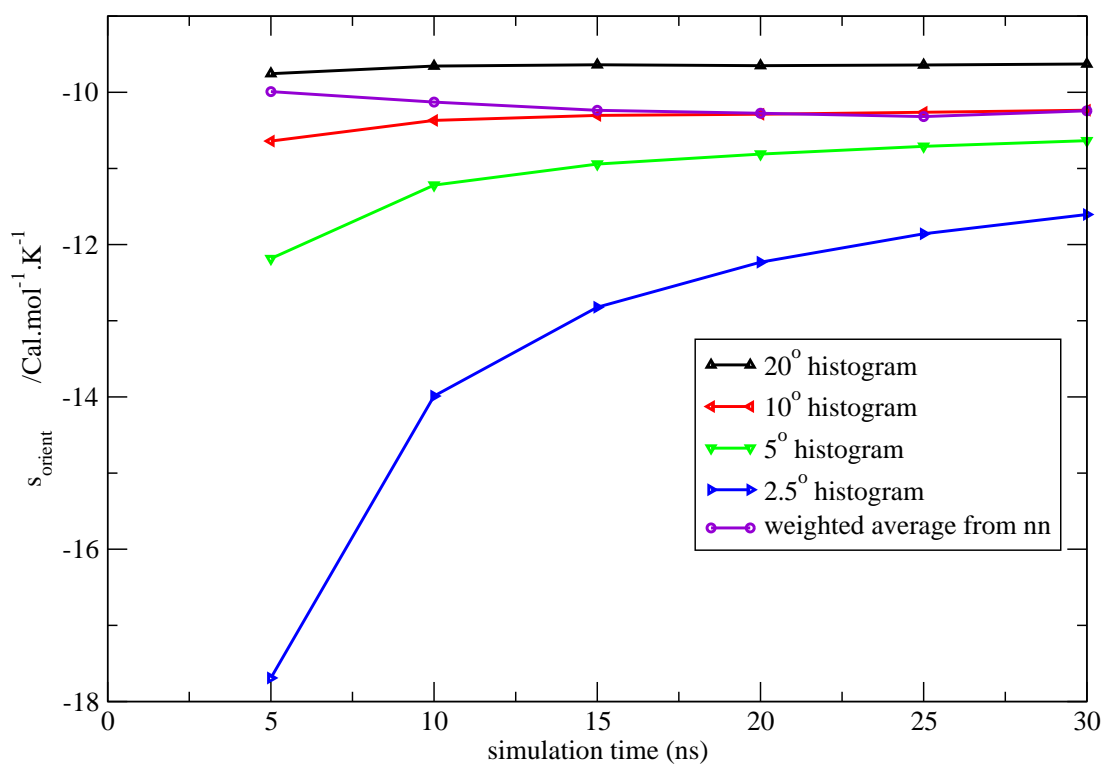


Figure 2.8: Total orientational excess entropy as a function of simulation time from the NN method and histogram method with different bin width for the SPC model.

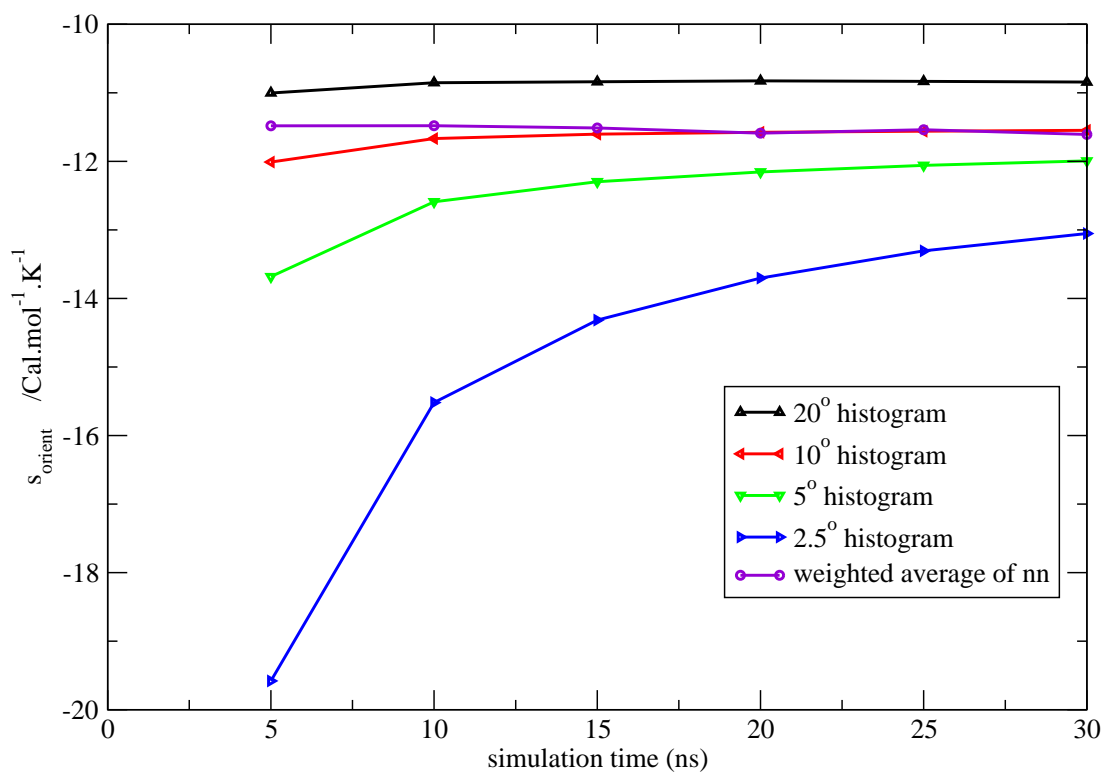


Figure 2.9: Total orientational excess entropy as a function of simulation time from the NN method and histogram method with different bin width for the SPC/E model.

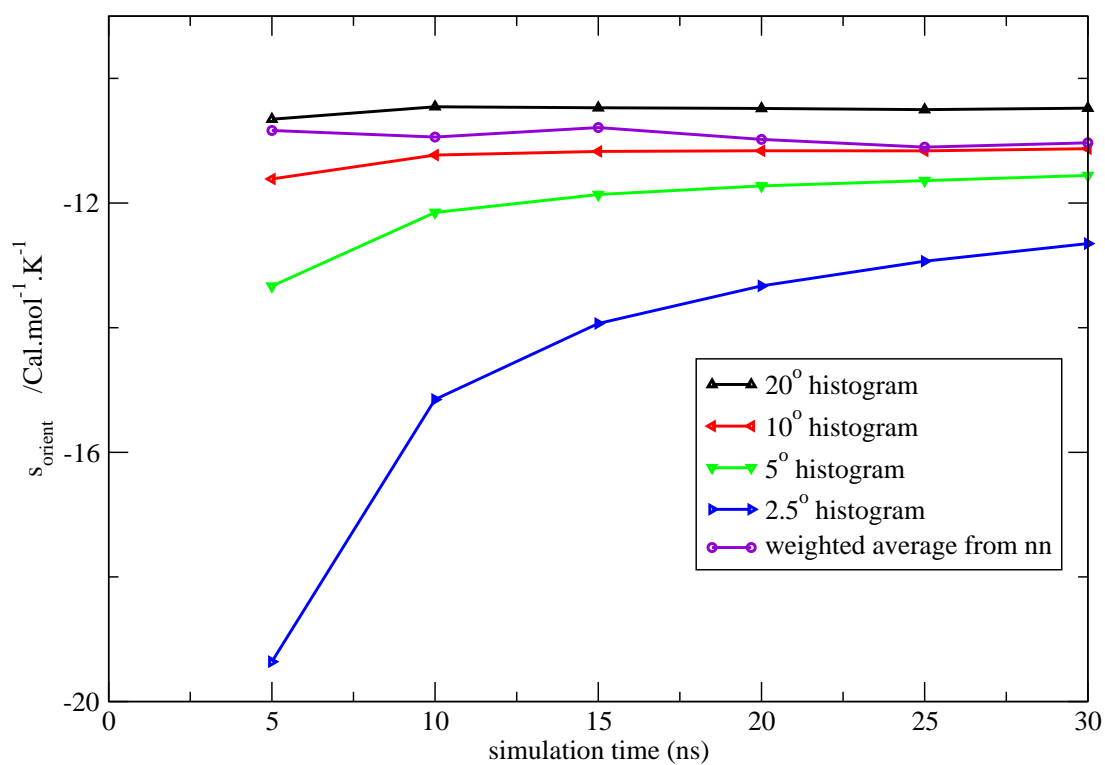


Figure 2.10: Total orientational excess entropy as a function of simulation time from the NN method and histogram method with different bin width for the TIP4P model.

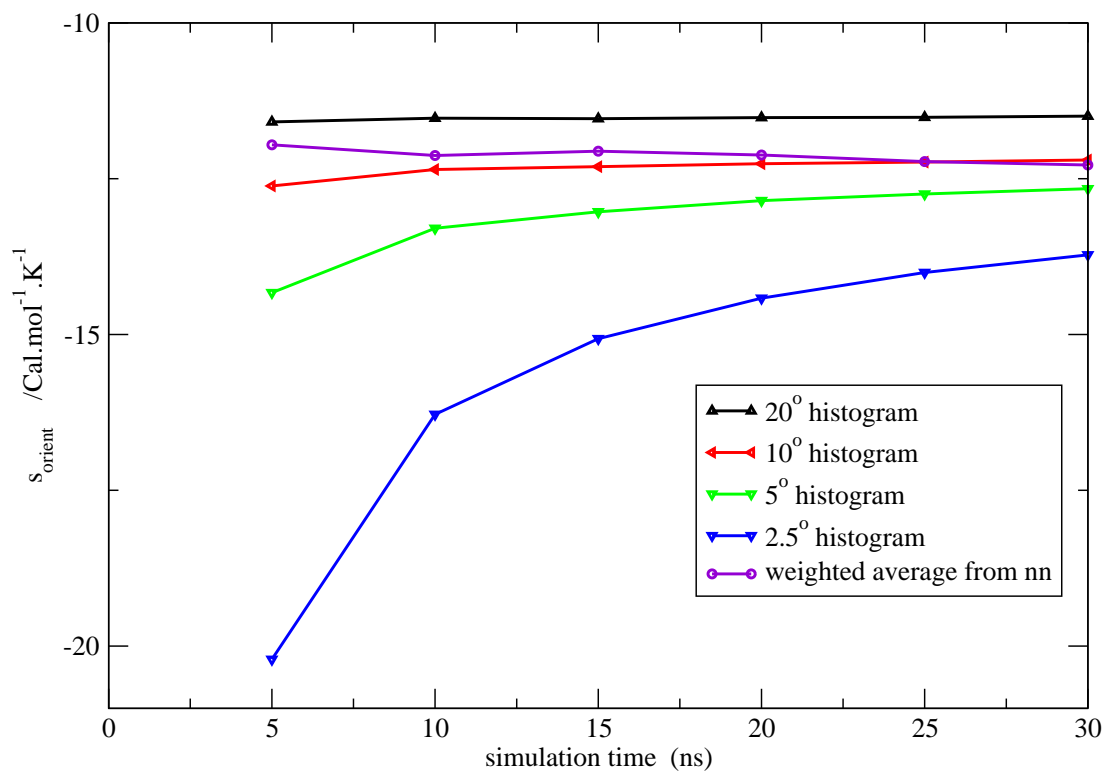


Figure 2.11: Total orientational excess entropy as a function of simulation time from the NN method and histogram method with different bin width for the TIP4P-Ew model.

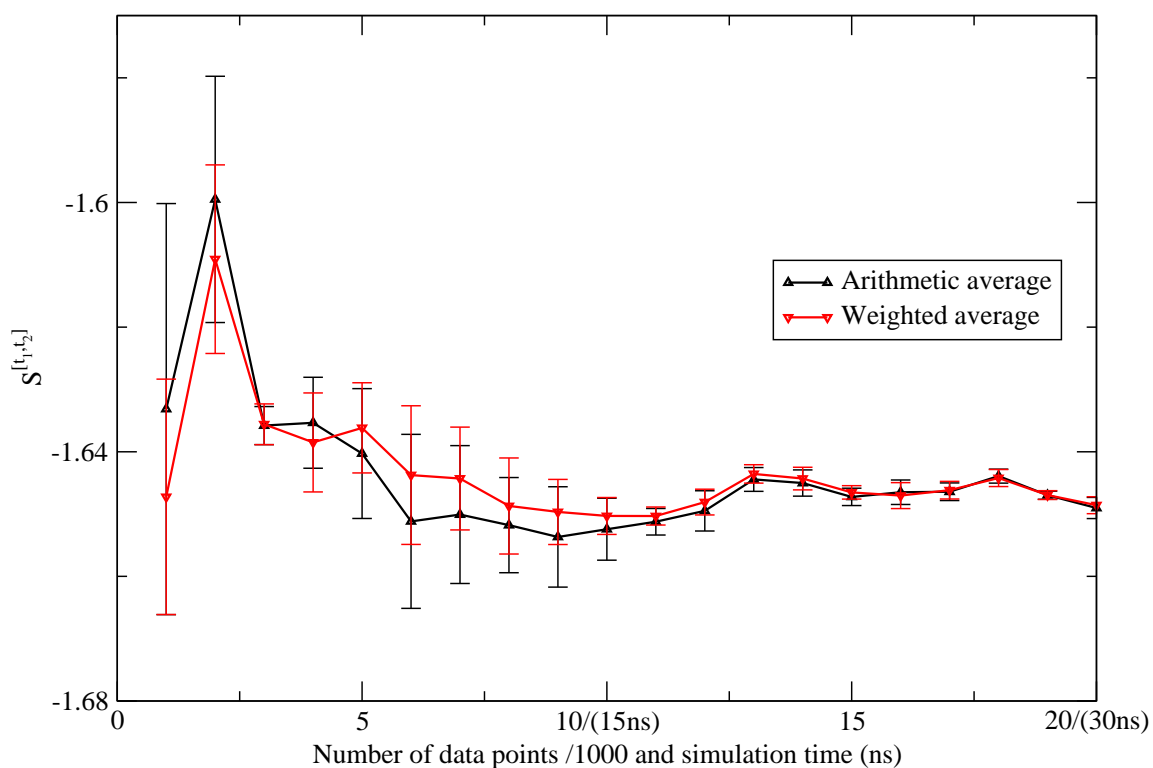


Figure 2.12: Comparison between the arithmetic average and the weighted average of the NN estimates for the first shell Shannon entropy $S^{[t_1, t_2]}$ for the TIP3P model. They are within the error bar of each other.

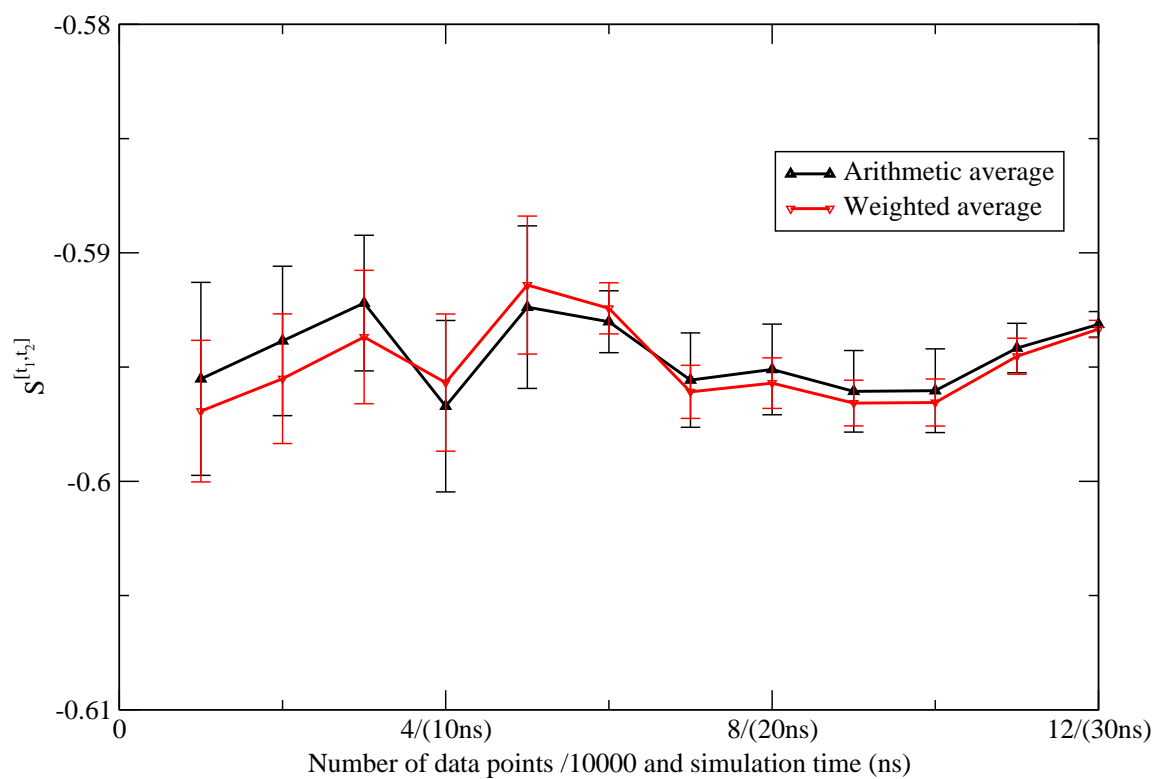


Figure 2.13: Comparison between the arithmetic average and the weighted average of the NN estimates for the second shell Shannon entropy $S^{[t_1, t_2]}$ for the TIP3P model. They are within the error bar of each other.

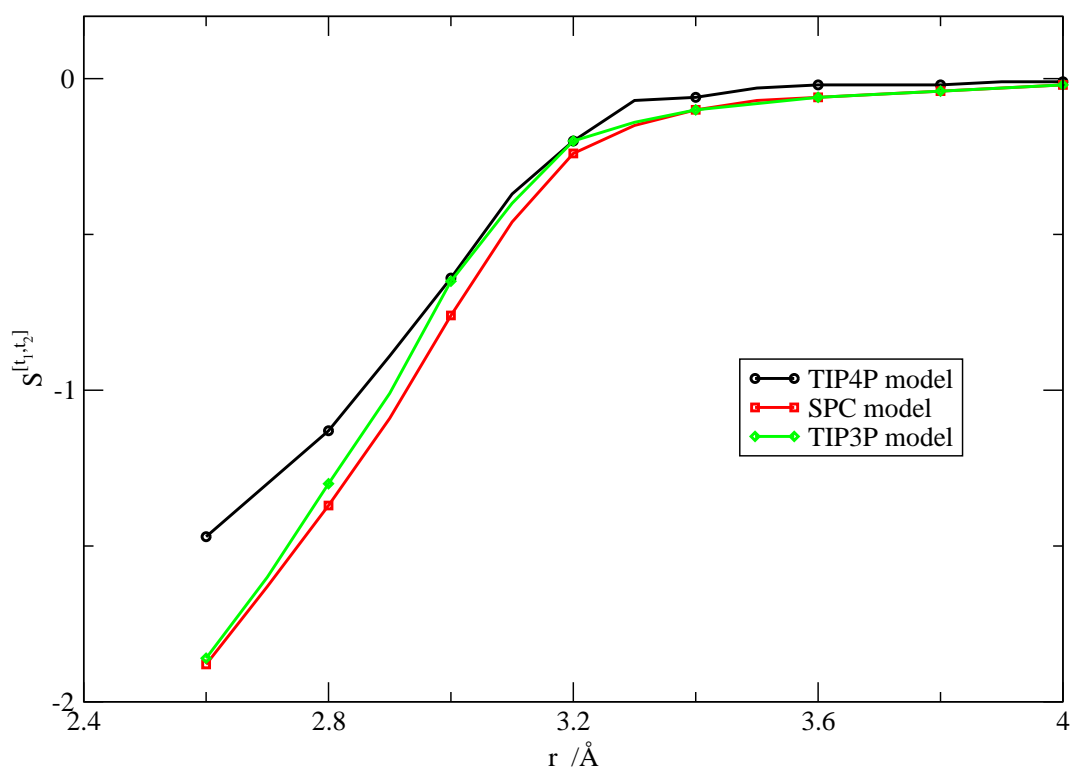


Figure 2.14: Orientational Shannon entropy $S^{[t_1, t_2]}$ as a function of r for the various water models.

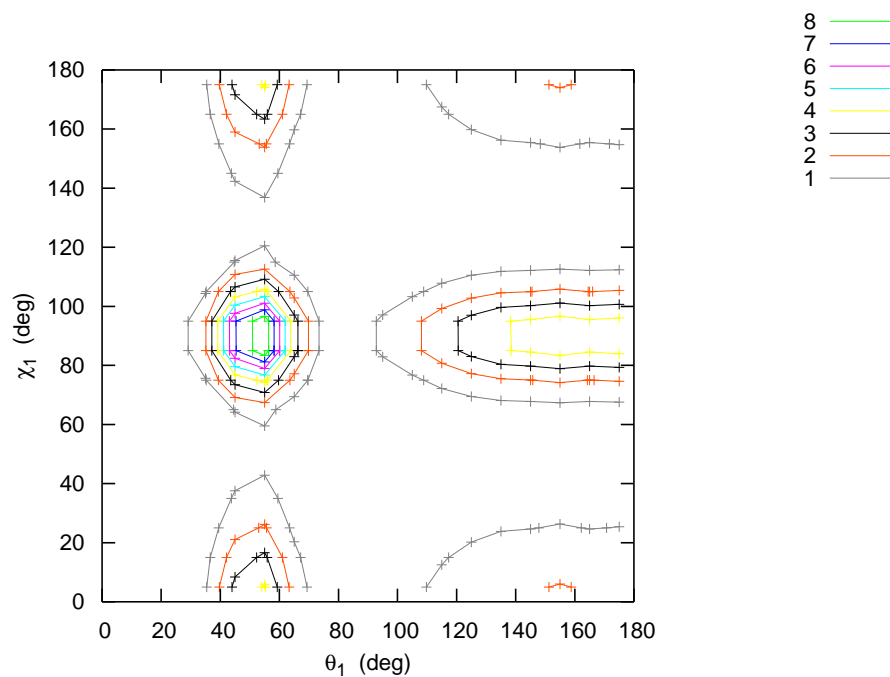


Figure 2.15: Products of one dimensional marginal distribution function $g(\theta_1) * g(\chi_1)$ for the TIP4P model in the first shell.

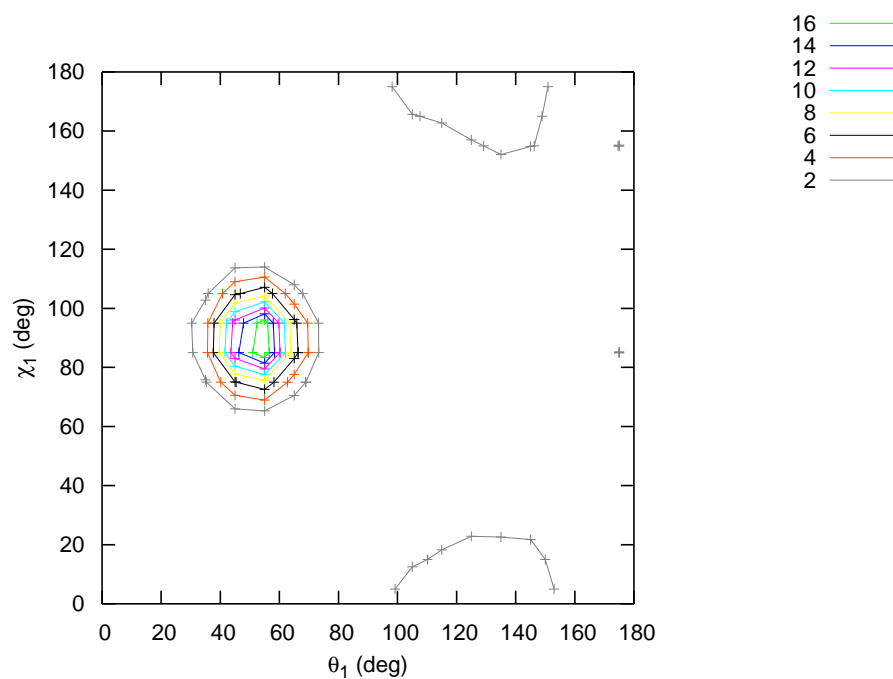


Figure 2.16: Two dimensional marginal distribution function $g(\theta_1, \chi_1)$ for the TIP4P model in the first shell.

Chapter 3

A displaced-solvent functional analysis of model hydrophobic enclosures

Abstract

Calculation of protein-ligand binding affinities continues to be a hotbed of research. Although many techniques for computing protein-ligand binding affinities have been introduced—ranging from computationally very expensive approaches, such as free energy perturbation (FEP) theory; to more approximate techniques, such as empirically derived scoring functions, which, although computationally efficient, lack a clear theoretical basis—there remains pressing need for more robust approaches. A recently introduced technique, the displaced-solvent functional (DSF) method, was developed to bridge the gap between the high accuracy of FEP and the computational efficiency of empirically derived scoring functions. In order to develop a set of reference data to test the DSF theory for calculating absolute protein-ligand binding affinities, we have pursued FEP theory calculations of the binding free energies of a methane ligand with 13 different model hydrophobic enclosures of varying hydrophobicity. The binding free energies of the methane ligand with the various hydrophobic enclosures were then recomputed by DSF theory and compared with the

FEP reference data. We find that the DSF theory, which relies on no empirically tuned parameters, shows excellent quantitative agreement with the FEP. We also explored the ability of buried solvent accessible surface area and buried molecular surface area models to describe the relevant physics, and find the buried molecular surface area model to offer superior performance over this dataset.

3.1 Introduction

Calculation of relative and absolute protein-ligand binding affinities continues to be an active hotbed of research in the field of computational biophysics.[48; 49; 50; 51] Although many techniques for computing protein-ligand binding affinities have been introduced—ranging from computationally very expensive *ab initio* approaches, such as free energy perturbation (FEP) theory; to more approximate techniques, such as empirically derived scoring functions, which, although computationally efficient, lack a clear theoretical basis—there remains a pressing need for more robust approaches. A recently introduced technique, the displaced-solvent functional (DSF) method was developed to bridge the gap between the high accuracy of FEP and the computational efficiency of empirically derived scoring functions.[2] This technique proceeds by first using explicitly solvated molecular dynamics simulations of a protein conformation which is complementary to a given ligand series (or, in some cases, a protein-ligand complex which can be used to build the remaining members of the series) to map out the approximate thermodynamic properties of water molecules solvating various regions of the protein active site; second, constructing a DSF to compactly represent this information; and third, computing the relative binding affinities of congeneric ligands for the given receptor by correlating the relative binding affinities of the congeneric ligands with the excess chemical potential of the solvent that is evacuated from the active site by the binding of the ligand.

This method has shown great promise in a number of pharmaceutically relevant applications such as accurately describing the relative binding thermodynamics of proteases, kinases, PDZ domain, and GPCR inhibitors; elucidating the role of hydration in kinase binding specificity; and offering novel qualitative insights into PCSK9-peptide binding

kinetics.[2; 3; 52; 53; 54; 55; 56; 57] However, despite the wide range successful applications of the technique to describe and explain experimental binding data, the physical-chemical basis of the DSF method has not yet been fully clarified. In this Chapter, the DSF approach is derived from first principles and the physical-chemical basis of the technique is clarified. Further, this derivation elucidates the key approximations of the method, which facilitates an understanding of when the technique is expected to succeed and fail. In order to develop a set of reference data to test the DSF theory for calculating absolute protein-ligand binding affinities, we have pursued FEP theory calculations of the binding free energies of a methane ligand with 13 different types of model hydrophobic enclosures of varying hydrophobicity. The binding free energies of the methane ligand with the various hydrophobic enclosures were then recomputed by the DSF theory presented herein and the results of the calculations were compared with the FEP reference data. We find that the DSF theory predictions, which rely on no empirically tuned parameters, show excellent quantitative agreement with the FEP results (root-mean-square error of 0.40 kcal/mol and an R^2 value of 0.95). Thus, DSF theory may offer, for systems that satisfy the necessary approximations, a method of calculating absolute binding affinities with FEP-like accuracy at only a small fraction of the computational expense. A further point is that the DSF approach can be unambiguously converged with current hardware capabilities, whereas convergence becomes quite challenging for FEP and related methods when applied to complex problems like protein-ligand binding (as opposed to the model systems studied in this paper).

3.2 Methods

3.2.1 Derivation of the displaced solvent functional approach to computing protein ligand binding free energies

It is well known that the binding free energy of a small molecule for its cognate protein receptor can be computed as

$$\Delta G_{bind}^o = -RT \ln \left[\frac{C_o}{8\pi^2} \frac{\int e^{-[(U(r_{PL})+W(r_{PL}))/RT]} d\mathbf{r}_{PL}}{\int e^{-[(U(r_P)+W(r_P))/RT]} d\mathbf{r}_P \int e^{-[(U(r_L)+W(r_L))/RT]} d\mathbf{r}_L} \right] \quad (3.1)$$

where the subscript P represents the protein in the unbound state, the subscript L

represents the ligand in the unbound state, the subscript PL represents the protein and ligand in their bound state, R is the gas constant, C_o is the standard concentration, U is the interaction energy term, and W represents the solvation free energy terms.[48] From this expression one can readily derive

$$\begin{aligned} \Delta G_{bind}^o &= \langle U_{PL} \rangle_{PL} - \langle U_P \rangle_P - \langle U_L \rangle_L \\ &+ \langle W_{PL} \rangle_{PL} - \langle W_P \rangle_P - \langle W_L \rangle_L - T\Delta S_{config}^o \end{aligned} \quad (3.2)$$

where the brackets ($\langle \rangle$) imply Boltzmann weighted averages over the specified ensemble, the changes of the configurational entropies of the protein and the ligand after binding have been grouped in a single term ($-T\Delta S_{config}^o$), and the terms related to the change in the interaction energies (U) and solvation free energies (W) of the protein and the ligand are enumerated explicitly. We note here that the $-T\Delta S_{config}^o$ term may be made arbitrarily small in equation 3.2 by first computing the free energy of restraining internal and relative degrees of freedom of the protein and the ligand to some appropriately chosen reference state by FEP, thermodynamics integration, or any other suitable ab initio approach, and then computing the binding free energy of the protein and ligand after these restraints have been removed.[58; 59]

Equation 3.2, although complete, has poor convergence properties since it is a series of very large terms that sum to a very small number. Thus, each individual term must be computed to very high accuracy and precision. This may in practice be more difficult than sampling Equation 3.1 directly, for example by FEP. However, we have made a series of observations in our recent work[2; 3] that suggest a path to improve the convergence of this expression.

The first observation is that the protein-ligand interaction energy (U_{PL}) can be expanded into an intra-protein term, a protein-ligand interaction term, and an intra-ligand term:

$$\langle U_{PL} \rangle = \langle U_P \rangle_{PL} + \langle U_{P-L} \rangle_{PL} + \langle U_L \rangle_{PL} \quad (3.3)$$

where the first term (U_P) is the intra-protein interaction energy, the second term (U_{P-L}) is the protein-ligand interaction energy, and the third term (U_L) is the intra-ligand interac-

tion energy. Therefore,

$$\begin{aligned} \Delta G_{bind}^o &= \langle U_P \rangle_{PL} + \langle U_{P-L} \rangle_{PL} + \langle U_L \rangle_{PL} - \langle U_P \rangle_P - \langle U_L \rangle_L \\ &\quad + \langle W_{PL} \rangle_{PL} - \langle W_P \rangle_P - \langle W_L \rangle_L - T\Delta S_{config}^o \end{aligned} \quad (3.4)$$

We will assume in this work that the loss of conformational entropy of the protein and ligand is compensated by the ligand and the strain energy incurred by the protein and ligand upon binding. For example a ligand with freely rotatable bonds binding to a protein will generally induce little protein strain energy, but will lose a great deal of conformational entropy upon binding. Conversely, a highly rigid ligand, which will avoid such entropic penalties, will likely require substantial induced fit of the protein, which will in turn increase the strain energy of the protein upon binding. Posed formally, this argument suggests

$$\langle U_P \rangle_{PL} + \langle U_L \rangle_{PL} - \langle U_P \rangle_P - \langle U_L \rangle_L - T\Delta S_{config}^o \approx 0 \quad (3.5)$$

In turn, equation 3.4 may be rewritten as

$$\begin{aligned} \Delta G_{bind}^o &\approx \langle U_{P-L} \rangle_{PL} + \langle W_{PL} \rangle_{PL} - \langle W_P \rangle_P - \langle W_L \rangle_L \\ &\quad + \delta_{strn} [\langle U_P \rangle_{PL} + \langle U_L \rangle_{PL} - \langle U_P \rangle_P - \langle U_L \rangle_L - T\Delta S_{config}^o] \end{aligned} \quad (3.6)$$

where switching function δ_{strn} allows equation 3.6 to be exact for $\delta_{strn} = 1$, and approximately correct for $\delta_{strn} = 0$. Equation 3.6 may be recognized as equivalent to the MM-GBSA method, where the protein and ligand strain energies and the change in the configurational entropy are neglected when $\delta_{strn} = 0$, although various formulations have emerged in the literature.[60; 61; 62] Note, the $\delta_{strn} = 0$ approximation will be exactly satisfied by the model enclosure studied herein, but is expected to apply generally to any series of congeneric ligands binding to a given protein receptor. The reason we expect the $\delta_{strn} = 0$ approximation to be a reasonable approach to treating a series of congeneric ligands is that small modification of the ligand scaffold can be loosely understood to either make the scaffold slightly more or slightly less rigid, thereby changing the associated entropic cost of the protein binding the ligand. Those modification that make the ligand more rigid will lead to a less unfavorable binding entropy, but will also likely increase the

protein strain energy, since the protein must now deform to accommodate a more rigid object. Conversely, small modifications which increase the flexibility of the ligand will reduce the protein strain energy, since less deformation of the protein active site will be required upon binding the ligand, but will increase the entropic penalty of the binding process. It is this hypothesized general compensation of the strain energy with the loss of conformational entropy that should lead to the general applicability of the $\delta_{strn} = 0$ approximate form of Equation 3.6 to congeneric series.

The next series of approximations requires us to restrict our investigations to complementary ligands—ie, ligands that form hydrogen bonds with the protein receptor where appropriate, hydrophobic contacts otherwise, and sterically “fit” within the accessible volume of the active site of the receptor. Such ligands will form interactions with the surrounding protein similar to the interactions the ligand made with the bulk solvent—i.e hydrogen bonds where appropriate and van der Waals contacts otherwise, be they with the protein active site or with the solvating water. With this in mind, we may rewrite the solvation free energy terms as

$$\begin{aligned}\Delta \langle W_{PL} \rangle_{P,L;PL} &= \langle W_{PL} \rangle_{PL} - \langle W_P \rangle_P - \langle W_L \rangle_L \\ &= \Delta \langle W_{PL} \rangle_{P,L;PL}^{cav} + \Delta \langle W_{PL} \rangle_{P,L;PL}^{chrg}\end{aligned}\quad (3.7)$$

where $\Delta \langle W_{PL} \rangle_{P,L;PL}$ is the difference in the solvation free energy of the free ligand and protein versus the complex, $\Delta \langle W_{PL} \rangle_{P,L;PL}^{cav}$ is the free energy of growing the repulsive core of the ligand in the bulk versus within the protein active site, and $\Delta \langle W_{PL} \rangle_{P,L;PL}^{chrg}$ is the difference in the free energy of charging the ligand-solvent dispersion and electrostatic interactions in the bulk versus within the protein active site. Such a separation of the charging and cavitation terms is common in FEP studies of protein-ligand binding.[63; 64]

With the introduction of this notation, we find

$$\begin{aligned}\Delta G_{bind}^o &\approx \langle U_{P-L} \rangle_{PL} + \langle W_{PL} \rangle_{P,L;PL}^{cav} + \langle W_{PL} \rangle_{P,L;PL}^{chrg} \\ &+ \delta_{strn} [\langle U_P \rangle_{PL} + \langle U_L \rangle_{PL} - \langle U_P \rangle_P - \langle U_L \rangle_L - T\Delta S_{config}^o]\end{aligned}\quad (3.8)$$

We now introduce a rather aggressive approximation

$$\langle U_{P-L} \rangle_{PL} \approx -\Delta \langle W_{PL} \rangle_{P,L;PL}^{chrg} + \delta_{sie} \left[\langle U_{P-L} \rangle_{PL} + \Delta \langle W_{PL} \rangle_{P,L;PL}^{chrg} \right] \quad (3.9)$$

where an exact result is obtained for $\delta_{sie} = 1$, but an approximate result is generated for $\delta_{sie} = 0$. The rationale for this approximation can be explained as followed: $\Delta \langle W_{PL} \rangle_{P,L;PL}^{chrg}$ is the free energy difference in turning on the attractive and electronic interaction between the ligand and the solvent in bulk water versus in the active site of protein (see Figure 3.1), which is the interaction between the ligand and the solvent that would be excluded by the protein (depicted by dashed line in figure 3.1); $\langle U_{P-L} \rangle_{PL}$ is the interaction energy between the ligand and the protein in the complex (right). For complementary ligands binding to the protein receptor, the two terms would be expected to be similar in magnitude: (1) for polar ligands that make strong interactions with the protein receptor such as a salt bridge, the interaction of the ligands with water would also be strong; (2) for apolar ligands that make weak dispersion interactions with the protein, the interactions between the ligands and water would also be weak. We note the approximation described in equation 3.9 as “aggressive” in the sense that it would be expected to be generally false for an arbitrary ligand binding to an arbitrary receptor. Thus, by employing the approximation described by equation 3.9, we would only expect the following treatment to well describe ligands that satisfy the underlying assumptions, ie, the ligand form hydrogen bonds where appropriate and hydrophobic contacts otherwise. However, with the above caveat notes, we may approximate the binding free energy as

$$\begin{aligned} \Delta G_{bind}^o &\approx \langle W_{PL} \rangle_{P,L;PL}^{cav} + \delta_{sie} \left[\langle U_{P-L} \rangle_{PL} + \Delta \langle W_{PL} \rangle_{P,L;PL}^{chrg} \right] \\ &+ \delta_{strn} \left[\langle U_P \rangle_{PL} + \langle U_L \rangle_{PL} - \langle U_P \rangle_P - \langle U_L \rangle_L - T\Delta S_{config}^o \right] \end{aligned} \quad (3.10)$$

where our identified approximate equivalence between the relative protein-ligand direct interaction energy and the solvation-charging free energies has been explicitly noted in the grouping of the terms. Equation 3.10 suggests that the binding free energy may be approximated by computing the relative free energies of forming a cavity isosteric to the ligand in the protein active site, versus forming the same cavity in the bulk fluid.

Our remaining task is to develop a computationally efficient procedure to approximate the $\langle W_{PL} \rangle_{P,L;PL}^{cav}$ term. This term corresponds to the difference in the free energy of growing

the repulsive ligand cavity within the protein active site versus growing the ligand cavity in the bulk, or equivalently dragging the ligand cavity from the bulk through the volume of the system into the active site of the protein. The $\langle W_{PL} \rangle_{P,L;PL}^{cav}$ term may be exactly expanded as

$$\begin{aligned} \langle W_{PL} \rangle_{P,L;PL}^{cav} &= \left(G_{IST}^{PLcav} - G_{IST}^P \right) - \left(G_{IST}^{Lcav} - G_{IST}^{H_2O(l)} \right) \\ &= \Delta G_{IST}^{P,PLcav} - \Delta G_{IST}^{H_2O(l),Lcav} = \Delta \Delta G_{IST}^{Lcav} \end{aligned} \quad (3.11)$$

where G_{IST}^X is the inhomogenous solvation theory (IST) [4] integral over the system designated by superscript X, ie

$$\begin{aligned} G_{IST}^X &= E_{IST} - TS_{IST}^X \\ E_{IST}^X &= (E^K + E^{sw} + E^{ww})^X = \frac{3}{2}N_w kT + \rho \int g_{sw}^X(r) u_{sw}^X(r) dr + \frac{\rho^2}{2} \int g_{ww}^X(r_1, r_2) u_{ww}^X(r_1, r_2) dr_1 dr_2 \\ S_{IST}^X &= \left(S^{id} + S^{(1)} + S^{(2)} \dots \right)^X = \left[\frac{5}{2}N_w k - kN_2 \ln(\rho \kappa^3) \right] - k\rho \int g_{sw}^X(r) \ln g_{sw}^X(r) dr \\ &\quad - \frac{1}{2}k\rho^2 \int g_{sww}^X(r_1, r_2) [\ln \delta g_{sww}^X(r_1, r_2) - \delta g_{sww}^X(r_1, r_2) + 1] dr_1 dr_2 \dots \\ \delta g_{sww}^X(r_1, r_2) &= \frac{g_{sww}^X(r_1, r_2)}{g_{sw}^X(r_1)g_{sw}^X(r_2)} \end{aligned} \quad (3.12)$$

where g_{sw} , g_{ww} , and g_{sww} are the solute-water, water-water, and solute-water-water correlation functions; u_{sw} and u_{ww} are the solute-water and water-water interaction energy terms; r is the solvent degrees of freedom of system X; ρ is the density of the bulk fluid, and k is the Boltzmann constant.

Another simplification can be made by noting that the IST integrals appearing in equation 3.12 can be decomposed into two contributions: the contribution coming from the integral over the space of ligand cavity and the contribution coming from the integral over the rest of the space. So the ΔG_{IST} integrals appearing in equation 3.11 (be they in the bulk fluid or the protein active site) can also be decomposed into the corresponding two contributions: (1) the solvation free energies of N_w water molecules that were formerly solvating the protein active site and are evacuated into solution by the growth of the ligand cavity ($\Delta G_{IST, N_w, solv}$) (which comes from the integral over the ligand cavity part) (2) the

contribution from the solvent located at the L cavity surface ($\Delta G_{IST,surf}$) (which comes from the integral over the rest of the space) This decomposition of the total IST integrals into $\Delta G_{IST,surf}$ and $\Delta G_{IST,Nw,solv}$ terms may be clarified by inspecting the graphical depiction of the decomposition to be found in figure 3.2. It is also worth noting that in this notation $\Delta G_{IST}^{H_2O(l),Lcav} = \Delta G_{IST,surf}^{H_2O(l),Lcav}$ exactly, since the water is evacuated from a bulk environment to a bulk environment by the growth of the ligand cavity (ie, $\Delta G_{IST,Nw,solv}^{H_2O(l),Lcav} = 0$ strictly). Therefore,

$$\begin{aligned}
 \Delta\Delta G_{IST}^{Lcav} &= \left(\Delta G_{IST,surf}^{P,PLcav} + \Delta G_{IST,Nw,solv}^{P,PLcav} \right) - \Delta G_{IST,surf}^{H_2O(l),Lcav} \\
 &= \left(\Delta G_{IST,surf}^{P,PLcav} - \Delta G_{IST,surf}^{H_2O(l),Lcav} \right) - \Delta G_{IST,Nw,solv}^{P,PLcav} \\
 &= \Delta\Delta G_{IST,surf}^{Lcav} + \Delta G_{IST,Nw,solv}^{P,PLcav}
 \end{aligned} \tag{3.13}$$

where the “surf” term is the difference in the free energetic cost of the fluid reorganizing its configuration around the surface of the ligand cavity when the cavity is bound to the protein versus free in solution, and the “Nw, solv” term corresponds to the difference in the local IST integral free energy of the Nw water occupying the active site of the protein versus the IST integral free energy of the same Nw water molecules in the bulk fluid. Our final approximation is to assume that for small ligands that are expected to displace only one or a few water molecules deep within the protein active site, the “Nw solv” term should dominate this expression. Therefore, our final approximation to the binding free energy of the complex is

$$\begin{aligned}
 \Delta G_{bind}^o &\approx \Delta G_{IST,Nw,solv}^{P,PLcav} + \delta_{surf} \Delta\Delta G_{IST,surf}^{Lcav} + \delta_{sie} \left[\langle U_{P-L} \rangle_{PL} + \Delta \langle W_{PL} \rangle_{P,L;PL}^{chrg} \right] \\
 &+ \delta_{strn} \left[\langle U_P \rangle_{PL} + \langle U_L \rangle_{PL} - \langle U_P \rangle_P - \langle U_L \rangle_L - T\Delta S_{config}^o \right]
 \end{aligned} \tag{3.14}$$

where difference in the IST “surf” integrals are approximated as negligible when δ_{surf} is set to zero. Thus, our remaining task is to develop a numerical estimate the “Nw, solv” term.

Interestingly, a possible candidate estimator of $\Delta G_{IST,Nw,solv}^{P,PLcav}$ was previously introduced in reference [2], although its connection to the more rigorous expressions for computing protein-ligand binding affinities was not fully understood at the time of its introduction.

In the so called, displaced-solvent functional (DSF) approach, the local values of the IST integrals are computed for regions of high solvent occupancy in the active site, denoted by hydration sites. Note, that the volume of each hydration site is chosen such that the number of hydration sites will correspond to the N_w water molecules that are evacuated from the protein active site to the bulk fluid upon the binding of the ligand. This estimator itself was based on the following assumptions: (1) if atoms of a ligand overlapped with a hydration site, they displace the water from that site; and (2) the less energetically or entropically favorable the expelled solvent, the more favorable its contributions to the binding free energy. Thus, the relative binding free energy of the ligand is approximated as

$$\begin{aligned} \Delta G_{IST,N_w,solv}^{P,PLcav} &\approx \Delta G_{bind}^{DSF} = \sum_{lig,hs} (E_{bulk} - E_{hs}) \Theta(R_{co} - |r_{lig} - r_{hs}|) \\ &+ T \sum_{lig,hs} S_{hs}^e \left(1 - \frac{|r_{lig} - r_{hs}|}{R_{co}}\right) \Theta(R_{co} - |r_{lig} - r_{hs}|) \\ &= \sum_{lig,hs} \Delta G_{hs} \left(1 - \frac{|r_{lig} - r_{hs}|}{R_{co}}\right) \Theta(R_{co} - |r_{lig} - r_{hs}|) \end{aligned} \quad (3.15)$$

where ΔG_{bind}^{DSF} is the predicted binding free energy of the ligand, R_{co} is the distance cutoff for a ligand atom beginning to displace a hydration site, E_{hs} is the system-interaction energy of water in a given hydration site, S_{hs}^e was the excess entropy of water in a given hydration site, ΔG_{hs} is the computed free energy of transferring the solvent in a given hydration site from the active site to the bulk fluid, and Θ is the Heaviside step function. We also capped the contribution from each hydration site, such that it would never contribute more than ΔG_{hs} to ΔG_{bind}^{DSF} no matter how many ligand atoms were in close proximity to it. The value R_{co} might be considered a free parameter. However, an approximate value was adopted by noting that the radius of a carbon atom and a water oxygen atom are both approximately 1.4 Å, thus suggesting contact distances between a water oxygen atom and a ligand carbon atom less than $0.8 \cdot (1.4 \text{ Å} + 1.4 \text{ Å}) = 2.24 \text{ Å}$ are statistically improbable due to the stiffness of the Van der Waals potential. From the preceding approximate theory we infer that this approach should yield quantitatively accurate predictions of protein-ligand binding free energies versus the FEP reference data when the ligand is complementary to the protein active site and the reorganization entropies and energies of the protein and the

ligand are small compared to the other terms contributing to binding.

Here however, the preceding theory also suggests an alternative but related approach to adapting the DSF method to compute the binding free energy of a united atom methane molecule to a model hydrophobic enclosure. Here since the united atom methane molecule is itself simply a sphere that will occupy a known position in the binding site, we may simply collect statistics from the water molecules observed to occupy the volume that will be later occupied by the binding methane. Thus, clustering is unnecessary. From this data the energetic and entropic properties of the solvating water can be readily obtained via an application of inhomogeneous solvation theory. Lastly, it would in principle be possible to approximate the binding free energy of the methane molecule via the one evacuated-site-one-evacuated-water approximation introduced in reference [2]. However, we may also identify an approximate scaling that makes use of the known volume of the methane particle. In particular, if the methane particle is assumed to have a van der Waals radius of 1.865 Å, then the expectation value of the number of water molecules expected to exist within that volume is

$$N_{eff} = \rho_{bulk} \left(\frac{4}{3} \pi R_{methane}^3 \right) \approx 0.85 \quad (3.16)$$

where N_{eff} is the effective number of water molecules expected to be displaced by the bound methane assuming the entire system remains at bulk density, ρ_{bulk} is the density of liquid water, and $R_{methane}$ is the Van der Waals radius of the methane particle. Clearly, the number density of water in the active site depends on the environment of the specific enclosure, and in general would be different from bulk. However, the effective volume that is displaced by the binding methane is also different for different enclosures. Taking the situation of methane between two hydrophobic plates for example, considering the solvent-excluded volume consisting of the inward-facing surface of the probe ball with radius 1.4Å (size of water), in the bulk water the volume displaced by methane is just the van der Waals volume of methane, but the four corners are also excluded by the methane in between the two plates (see figure 3.3). It is well known that the number density of water in the hydrophobically enclosed region is smaller than bulk water because of dewetting. Thus the more enclosed the enclosures are, the smaller the number density of water in the active site, and the larger the effective volume displaced by the methane. These two competing

factors make the approximation introduced in equation 3.16 to be appropriate for all the enclosures. In principle, the exact number of excluded water molecules could be identified by the difference in the average number of water molecules surrounding the enclosure in the presence and absence of the bound methane, but this might require excellent statistics to converge.

To numerically test the validity of the preceding theory, we have constructed a series of model hydrophobic enclosures, as depicted in figure 3.4, and computed the binding free energy of a methane ligand for these hydrophobic enclosures both with FEP theory and the proposed DSF theory. The binding free energies of methane for the described enclosures, as computed by FEP, lie over a 5 kcal/mol range, which would correspond to 4 orders of magnitude of binding affinity. Thus, the ability to accurately predict such free energy differences would be expected to have great utility in a drug-design setting.

A final important point, not relevant to the present model systems but relevant when considering realistic problems such as protein-ligand binding, is the necessity in such real problems for integrating over the solute coordinates. For example, fluctuations of the protein-ligand complex at room temperature can be significant, and in principle this affects the water structure in the active site. In our DSF approach to date, we have employed a single “representative” structure for the protein structure (by harmonically restraining the coordinates to a target structure during the DSF molecular dynamics simulation) rather than allowing the solute phase space to be fully explored. For the model hydrophobic enclosures, there is no issue with averaging over solute configurations because the model enclosures are specified as rigid from the beginning.

In the context of our DSF methodology, the interesting question is how good an approximation the harmonically restrained simulation is to the fully fluctuating solute when estimating the free energy changes resulting from solvent displacement by the ligand. A heuristic argument that the approximation is reasonable if it is assumed that, for relatively modest fluctuations of the complex (as opposed to major conformational changes), the solvation in the active site “follows” the solute atoms in essence an adiabatic approximation in which the solvation structure readjusts quickly to typical excursions of solute atoms from the central configuration. If this is in fact the case, then the free energy of displacement of

a given water molecule at all accessible solute configurations can be approximated by the displacement free energy at the central configuration. This is not a rigorous or controlled approximation, but it appears to work reasonably well based on a range of examples that we have investigated to date. We do not consider this point further in the present paper, as our focus is on a series of rigid solutes; however, in future work, explicit investigation of this hypothesis, based on computing DSFs for different solute configurations and comparing them, will be pursued.

3.2.2 Simulation details

3.2.2.1 DSF analysis

To generate the data required to apply the DSF method of computing protein-ligand binding free energies to the model hydrophobic enclosures, each of the thirteen hydrophobic enclosures depicted in figure 3.4 were subjected to explicitly solvated molecular dynamics with the Desmond molecular dynamics program.[31] The Maestro System Builder[65] utility was used to insert each enclosure into a cubic water box with a 10 Å buffer. The SPC water model[38] was used to describe the solvent, and the united atom methane molecules that formed the “atoms” of the enclosures were uniformly represented with $\sigma = 3.73\text{\AA}$ and $\epsilon = 0.294$ kcal/mol Lennard Jones parameters. The atoms of the enclosures were constrained to their initial positions throughout their dynamics, and only the solvent degrees of freedom were sampled. The energy of the system was minimized, and then equilibrated to 298 K and 1 atm with Nose-Hoover[33; 34] temperature and Martyna-Tobias-Klein[35] pressure controls over 500 ps of molecular dynamics. A cutoff distance of 9 Å was used to model the Lennard Jones interactions, and the particle-mesh Ewald method[37] was used to model the electrostatic interactions. Following the equilibration, a 20 ns production molecular dynamics simulation was used to obtain statistics of the water solvating the enclosures, and configurations of the system were collected every 1.002 ps.

Following the previously developed approach,[2; 3] the position the ligand would occupy in the enclosures was used to define the active site volume. Here, a 1 Å cutoff distance from the center of where the ligand center would be was used to define the solvent volume of interest. A water molecule was identified to be in the active site when its oxygen lay

within the sphere, and otherwise not. For each solvent molecule identified in this volume, we computed the system-interaction energy of the solvent molecule (ie, the interaction energy of the solvent molecule with the rest of the system), and recorded its orientation and position. From this data, we computed the average system-interaction energy of solvent occupying this volume, and the excess entropy of this solvent from an expansion of the entropy in terms of translational and orientational correlation functions.

The calculation of excess entropies of water in the hydration sites was processed in a two-step manner: (1) introduce an intermediate reference state with the same average number density as the hydration site we are studying but a flat translational and orientational distribution, and calculate the excess entropy of the hydration site water with respect to this intermediate reference state due to the local ordering of water in the hydration site (2) determine the entropy difference between the intermediate reference state and the bulk water that is due to the difference of number density. The entropy difference between water in the hydration site and the intermediate state was calculated through the integral introduced in equation 3.12, with $g_{sw}(r)$ defined with respect to the intermediate reference state number density. In order to integrate this entropy expansion, we adopted a k-th nearest neighbors approach as introduced in reference.[66]

To characterize the orientation of waters in the hydration site, we built the coordinate system such that the center of the hydration site was taken to be the origin, the z axis was perpendicular to the plate (take enclosure F, for example), and a second methane not lying on the z axis was arbitrarily chosen to define the direction of the x axis. The orientation of water in the hydration site was defined by six variables, $[r, \theta, \phi, \chi_\theta, \chi_\phi, \chi_\sigma]$, where $[r, \theta, \phi]$ are the typical spherical coordinates which define the position of the oxygen atom, and $[\chi_\theta, \chi_\phi, \chi_\sigma]$ are the three angles which define the orientation of the water around its oxygen (see figure 3.5). To clarify, $[\chi_\theta, \chi_\phi]$ are similar to the typical spherical coordinate angles $[\theta, \phi]$ which define the orientation of the dipole vector of water, and χ_σ defines the rotation of hydrogen around the dipole vector. For enclosures with rotational symmetry about the z axis, the distribution along ϕ angle is flat by symmetry, so we only need five angles to define the orientation of water. The calculation of the entropy difference is performed through the

following equation:

$$S_1 = -k \frac{1}{V\Omega} \int J(r, \theta, \phi, \chi_\theta, \chi_\phi, \chi_\sigma) g(r, \theta, \phi, \chi_\theta, \chi_\phi, \chi_\sigma) \ln g(r, \theta, \phi, \chi_\theta, \chi_\phi, \chi_\sigma) dr d\theta d\phi d\chi_\theta d\chi_\phi d\chi_\sigma \quad (3.17)$$

where $g(r, \theta, \phi, \chi_\theta, \chi_\phi, \chi_\sigma)$ is the solute water pair correlation function (PCF), and $J(r, \theta, \phi, \chi_\theta, \chi_\phi, \chi_\sigma)$ is the Jacobian associated with these variables. Here $g(r, \theta, \phi, \chi_\theta, \chi_\phi, \chi_\sigma)$ has the property that

$$\frac{1}{V\Omega} \int J(r, \theta, \phi, \chi_\theta, \chi_\phi, \chi_\sigma) g(r, \theta, \phi, \chi_\theta, \chi_\phi, \chi_\sigma) dr d\theta d\phi d\chi_\theta d\chi_\phi d\chi_\sigma = 1 \quad (3.18)$$

where V is the volume of the sphere and Ω is the total angular volume over angular variables $[\chi_\theta, \chi_\phi, \chi_\sigma]$, ie

$$\Omega = \int J(\chi_\theta, \chi_\phi, \chi_\sigma) g(\chi_\theta, \chi_\phi, \chi_\sigma) d\chi_\theta d\chi_\phi d\chi_\sigma \quad (3.19)$$

In line with reference[66] (Chapter 1) we approximate the total pair correlation function (PCF) through generalized Kirkwood superposition approximation (GKSA),[21] which allowed the entropy to be approximated by the summation and subtraction of one- and two-dimensional entropies, and calculated the one- and two-dimensional entropies through NN method.

The entropy difference between the reference state and bulk water can be simply calculated by recognizing the entropy expression for homogeneous ideal-gas:

$$S_{id} = \frac{3}{2} - k \ln(\rho\Lambda^3) \quad (3.20)$$

where Λ is the thermal wavelength. So the excess entropy of the second step is simply:

$$S_2 = -k \ln \left(\frac{\rho_{ref}}{\rho_{bulk}} \right) \quad (3.21)$$

where ρ_{ref} , ρ_{bulk} are the number density of the reference state and bulk water respectively.

The total excess entropy is the sum of S_1 and S_2 as defined by equation 3.17 and 3.21.

3.2.2.2 FEP analysis

The dynamics simulation used to perform the FEP analysis of the binding free energy of the methane ligand to the model hydrophobic enclosures were run under identical simulation

protocols as the DSF analysis. The ligand was turned on inside the model enclosures over 9 lambda windows with $\lambda = [0, 0.125, 0.25, 0.375, 0.50, 0.625, 0.75, 0.875, 1]$, where λ is the coupling parameter to turn on/off the interaction between the methane and the rest of the system with initial state and final state correspond to $\lambda = 0$ and $\lambda = 1$ respectively. At different λ windows, we performed molecular dynamics simulations, and calculated the energy difference between neighboring λ values for each configuration saved. In these simulations, the soft-core interactions were used for the Lenard-Jones potential.[67] Bennett acceptance ratio method[30] were then used to calculate the free energy difference between neighboring states. The sum of the free energy differences between neighboring states gave the solvation free energy of methane in question. The same procedure was followed to calculate the solvation free energy of methane in bulk water. The difference between the two solvation free energy gave the binding free energy to bring a methane from infinitely far to inside the hydrophobic enclosure. (We can also interpret the binding free energy as the potential of mean force between the methane and the enclosure.)

3.2.2.3 Buried surface area analysis

The solvent accessible surface area (SASA) and molecular surface area (MSA, or Connolly surface) of each enclosure with and without the bound methane was computed with the Connolly molecular surface package,[68] as was the SASA and MSA of the methane particle by itself. From this data the buried solvent accessible surface area upon methane-enclosure complexation was determined. The Lennard Jones interaction energy of the methane particle with the model enclosure was similarly computed. The buried surface area times the surface tension would give the solvent induced interaction energy, and together with the direct Lennard-Jones interaction energy, the total binding energy of methane with different enclosures can be calculated, as routinely estimated in various empirical methods to estimate the contribution of the nonpolar term to the binding energy.

3.3 Results

The binding free energies of methane for the model hydrophobic enclosures, as measured by FEP, are reported in table 3.1. It is found that the range of binding free energies of the methane ligand for the model enclosures is nearly 5 kcal/mol. Also reported in table 3.1 are the system-interaction energies and excess entropies of the water displaced by the methane ligand, the buried surface area upon complexation, (both SASA and MSA), the change of the Lennard Jones interaction energy between the methane particle and the enclosure upon complexation, the DSF prediction of the binding free energy of the complex, and the scaled DSF prediction that makes use of the scaling coefficient deduced from first principles in section 3.2. The R^2 value, mean-absolute-error (MAE) and the root-mean-square-error (RMSE) between the various predictions with the FEP-reference data are also listed in the last few rows of the table. Note here that the surface tension coefficients for the buried surface area/molecular mechanics predictions (Both SASA and MSA) were explicitly tuned to minimize the MAE of the predictions. Such explicit tuning yields significantly better results than could reasonably be expected to be obtained if such methods were employed with fixed coefficients across realistically variable data sets.

The DSF predictions show very high correlation with the FEP reference data, as indicated by the R^2 value of 0.95, (which can also be seen in figure 3.6) where the buried surface area/Lennard Jones interaction predictions show reduced correlations, as indicated by R^2 values of 0.92 for MSA/MM and 0.76 for SASA/MM respectively. The DSF method also allows for the decomposition of the binding free energy prediction into separate enthalpic and entropic components. Inspection of the data reported in table 3.1 indicates that the DSF predictions are dominated by the enthalpic contribution to the binding affinity, which by itself manifests a R^2 value of 0.94 versus the FEP reference data. Detailed analysis of these data indicates that, except for the first three systems, the binding of the methane molecule to these hydrophobic enclosures is mainly an enthalpy driven event, which is consistent with our knowledge about large length scale hydrophobicity.[1; 69; 70] Recent calorimetry data obtained for Major Mouse Urinary Protein by Homans et al,[71] appear to indicate such enthalpy driven hydrophobic binding events are witnessed in vivo, as well.

The inspection of the trajectory indicates the atomistic basis of the enthalpy driven effect is that water molecules that solvate such enclosures are forced to break hydrogen bonds. The effect is most obvious for hydrophobic enclosures L and M, where the solvent suffers a 7 kcal/mol reduction in system-interaction energy when occupying these enclosures, while almost no reduction in excess entropy versus bulk water. Conversely, the methane dimerization free energy described by methane binding to “enclosure” A is dominated by the entropic contribution, again consistent with entropy driven small length scale hydrophobic effect. This finding is analogous to the well characterized length scale dependence of the hydrophobic effect, while small hydrophobes are found to induce entropic ordering of the solvent, large hydrophobes are found to break water-water hydrogen bonds.[1; 69; 70] The enclosures L and M can thus be understood as manifesting extreme large-length scale hydrophobic character from the perspective of the solvating water.

Figure 3.6 plots the correlation of the DSF binding free energies versus the FEP reference data with and without the derived scaling coefficient deduced from the size of the methane ligand itself. As can be seen from the figure, both sets of predictions track the FEP reference data quite well. However, the scaled predictions have greater quantitative agreement with the FEP, which may be quantified by the mean-absolute error (MAE) and root-mean-square error (RMSE) metrics. Here the scaled predictions are found to have a MAE of 0.36 kcal/mol and a RMSE of 0.40 kcal/mol, while the unscaled predictions have a MAE of 0.66 kcal/mol and a RMSE of 0.84 kcal/mol. Thus, the deduced scaling coefficient appears to increase the quantitative accuracy of the approach, in line with the expectation of the theoretical analysis.

We also investigated to what extent a combined buried surface area/Lennard-Jones interaction energy model might be able to reproduce the binding affinities. Tuning the model to minimize the MAE of the fit, we obtained an optimal surface tension coefficient of $\gamma = 0.011 \text{ kcal/mol} * \text{\AA}^2$ for SASA and $0.044 \text{ kcal/mol} * \text{\AA}^2$ for MSA for these enclosures, which is somewhat smaller than the reported literature values.[72] These predictions versus the FEP reference data are reported in figure 3.7. It is found that MSA/MM performed much better compared with SASA/MM, which is indicated by much higher R^2 value, and smaller MAE and RMSE values. (Data listed in last 3 rows in table 3.1.) However, both of

them performed less well than the DSF predictions with the scaling coefficient correction, and much worse results would be expected with such an model in general, as noted above, since it would not benefit from explicit fitting to the reference data.

The better performance of MSA/MM versus SASA/MM is due to the better characterization of MSA for the topology of enclosures J, K, L, M. SASA/MM predicts enclosure J to be most hydrophobic, which corresponds to a methane molecule binding between two hydrophobic plates, because large swaths of formerly SASA on the faces of the plates are buried by the presence of the methane ligand for enclosure J, while for enclosures K, L, and M several methane molecules already lie between the plates in the absence of the binding ligand and thus some of the surface area that would be buried by the binding methane is already buried by the other particles. However, MSA can better characterize the curvature of these enclosures and predict the right order of binding affinity.

3.4 Conclusion

Calculations suggest that the DSF method of computing protein-ligand binding affinities may offer near-FEP accuracy at a substantially reduced computational expense for systems that satisfy the requisite approximations and should offer greater quantitative accuracy than competing implicit solvent methodologies. Further, the clear connection between the DSF method and more rigorous statistical mechanical expressions may offer a rational path to systematically improve the accuracy and rigor of the method by progressive inclusion of those counter-balancing terms currently approximated to exactly cancel. This previously opaque connection to the underlying theory facilitated the derivation of a scaling coefficient that was seen to increase the quality of the predictions of the method versus the FEP reference data. Lastly, the molecular detail afforded by the technique may offer insight into protein-ligand binding processes, such as highlighting the importance of the enthalpy in the binding of methane to such model enclosures, which may have been difficult to discern from only FEP or implicit modeling.

Table 3.1: The binding thermodynamics of methane for the various model hydrophobic enclosures as computed from DSF theory and FEP theory. E_{hs} was the hydration site system interaction energy, S_{hs}^e was the hydration site solute-water correlation entropy, $\Delta SASA$ was the buried solvent accessible surface area using a 1.4 Å radius probe, ΔE_{LJ} was the Lennard Jones interaction energy of the bound methane with the rest of the enclosure, ΔG_{bind}^{DSF} was the predicted binding free energy of the methane molecule for the model enclosure as computed from DSF theory, $N(N_{eff})$ was scaling coefficient derived by determining the expectation value of the number of water molecules occupying a volume in the bulk fluid equal to the volume of the methane probe molecule, and ΔG_{bind}^{FEP} was the predicted binding free energy of the methane molecule for the model enclosure as computed from FEP theory. Note that the standard deviation of the E_{hs} values reported below were found to be uniformly less than 0.4 kcal/mol (as obtained from block averaging), and the standard errors of the ΔG_{bind}^{FEP} values were uniformly less than 0.02 kcal/mol.

Model En- closure	E_{hs} (kcal/mol)	S_{hs}^e (kcal/mol*K)	$\Delta SASA$ (Å ²)	ΔMSA (Å ²)	E_{LJ} (kcal/mol)	ΔG_{bind}^{DSF} (kcal/mol)	$N \cdot \Delta G_{bind}^{DSF}$ (kcal/mol)	ΔG_{bind}^{FEP} (kcal/mol)
bulk	-19.8	0	0	0	0	0	0	0
A	-19.6	-1.2	-59.45	-3.84	0	-0.5	-0.46	-0.61
B	-18.9	-2.0	-118.9	-7.67	0	-1.5	-1.28	-1.15
C	-19.2	-1.8	-98.21	-10.49	0	-1.1	-0.97	-1.41
D	-18.7	-1.2	-91.32	-13.51	-1.41	-1.5	-1.26	-1.66
E	-17.7	-2.3	-151.15	-17.35	-1.41	-2.8	-2.39	-2.17
F	-17.3	-1.5	-117.52	-24.06	-1.41	-2.9	-2.5	-2.63
G	-16.0	-3.0	-156.39	-30.7	-1.41	-4.7	-4.00	-3.41
H	-15.6	-1.2	-132.41	-37.35	-1.41	-4.6	-3.92	-3.43
I	-15.6	-1.8	-143.71	-34.6	-1.41	-4.8	-4.05	-3.47
J	-17.8	-2.6	-182.65	-27.02	-2.82	-2.8	-2.41	-2.86
K	-15.5	-2.1	-175.59	-44.27	-2.82	-4.9	-4.17	-4.59
L	-13.0	0.3	-166.61	-64.21	-2.82	-6.8	-5.74	-5.24
M	-13.3	-0.1	-168.52	-61.51	-2.82	-6.6	-5.6	-5.45
R^2	0.94	0.16	0.76(a)	0.92(b)	0.73	0.95	0.95	N/A
MAE	0.61	N/A	0.54(a)	0.47(b)	1.41	0.66	0.36	N/A
RMSE	0.75	N/A	0.74(a)	0.58(b)	1.63	0.85	0.40	N/A

Note: (a): these values correspond to the correlation between the buried SASA/LJ interaction with optimized surfacetension coefficient ($\gamma = 0.044kcal/mol * \text{Å}^2$) and the FEP reference data.

(b):these values correspond to the correlation between the buried MSA/LJ interaction with optimized surface tension coefficient ($\gamma = 0.011kcal/mol * \text{Å}^2$) and the FEP reference data.

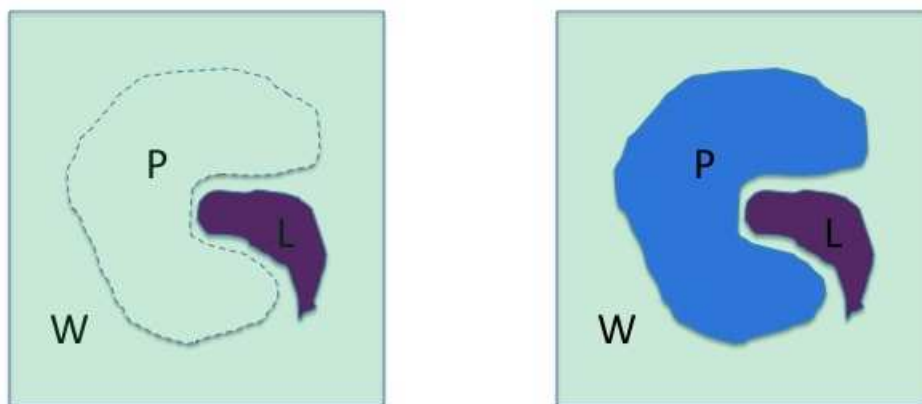


Figure 3.1: Cartoon depicting the relationship between $\Delta \langle W_{PL} \rangle_{P,L;PL}^{chg}$ and $\langle U_{P-L} \rangle_{PL}$. $\Delta \langle W_{PL} \rangle_{P,L;PL}^{chg}$ is the free energy difference in turning on the attractive and electronic interaction between the ligand and the solvent in the bulk water (left) versus in the active site of protein (right), which is the interaction between the ligand and the solvent that would be excluded by the protein (depicted by dashed line on the left). $\langle U_{P-L} \rangle_{PL}$ is the interaction energy between the ligand and the protein in the complex (right). For complementary ligands binding to the protein receptor, the two terms would be expected to be of similar magnitude.

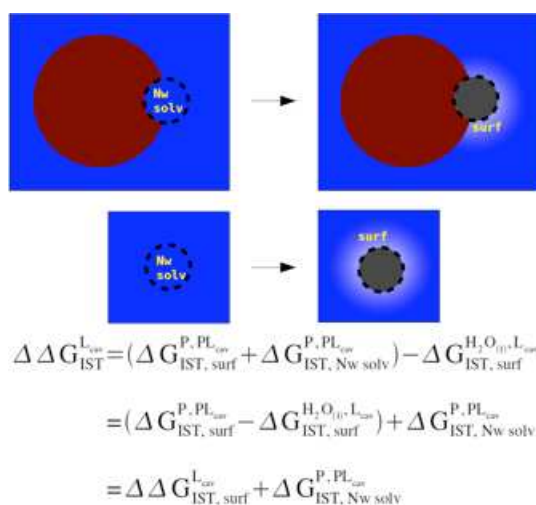


Figure 3.2: Cartoon depicting the spatial decomposition of the IST integral equations introduced in equations 3.11 to 3.14. The net “surf” term is the difference in the free energetic cost of the fluid reorganizing its configuration around the surface of the ligand cavity when the cavity is bound to the protein versus free in solution, and the net “Nw solv” term corresponds to the difference in the local IST integral free energy of the Nw water occupying the active site of the protein versus the IST integral free energy of the same Nw water molecules in the bulk fluid.

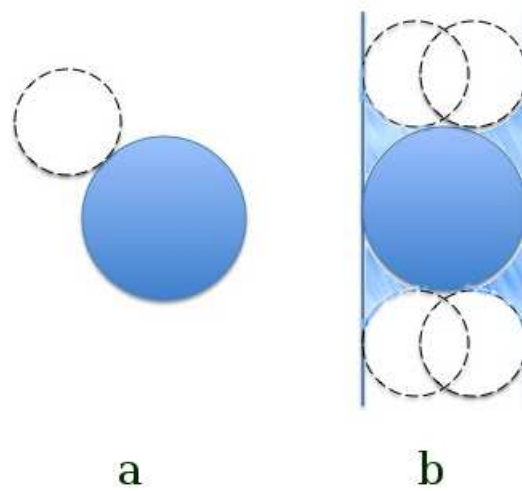


Figure 3.3: The effective volume displaced by a methane in the bulk (a) and in between two hydrophobic plates(b). The blue particle denotes a methane, and a dashed circle denotes a probe solvent molecular. The volume displaced by a methane in the bulk is just the van der Waals volume of the methane, but in between the two plates, the four corners are also displaced by the methane due to the finite volume of the probe ball.

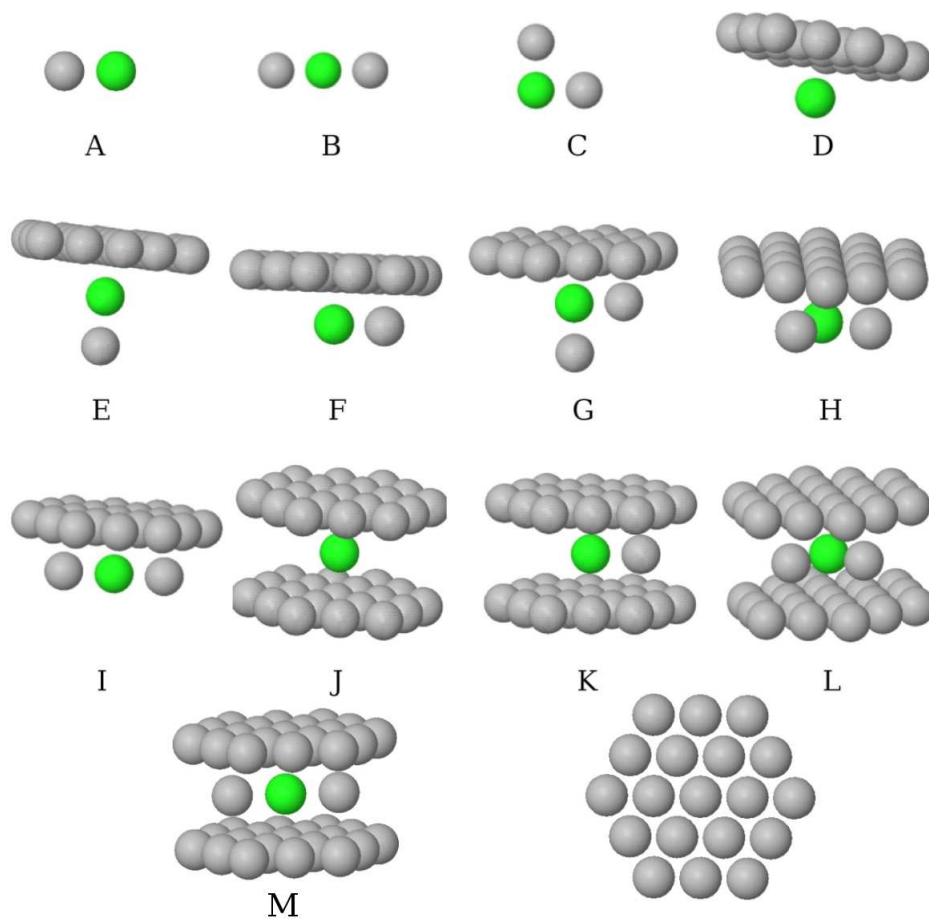


Figure 3.4: The 13 model hydrophobic enclosures are here depicted in gray. The location of the methane molecule when bound to the respective hydrophobic enclosures is here depicted in green. The geometry of the plate is depicted at the right bottom of this figure. The distance between the neighboring particles in the plate is 3.2 \AA , and the distance between the two plates is 7.46 \AA . All the others particles are at contact distance with linear (B, I and M) and triangle (C, G, H and L) geometries.

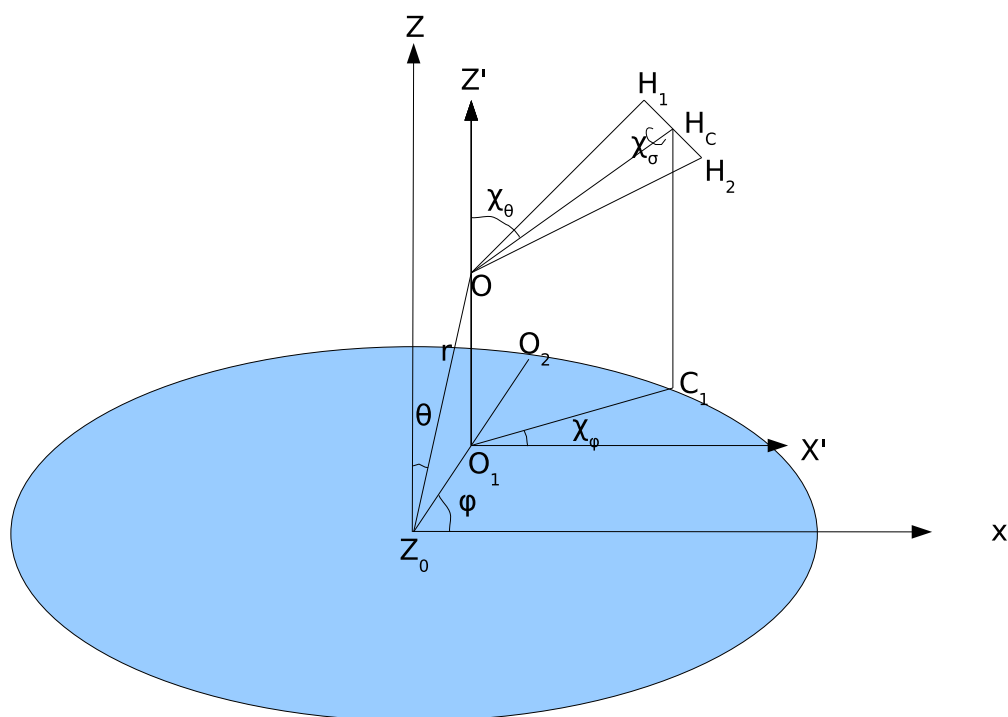


Figure 3.5: The coordinate system to characterize the position and orientation of water inside the hydration site. The z axis is perpendicular to the model hydrophobic plate, and the x axis is such defined that the other methane lie on the x axis. $[r, \theta, \phi]$ are the typical spherical coordinates which define the position of the oxygen atom, and $[\chi_\theta, \chi_\phi, \chi_\sigma]$ are three angles which define the orientation of the water around its oxygen.

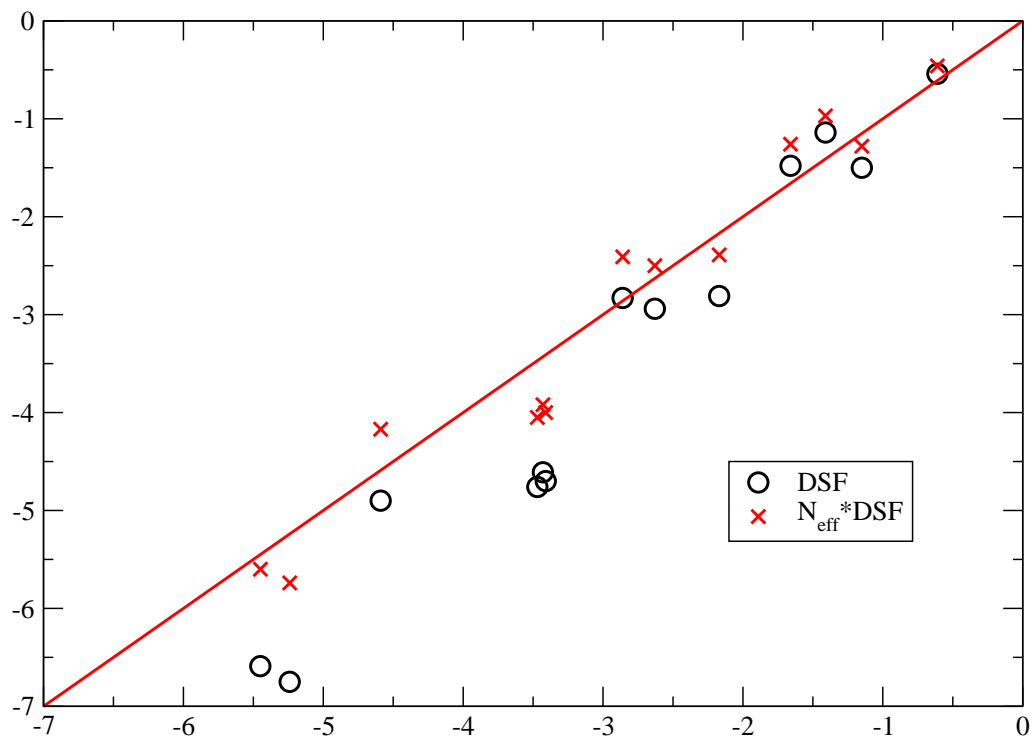


Figure 3.6: The correlation of the of the DSF predictions of the methane-enclosure binding free energies with the FEP reference data.

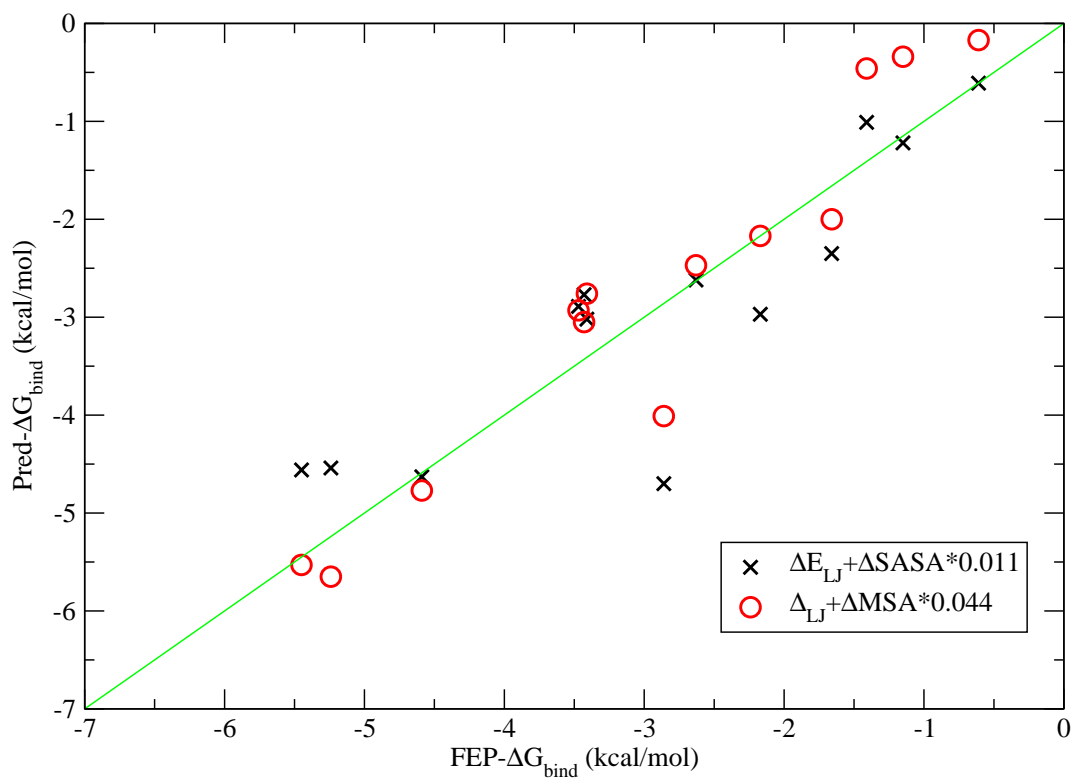


Figure 3.7: The correlation of buried surface area/molecular mechanics predictions of the methane-enclosure binding free energies with the FEP reference data. The water SASA surface tension coefficient ($0.011\text{kcal/mol} * \text{\AA}^2$) and MSA surface tension coefficient ($0.044\text{kcal/mol} * \text{\AA}^2$) were tuned to minimize the absolute average error of the predictions with respect to the reference data.

Chapter 4

Protein-Ligand binding: Contributions from wet and dry regions of the binding pocket

Abstract

Biological processes often depend on protein-ligand binding events, yet accurate calculation of the associated energetics remains as a significant challenge of central importance to structure-based drug design. Recently, we have proposed that the displacement of unfavorable waters by the ligand, replacing them with groups complementary to the protein surface, is the principal driving force for protein-ligand binding, and we have introduced the WaterMap method to account this effect. However, in spite of the adage “Nature abhors vacuum”, one can occasionally observe situations in which a portion of the receptor active site is so unfavorable for water molecules that a void is formed there. In this Chapter, we demonstrate that the presence of dry regions in the receptor has a nontrivial effect on ligand binding affinity, and suggest that such regions may represent a general motif for molecular recognition between the dry region in the receptor and the hydrophobic groups in the ligands. With the introduction of a term attributable to the occupation of the dry regions by ligand atoms, combined with the WaterMap calculation, we obtain excellent agreement

with experiment for the prediction of relative binding affinities for a number of congeneric ligand series binding to the MUP receptor. In addition, WaterMap when combined with the cavity contribution is more predictive than at least one specific implementation (described in ref. [2]) of the popular MM-GBSA approach to binding affinity calculation.

4.1 Introduction

The calculation of protein-ligand binding affinities is a central goal of computational structure based drug design methodologies. Many different approaches, ranging from rapid empirical scoring functions to rigorous free energy perturbation methods, have been employed.[48; 49; 51] At present, however, there is no method that is fully satisfactory from the point of view of both the expected accuracy and reliability, and the required computing resources.

In the first two chapters, we have introduced a new approach to estimating relative free energies of binding of a series of congeneric ligands, based on their measured displacement of quasi-localized water molecules with unfavorable free energies in the receptor active site.[3; 2] We refer to this approach as WaterMap. Molecular dynamics simulations are used to generate the positions of the relevant water sites, and inhomogeneous solvation theory is employed to estimate free energies of displacement of the various waters as compared to bulk solvent. Successful prediction of the relative binding free energies of a set of congeneric pairs of Factor Xa ligands, without the use of any adjustable parameters, was achieved, with a correlation coefficient considerably superior to an widely used alternative, the MM-GBSA approach which employed a continuum description of solvent.[2; 73] A number of other applications have recently appeared, all of which yield encouraging results with regard to the efficacy of relative ligand binding affinity predictions.[53; 54; 56]

Displacement of unfavorable waters by the ligand, replacing them with groups complementary to the protein surface, has been established as a principal driving force for protein-ligand binding in many systems, including a significant fraction of receptors of pharmaceutical interest.[74] However, one can also occasionally observe situations in which a portion of the receptor active site is so unfavorable for water molecules that a void is formed, i.e. in the molecular dynamics runs which generated the WaterMap, regions could be identified

where occupancy of water molecules was observed to be below a specified threshold. A number of proteins exhibiting a dry region in the binding pocket were discussed in ref. [75]

The presence of dry regions would be expected to have a nontrivial effect on ligand binding affinity, if the ligand places atoms in these regions (as would be highly favorable in terms of free energy if the ligand groups are complementary to the protein surface in the appropriate region). In the present paper, we investigate this issue quantitatively by obtaining from the literature a number of ligand series for ligands which bind to several proteins with dry regions, and developing a methodology to combine the WaterMap free energy difference with an additional term attributable to occupation by ligand atoms of the dry regions. Using a very simple model with essentially no adjustable parameters, excellent agreement with experiment is obtained, as compared to results derived from a WaterMap-only calculations, which fails to yield a plausible correlation of the theoretical predictions with experiment. The term for the dry region is straightforward to implement, and we expect to employ it routinely in future studies of binding affinity using this general type of approach.

In what follows, we describe the new methodology, and compare results for a number of ligand series for the combined method and WaterMap alone. In the Conclusion, we summarize our results and suggest future research directions.

4.2 Results and Discussions

We analyzed the hydration properties of the unliganded binding pockets for several holo-proteins, including the mouse major urinary protein (MUP, PDB ID 1znk),[76] the bovine apo-glycolipid transfer protein (GLTP, PDB ID 1wbe),[77] and the secretin pilot protein (PDB ID 1y9l),[78] and identified both the high occupancy hydration sites using the WaterMap program[3; 2] and the low occupancy cavity regions using the protocol described in the methods section. Fig. 4.1 displays the high occupancy hydration sites and the dry regions in the active site of MUP. As opposed to most proteins with well hydrated active sites, the active site of MUP is poorly hydrated, as indicated by a large dry region and only two active site water molecules, which is consistent with previous discussions.[79; 71;

80]

There are several ligands which bind to MUP.[76; 81; 82; 83] As indicated by X-ray diffraction data, MUP is rather rigid, and the structure remains essentially unchanged upon binding to these different ligands.[76] By superposition of each protein-ligand complex to the “apo” structure of the protein and accounting for the contribution to the binding affinity through displacing the active site solvent, which is the standard protocol of the WaterMap calculation, we get the WaterMap predicted binding affinity for each ligand. Fig. 4.2 plots the WaterMap predicted binding affinities versus the experimental results (circles in Fig. 4.2) for the ligands with experimental binding affinity data available from literature. The ligands are divided into four groups (indicated by four different colors in figure 4.2) based on their structure similarity and binding mode. The ligands in each group share the same scaffold and binding mode based on their PDB structures, and their experimental binding affinity data are from the same publication, and derived using the same method. (For the 2-sec-butyl-4,5-dihydrothiazole (SBT) series of ligands, PDB structure is only available for SBT-MUP complex; all the other structures in that group were obtained by removing the appropriate carbon atoms from ligand SBT.[83]) We see from Fig. 4.2 that, while WaterMap can explain the binding affinity difference between ligand PE9 and ligand HE2 (blue circles in Fig. 4.2), it can not explain the binding affinity differences among the other groups of ligands (red, green, and black circles in Fig. 4.2). To be specific, WaterMap predicts ligand HE4 to have zero binding affinity (because the ligand displaces none of active site solvents), which is much lower than the other two ligands OC9 and F09 in that group, while experimentally their binding affinity difference is much smaller. In addition, WaterMap predicts that the binding affinities for ligands OC9 and F09 are the same, while experimentally ligand F09 is 3.2 kJ/mol more favorable than ligand OC9. Similar deficiencies are observed for ligands IBMP and IPMP, as well as for all the ligands in the SBT series.

While the WaterMap calculation takes into consideration the binding affinity gain from ligand atoms displacing the energetically and entropically unfavorable hydration sites, the ligand atoms located in the dry region are not scored. It is well known that the solvation free energy of the ligand has two contributions: the free energy to create the cavity via

displacement of solvent, and the free energy to turn on the interactions between the ligand and the rest of the system.[84] While it engenders a large free energy penalty to create a cavity in bulk water in order to solvate the ligand, the free energy to create the cavity in the active site of the protein is almost zero if it is dry there. So the ligand gains much binding affinity if it is located in the dry region of active site, which we call the cavity contribution. We use the scoring function described in the methods section to take this effect into consideration. The physical basis of the method is that the free energy difference of “growing” one ligand heavy atom inside the active site of the protein versus that in bulk water is the gain in binding affinity from that atom.

Adding together the WaterMap contribution and the cavity contribution described above for each ligand, the overall predicted binding affinities versus experimental results are displayed in Fig. 4.2 (crosses in Fig. 4.2). It is quite obvious that after taking the cavity contribution into consideration, the binding affinity differences among different ligands in each group (indicated with different colors) are correctly predicted. For comparison, the MM-GBSA predictions for the binding affinities of these ligands were also calculated, and the WaterMap combined with cavity predictions works much better than MM-GBSA predictions for all four congeneric groups (see Fig 4.3). (The WaterMap and cavity contribution to the binding affinities and the MM-GBSA predictions are given in Table4.1) If we fit the predicted results against the experimental data among each group with a line, the slopes of the lines for the four groups are of similar magnitude, but the intercepts are different. This behavior is expected. The different intercepts among the groups indicate the different strain and conformational energy and entropy changes upon protein-ligand complexation for different ligand scaffolds, which are not taken into account in this analysis and which is also part of the reason the predicted binding affinities much larger in magnitude than the experimental ones. The fact that we only take into account the favorable effects in binding either from water displacement or from favorable ligand-cavity interaction, but not the unfavorable effects such as loss of conformational entropy and part of the desolvation penalty also makes the predicted binding affinities much larger than experimental results. However, the ability of the current analysis method in rank-ordering a series of congeneric ligands makes it useful and important in lead optimization. This is clearly demonstrated in Fig.

4.4 where the predictions of the relative binding affinities among congeneric ligand pairs for the three methods versus experimental data are plotted, and the WaterMap combined with cavity predictions work much better than WaterMap alone and MM-GBSA method.

As an example of how the WaterMap and cavity contributions complement each other to rank-order a pair of congeneric ligands, Fig. 4.5 displays the structures of ligand HE4 (colored green) and ligand OC9 (colored blue) in the binding pocket of MUP. While ligand OC9 displaces one of the two principal hydration waters (red spheres in Fig. 4.5), ligand HE4 does not have any overlap with the two hydration waters. This is consistent with experimental results that one more ordered water molecule is present within the binding pocket of HE4-MUP complex.[76] And this is the reason why the WaterMap calculation predicts zero binding affinity for ligand HE4 and -44.3KJ/mol for ligand OC9, while experimentally ligand OC9 is only 3.1 KJ/mol more favorable than ligand HE4. However, most of the atoms of ligand HE4 are located in the dry region (white networks in Fig. 4.5), which leads to a more favorable cavity contribution to the binding affinity for ligand HE4 than for ligand OC9 (-78.6KJ/mol for HE4 versus -49.1KJ/mol for OC9). So the overall binding affinity difference predicted agrees well with experimental data.

Fig. 4.6 (a) displays the structures of ligand OC9 (colored blue) and ligand F09 (colored green) in the binding pocket of MUP. Both ligands have similar structure in the hydration water part of the pocket, so the WaterMap calculation predicts their binding affinities to be the same. However, experimentally ligand F09 is 3.2 KJ/mol more favorable than ligand OC9.[76] Looking at their structures in the dry region, it is quite clear that ligand F09 has one more atom located in the dry region, which leads to the more favorable binding of ligand F09 than ligand OC9. Similar behavior is observed for ligand IBMP and ligand IPMP (Fig. 4.6 (b)): one more atom of ligand IBMP in the dry region leads to the more favorable binding of ligand IBMP as compared to ligand IPMP. The binding affinity difference among the SBT series of ligands are all due to the cavity contributions.

The molecular recognition between the dry region in the binding pocket and the hydrophobic groups in the ligands is not unique for MUP. Fig. 4.7 provides another two examples where ligands with hydrophobic groups bind to the dry region of MUP receptor. In a previous work, Siebert and Hummer also observed a strong correlation between the

location of conserved nonpolar groups of ligands and the low water occupancy regions in the binding surface of the IQN17 peptide, a soluble analogue of the N-peptide coiled coil.[85] Fig. 4.8 displays the active sites of GLTP and the secretin pilot protein. In both cases, there is a large dry region in the binding pocket and a large portion of the hydrophobic groups of the ligand is located in that dry region, consistent with previous studies.[75] So the dry region in the receptor and the hydrophobic groups in the ligands may represent a general motif for molecular recognition. For GLTP, the ligand is a alkane chain and the whole binding pocket is dry except the entrance. There are no principal hydration sites identified by the WaterMap calculation for this system. For secretin pilot protein, the tail of the ligand is a carboxylic group, and only the middle part of the binding pocket is dry. There are two principal hydration waters near the entrance of the pocket identified by the WaterMap calculation.

4.3 Conclusion

We have augmented our WaterMap scoring function for computing free energy differences between congeneric ligands with a new term which models the free energy gain from ligand atoms occupying dry regions of the receptor. The results of the new scoring function are highly satisfactory for the data sets that we have examined, and require no adjustable parameters. Hence, our expectation is that this model will prove successful in other systems where dry regions exist.

This paper represents an initial effort to improve the core functionality embodied in the current WaterMap scoring function. There are clearly other augmentations that need to be made before the method can robustly handle a wide variety of test cases, most prominently an approach to treating protein-ligand interactions, particularly when these are not fully complementary, is required. Our objective is to systematically add new functionality, building on the success of the core approach, and render the method increasingly more accurate and reliable, while retaining the favorable computational properties that characterize the current methodology.

4.4 Systems and Simulations

The starting structures for the mouse major urinary protein (MUP), the bovine apoglycolipid transfer protein (GLTP) and the secretin pilot protein are taken from PDB with PDB ID 1znk, 1wbe, 1y9l respectively.[76; 77; 78] All the nonprotein molecules were then removed and protein preparation wizard[86] was used to modify the structures of the proteins for simulation. Protonation states were assigned assuming the systems are at pH 7.0. The proteins without the ligands, which we refer to the “apo” proteins, were inserted into water boxes using Maestro,[65] and water molecules that sterically overlapped with the proteins were removed. The size of each system was chosen to accommodate a minimum of 10 Å of water between the protein surface and the box walls. Counter ions were added to maintain electric neutrality. The systems were then relaxed and equilibrated for a series of minimizations and short molecular dynamics simulations using the standard relaxation protocol in Desmond.[31] To ensure equilibration between water in the binding pocket and bulk water, grand canonical Monte Carlo method is used to sample both the number of water molecules in the pocket and their positions using the solvate-pocket utility in Desmond during equilibration.[31]

The production simulations were done in NPT ensemble with a constant temperature of 300K and 1 atmospheric pressure.[33; 34; 35] The OPLS-AA force field was used for the protein, and the TIP4P water model[32] was used for the solvent, with a cut-off of 9 Å for Lennard-Jones interactions and a Particle-Mesh Ewald for electrostatic interactions.[37] During the simulation the protein heavy atoms were harmonically restrained to their initial positions. Data were taken from 10 ns production simulations for MUP and 2ns for GLTP and secretin pilot protein. Running the simulation for longer time does not change the results.

4.5 Methods

Our analytical effort focused on the hydration properties of the active sites of the “apo” proteins. The active site was defined as the region within 10 Å of where the ligand heavy atom would be but not closer than 2.8 Å to any heavy atom of the protein. We refer to this

region in the following as the binding pocket.

4.5.1 WaterMap calculation

The high occupancy principal hydration sites inside the binding pocket were identified and their associated enthalpy and entropy were calculated using the WaterMap program developed in our group.[3; 2] To be specific, water molecules inside the binding pocket were clustered into high occupancy hydration sites each of which is a sphere of 1 Å radius, and the enthalpy and entropy for each principal hydration site water were calculated using the inhomogeneous solvation theory.[4] Details of the implementation of the method are discussed in ref.[2].

4.5.2 Cavity calculation

The binding pocket is covered by a 3D grid with 1 Å spacing in each dimension. For each frame during the simulation, the positions of water oxygen atoms inside the binding pocket were recorded. If any water oxygen atom is closer than 3.3 Å to a grid point, that grid point is regarded as being occupied; otherwise the grid point is regarded as being unoccupied. In general one would have to have chosen different radii for different atom types, but here we constrained the heavy atoms of the protein so that only water molecules can enter into the cavities. Note here that, there may be more than one water molecule simultaneously occupying the same grid point, and that a given water molecule may simultaneously occupy several grid points. The probability, P_0 , for a grid point to be unoccupied is calculated and if it is ≥ 0.5 the cavity is considered to be dry. In fact, from the simulation, grid points which are identified as dry are found to be physically close to each other, and we draw a white line between neighboring dry grid and in this way identify the dry region displayed in the corresponding figures.

Note that, in bulk water there are on average 4.6 water molecules in a spherical volume of radius 3.3 Å, and the probability of this cavity being unoccupied by water is $P_0 \approx 10^{-4}$. Here, 3.3 Å is the size of the united atom methane. Both the hydration free energy of a methane particle and the potential of mean force (PMF) between two methane particles in neat water can be understood from information theory with a cavity of 3.3 Å radius.[87]

Thus grid points in the binding pocket of the protein with $P_0 \geq 0.5$ are clearly dry.

4.5.3 Protein-ligand binding affinity analysis

The binding affinity of each ligand to the protein receptor is decomposed into the WaterMap contribution and the cavity contribution. We conducted a structure alignment between the holo-protein-ligand complex from the PDB structure and the “apo” protein simulated. The WaterMap contribution was calculated through the displaced solvent functional introduced in ref.[2] Ligand heavy atoms close to the principal hydration sites were assigned a score, depending on the distance from the heavy atom to the hydration site and free energy difference between water in that hydration site and bulk water. Details of the Functional is in ref.[2]

For the cavity contribution, the probability, P_0 , to observe an empty spherical region with radius 3.3 Å centered on each ligand heavy atom was calculated. If the probability of the cavity being unoccupied by water is greater than 0.5, the binding affinity gain for the ligand atom occupying that dry cavity is

$$\Delta G = -kT \ln(P_0) - 2.36 \text{ (kcal/mol)} \quad (4.1)$$

Here P_0 is the probability of the cavity being unoccupied, and $-kT \ln(P_0)$ is the free energy to create a cavity of radius 3.3 Å inside the active site, and 2.36 kcal/mol is the solvation free energy of methane. As mentioned above, the free energy to “grow” a ligand heavy atom inside the binding pocket is the sum of the free energy to create a cavity and the free energy to turn on the interactions between that atom and the rest of the system. If the atom is in the dry region, and if the atom is nonpolar (from the simulation we found that all the ligand heavy atoms located in the dry region are nonpolar), then the free energy to turn on the interactions between that atom and the rest of the system is almost zero. (Lennard-Jones interactions are short ranged, and there are no surrounding water molecules if it is in the dry region.) So the two terms in Eq. (4.1) are approximately the free energy to “grow” a ligand heavy atom inside the binding pocket and that in bulk water, and their difference gives the contribution to the binding affinity from that atom. Here, we assume that the size of each ligand heavy atom is comparable to the size of a united atom methane. The total

cavity contribution is a summation over all ligand heavy atoms located in the dry region. Note that some of the ligand heavy atoms may have partially overlapped cavities. We treat them as independent of each other, which is equivalent to assuming pairwise additivity.[84] The error for this approximation is relatively small, because they overlap both in the active site and in bulk water, and there is a large cancellation for the effects. Even in the extreme case, where the cavity for a ligand heavy atom is fully overlapping with existing cavities, the free energy to create that additional cavity, which is 0 in this case, is not quite different from $-kT \ln(0.5) = 0.4$ kcal/mol, the maximum free energy to create that cavity in the dry region, and the error in the real case is much smaller than this number.

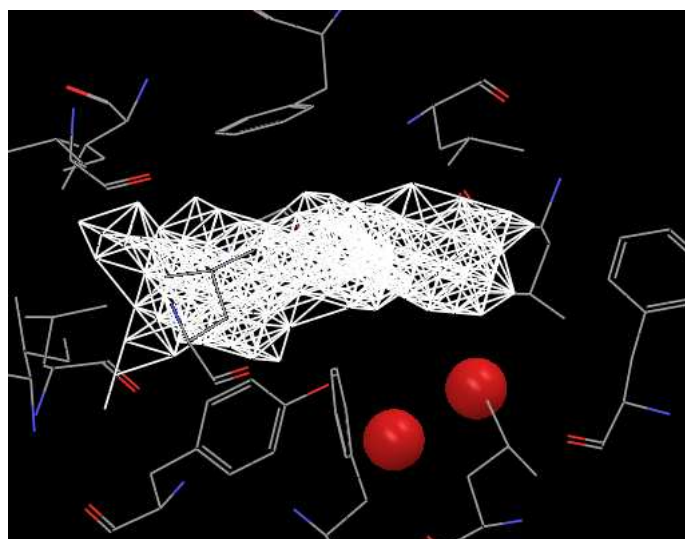


Figure 4.1: The principal hydration sites and the dry region in the binding pocket of MUP. The two principal hydration site waters are displayed in red sphere, and the dry region is displayed by white dots connected with white lines. The side chains surrounding the binding pocket are also displayed. A large region of the binding pocket is dry.

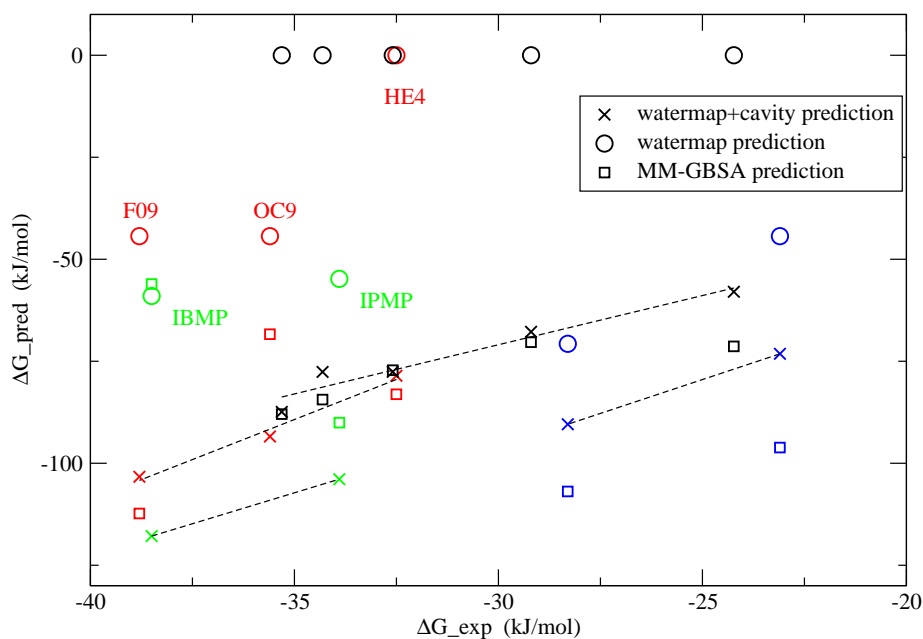


Figure 4.3: The Watermap, Watermap+cavity and MM-GBSA predictions for the binding affinities of different ligands to the MUP receptor versus the experimental data. The Watermap predictions are displayed as circles, Watermap+cavity predictions are displayed as crosses, and MM-GBSA predictions are displayed as rectangles. The ligands belonging to different groups are indicated by different colors. While the Watermap and MM-GBSA predictions fail to rank-order most of the congeneric ligands, Watermap+cavity predictions correctly rank-order all the congeneric ligands in each group.

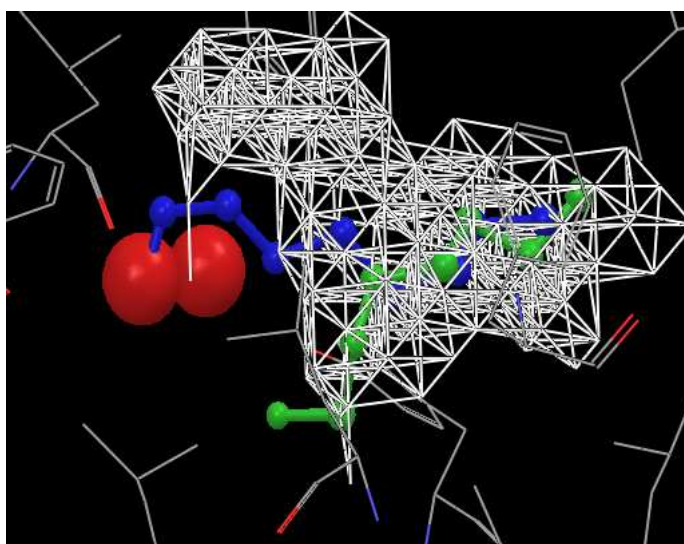


Figure 4.5: Ligands HE4 (green) and OC9 (blue) in the binding pocket of MUP. Ligand OC9 displaces one of the principal hydration water, while ligand HE4 does not. So the WaterMap predicted binding affinity for ligand OC9 is much more favorable than for ligand HE4, much larger than the experimentally measured binding affinity difference. However, a large portion of ligand HE4 is located in the dry region, so the cavity contribution is more favorable for ligand HE4. Combined with WaterMap and cavity contribution, the experimentally measured binding affinity difference can be easily explained.

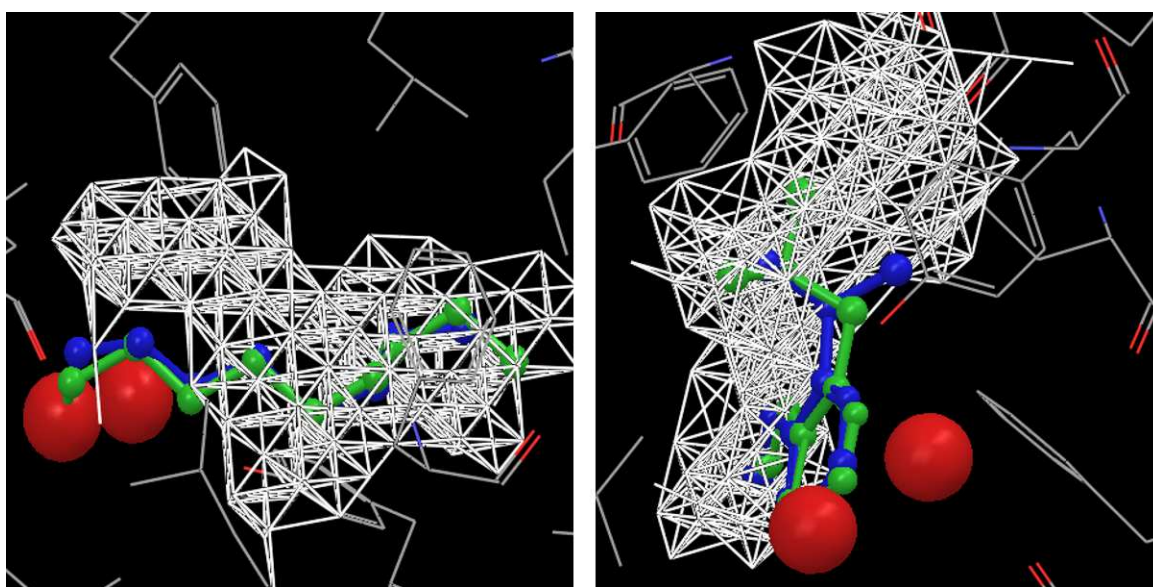


Figure 4.6: (a) Ligands OC9 (blue) and F09 (green) in the binding pocket of MUP. They have similar structure in the principal hydration site, so the WaterMap predicts their binding affinities are the same. However, ligand F09 has one more atom located in the dry region, which leads to the stronger binding of ligand F09 than ligand OC9, verified by experimental data. (b) Ligand IBMP (green) and IPMP (blue) in the binding pocket of MUP. Ligand IBMP has one more atom located in the dry region of the pocket, leading to stronger binding of IBMP than IPMP.

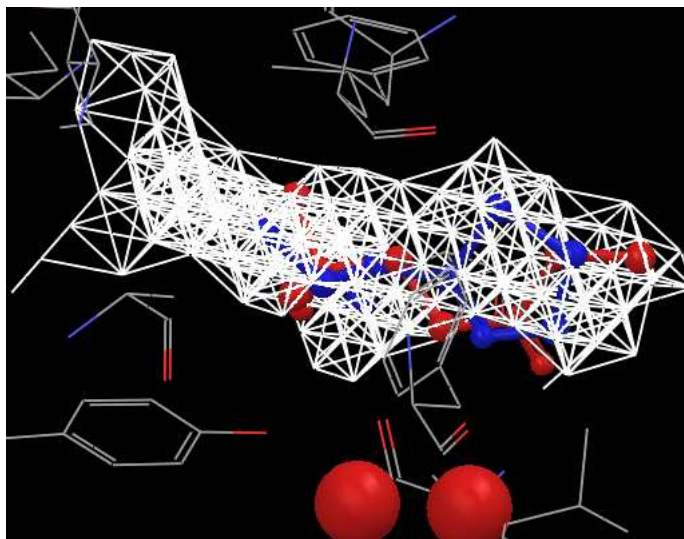


Figure 4.7: Ligands LTL (red) and TZL (blue) binding to the MUP receptor. Large portions of the ligand atoms are located in the dry cavity region.

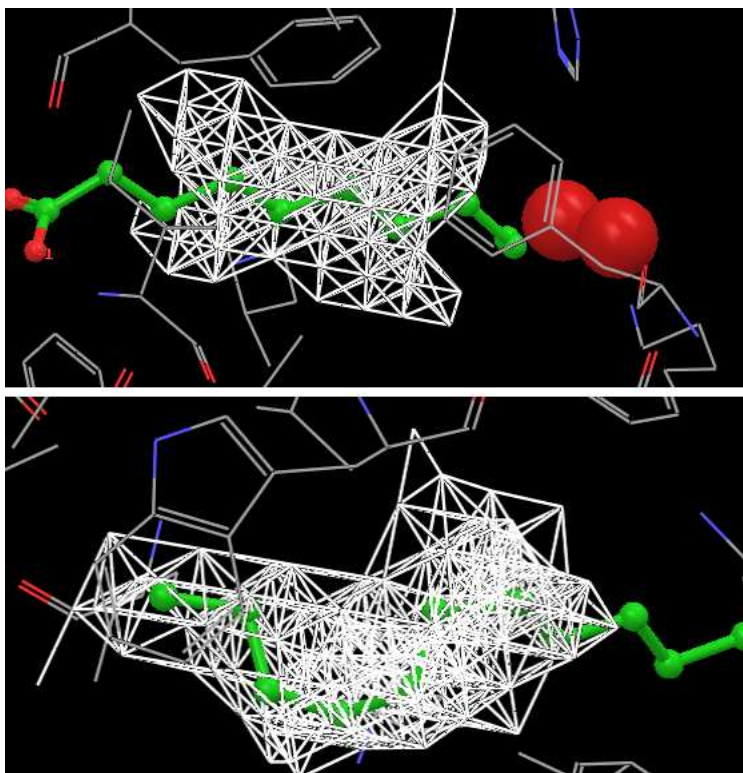


Figure 4.8: The binding pockets of the secretin pilot protein (upper) and GLTP (below). In both cases, there is a large dry region in the binding pocket and a large portion of the hydrophobic groups of the ligands are located in that dry region. For GLTP, the ligand is an alkane chain and the whole binding pocket is dry except the entrance. For secretin pilot protein, the tail of the ligand is a carboxylic group, and only the middle part of the binding pocket is dry. There are two principal hydration waters near the entrance of the pocket identified by the WaterMap calculation.

Table 4.1: Watermap and cavity contributions to the binding affinities for different ligands binding to MUP receptor

Binding affinities	PE9	HE2	HE4	OC9	F09	IBMP	IPMP	SBT	PT	IPT	ET	MT
Exp	-23.1	-28.3	-32.5	-35.6	-38.8	-38.5	-33.9	-35.3	-34.3	-32.6	-29.2	-24.2
WaterMap	-44.4	-70.7	0.0	-44.4	-44.4	-59.0	-54.8	0.0	0.0	0.0	0.0	0.0
Cavity	-28.8	-19.7	-78.6	-49.1	-59.0	-58.9	-49.1	-87.4	-77.6	-77.6	-67.8	-58.8
Total	-73.2	-90.4	-78.6	-93.5	-103.4	-117.9	-103.9	-87.4	-77.6	-77.6	-67.8	-58.8
MM-GBSA	-96.2	-106.9	-83.1	-68.4	-112.3	-56.1	-90.0	-88.0	-84.4	-77.2	-70.3	-71.3

Note1: Free energies in kJ/mol. For ligand PE9, the PDB structure (PDB ID 1ZND) contains two ligands (with two binding modes). However, experimental ITC data indicate a binding stoichiometry of approximately 1 for PE9,[76] so only the binding mode with stronger binding affinity was analyzed. The predicted binding affinities for the two binding modes agree with experimental data.

Note 2: Ligands PE9 and HE2 bind in a similar orientation whereas ligands HE4, OC9, and F09 bind in an alternate orientation. So they are considered as two groups.[76]

Note 3: For SBT series of ligands, PDB structures are only available for SBT/MUP complex, and structures of other ligands were obtained by removing the appropriate carbon atoms from SBT.

Part II

Development of FEP/REST

Chapter 5

Introduction of the FEP/REST method

In the late stage drug design projects, when important decisions about how to modify and refine the lead molecule is made, highly accurate and reliable binding affinity results are required, and explicit solvent model free energy perturbation (FEP) molecular dynamics simulation represents one of the most rigorous methods to calculate Protein-Ligand binding affinities. In spite of the potentially large impact FEP may have on drug design projects, practical applications in an industrial context have been limited over the past decade. High accuracy and reliability in the methodology are required for developing FEP into a true engineering platform for drug candidate optimization, but neither has yet been demonstrated by existing implementations.

Two types of challenges stand in the way of applying FEP into real drug design projects. Firstly, converging explicit solvent simulations to the desired precision is far from trivial, even with the immense computing power that is currently available using low cost multiprocessor clusters or cloud computing platforms of various types. This problem is more severe when the protein or the ligands adopt different conformations upon alchemical transition from one ligand to another or upon the protein-ligand binding process, and there are large energy barriers separating the relevant conformations. In these cases, the protein or the ligands may remain kinetically trapped in the starting conformation for a long time during

brute force MD simulation, and the calculated binding affinities are dependent on the starting conformation to do the simulation, giving rise to the well known quasi-nonergodicity problem in FEP. Secondly, errors in the potential energy models must be reduced to the point where they lead to errors in a converged calculation that are smaller than the desired errors in relative binding affinities compared to experiment, typically on the order of 0.5 kcal/mole. At the present stage, it is more imperative to solve the sampling problem in FEP, since before the precision of the free energy results is promised, it is impossible to study and improve the accuracy of the force field.

Our strategy to solve the sampling problem in FEP simulation is to combine enhanced sampling technique with normal FEP. The general method to get enhanced sampling is temperature replica exchange method (TREM). In TREM, a number of replicas are started simultaneously each on a different temperature, and attempts to exchange configuration between neighboring replicas are made during regular intervals of the simulation. If the temperature of the highest level replica is high enough that it can sample different regions of phase space, through replica exchange the lower level replica can also sample different regions of phase space. However, the number of replicas required in normal TREM is very large, (proportional to \sqrt{f} , where f is the number of degrees of freedom of the whole system) which limits the application of TREM to large systems like protein ligand binding.

In this section, a new efficient enhanced sampling technique is introduced, and it was combined with normal FEP to solve the sampling problem in brute force FEP simulation. In Chapter 6, the details of the proposed enhanced sampling method called REST2 are presented, and the connections with previous enhanced sampling methods, and the reasons why it is more efficient are discussed. In REST2, we separate the simulation system into two regions, the “hot” region, (the region we are interested in, usually including the ligands and protein residues surrounding the binding pocket) and the “cold” region. We scale the Hamiltonian of the system in such a way that the effective temperature of the “hot” region is increased and the effective temperature of the “cold” region is at temperature T_0 for all replicas. In this way, a small number of replicas are sufficient to maintain a good exchange efficiency. Example application of REST2 on two protein systems, trpcage and β -hairpin, problematic for previous version of enhanced sampling method, demonstrated the efficiency

of REST2 on sampling the different conformations of the protein.

In Chapter 7, we combine the enhanced sampling technique introduced in Chapter 6 with normal FEP to solve the sampling problem in FEP. Previous efforts trying to combine enhanced sampling into FEP require 2-d replica exchange protocol, one in Hamiltonian axis like normal FEP, the other in the “boosting potential” axis, where the boosting potential will cancel the potential of mean force (PMF) along the slow degree of freedom to get enhanced sampling. This method requires a large number of replicas and requires the prior known knowledge of the slow degree of freedom. In the FEP/REST method we designed here, the enhanced sampling technique REST2 is combined with normal FEP through one dimensional replica exchange protocol. In this way, the computational expense of FEP/REST is comparable with normal FEP, and it can be easily used in real Protein-Ligand binding problems of medicinal interest where nothing is known about the slow degree of freedom. Application of FEP/REST on two modifications, the T4L/L99A and the Thrombin systems, both leading to large structural reorganizations, one in the protein and the other in the ligands, demonstrates the superior convergence of the free energy as indicated both by consistency of the results (independence from the starting conformation) and agreement with experimental binding affinity data.

Chapter 6

Replica Exchange with Solute Scaling: A more efficient version of Replica Exchange with Solute Tempering

Abstract

A small change in the Hamiltonian scaling in replica exchange with solute Tempering (REST) is found to improve its sampling efficiency greatly especially for the sampling of aqueous protein solutions in which there are large scale solute conformation changes. Like the original REST (REST1), the new version (which we call REST2) also bypasses the poor scaling with system size of the standard temperature replica exchange method (TREM), reducing the number of replicas (parallel processes) from what must be used in TREM. This reduction is accomplished by deforming the Hamiltonian function for each replica in such a way that the acceptance probability for the exchange of replica configurations does not depend on the number of explicit water molecules in the system. For proof of concept, REST2 is compared with TREM and with REST1 for the folding of the trpcage and β -hairpin in water. The comparisons confirm that REST2 greatly reduces the number of CPUs required

by regular replica exchange and greatly increases the sampling efficiency over REST1. This method reduces the CPU time required for calculating thermodynamic averages and for the *ab initio* folding of proteins in explicit water.

6.1 Introduction

Sampling the conformational space of complex biophysical systems, such as proteins, remains a significant challenge, because the barriers separating the local energy minima are usually much higher than $k_B T$, leading to kinetic “trapping” for long periods of time and quasi-ergodicity in the simulations. The Temperature Replica Exchange Method (TREM) has attracted attention recently as a means for overcoming the problem of quasi-ergodicity.[88; 89; 90; 91; 92; 93] However, the number of replicas required to get efficient sampling in normal TREM scales as \sqrt{f} , where f is the number of degrees of freedom of the whole system, which often limits the applicability of TREM for large systems. To overcome this problem, we recently devised the method “Replica Exchange with Solute Tempering” (REST1),[94] in which only the solute biomolecule is effectively heated up while the solvent remains cold in higher temperature replicas, so that the number of the replicas required is much reduced. It has been shown that the required number of replicas in REST1 scales as $\sqrt{f_p}$, where f_p is the number of degrees of freedom of the solute, and the speedup versus the TREM, in terms of converging to the correct underlying distribution, is $O(\sqrt{(f/f_p)})$ for small solutes like alanine dipeptide.[94] However, when applying REST1 to large systems involving large conformational changes, like the trpcage and β hairpin, it was found that REST1 can be less efficient than TREM.[95] For example, we observed that the lower temperature replicas stayed in the folded structure, the higher temperature replicas stayed in the extended structure, and the exchange between those two conformations was very low.[95] Moors *et. al.*[96] and Terakawa *et. al.* [97] independently modified our REST1 scaling factor for E_{pw} so that approach could be easily run in GROMACS. Moors *et. al.* included only part of the protein in the “hot region”, keeping the rest of it “cold” and called their method “Replica Exchange with Flexible Tempering” (REFT). Interestingly, they observed an improved sampling efficiency in sampling a particular reaction coordinate involving the opening and closing of

the binding pocket in T4 Lysozyme and suggested that the improved sampling efficiency for their method over REST1 occurred because in REST1 all of the protein degrees of freedom contribute to the acceptance probability for replica exchange whereas in REFT only those degrees of freedom involved in the opening and closing of the pocket contribute. Thus the acceptance probability for replica exchange is larger in REFT than in REST1. As we shall see, this is not the only reason for the observed improvement.

In this paper we use the modified scaling of the Hamiltonians suggested by Moors *et. al.*[96] and Terakawa *et. al.* [97] instead of the original scaling of our REST1, to see if it samples the folded and unfolded conformations of proteins more efficiently than REST1, although all of the protein degrees of freedom are allowed to be hot in this study. For simplicity we call REST with this new scaling REST2. Application of REST2 to the trpcage and the β -hairpin systems, the same systems that were problematic when sampled by REST1, indicates that REST2 is much more efficient than REST1 in sampling the conformational space of large systems undergoing large conformation changes. In what follows, we will present the scaling, its connection with our original scaling of REST1, and the results for the trpcage and β -hairpin systems for REST1, REST2, and TREM. We will also discuss the reasons for the improvement found with REST2.

6.2 Methodology

In REST1, the total interaction energy of the system was decomposed into three components: the protein intra-molecular energy, E_{pp} ; the interaction energy between the protein and water, E_{pw} ; and the self interaction energy between water molecules, E_{ww} . Replicas running at different temperatures then evolve through different Hamiltonians involving relative scalings of these three components. To be specific, the replica running at temperature T_m has the following potential energy:

$$E_m^{REST1}(X) = E_{pp}(X) + \frac{\beta_0 + \beta_m}{2\beta_m} E_{pw}(X) + \frac{\beta_0}{\beta_m} E_{ww}(X). \quad (6.1)$$

Here, X represents the configuration of the whole system, $\beta_m = 1/k_B T_m$ and T_0 is the temperature that we are interested in. The potential for replica running at T_0 reduces to the normal potential.

Imposing the detailed balance condition, the acceptance ratio for the exchange between two replicas m and n depends on the following energy difference:

$$\Delta_{mn}(REST1) = (\beta_m - \beta_n)[(E_{pp}(X_n) + \frac{1}{2}E_{pw}(X_n)) - (E_{pp}(X_m) + \frac{1}{2}E_{pw}(X_m))]. \quad (6.2)$$

Note that the water self interaction energy, E_{ww} , does not appear in the acceptance ratio formula, and this is the reason why only a relatively small number of replicas are sufficient to achieve good exchange probabilities in REST1.

In REST1, both the potential energy and the temperature are different for different replicas. According to the law of corresponding states, the thermodynamic properties of a system with potential energy E_m at temperature T_m , are the same as those for a system with potential energy $(T_0/T_m)E_m$ at temperature T_0 . So instead of using different potential energies and different temperatures for different replicas, we can run all the replicas at the same temperature albeit on different potential energy surfaces using the Hamiltonian Replica Exchange Method (H-REM).[98; 99] To be specific, in REST2, all of the replicas are run at the same temperature T_0 , but the potential energy for replica m is scaled differently,

$$E_m^{REST2}(X) = \frac{\beta_m}{\beta_0} E_{pp}(X) + \sqrt{\frac{\beta_m}{\beta_0}} E_{pw}(X) + E_{ww}(X). \quad (6.3)$$

In REST1, enhanced sampling of the protein conformations is achieved by increasing the temperature of the protein, but between attempted exchanges with neighboring replicas, replica m moves on the full intramolecular protein potential energy surface with high energy barriers, although the other energy terms are scaled. In REST2, enhanced sampling is achieved through scaling the intra-molecular potential energy of the protein by (β_m/β_0) , a number smaller than 1, so that the barriers separating different conformations are lowered. Thus between attempted replica exchanges replica m moves on a modified potential surface where the barriers in the intra protein force field are reduced by the scaling. We call REST with this new scaling ‘‘Replica Exchange with Solute Scaling’’ (REST2). Thus REST1 and REST2 arrive at the final distribution at temperature T_0 by different but rigorously correct routes. The acceptance criteria for replica exchanges are different in REST1 and REST2 but the Hamiltonians for the MD trajectories are also different in such a way that the long time sampling at T_0 should converge to the same ensemble for REST1 and REST2,

albeit with different rates of convergence for the two methods. In REST2, the differences between different replicas are the different scaling factors used, but to make connections with REST1 we will keep using the term “temperature” for replica m which means the effective temperature of the protein with the unscaled potential energy.

Note that the scaling factor used in REST2 for the interaction energy between the solute and water for replica m is $\sqrt{(\beta_m/\beta_0)}$, which is different from $(\beta_0 + \beta_m)/2\beta_0$ used in REST1 (Eq. (6.1)). The interaction energy in Eq. (6.3) can be easily achieved by scaling the bonded interaction energy terms, the Lennard-Jones ϵ parameters, and the charges of the solute atoms by (β_m/β_0) , (β_m/β_0) , and $\sqrt{(\beta_m/\beta_0)}$ respectively, and the scaling factor for the E_{pw} term, $\sqrt{(\beta_m/\beta_0)}$, follows naturally from standard combination rules for LJ interactions. This minor change of the scaling factor for E_{pw} term, suggested in the original REST paper but not appreciated at that time, proves to be important for the better performance of the REST2. In addition, we find that scaling the bond stretch and bond angle terms does not help the sampling, so in practice only the dihedral angle terms in the bonded interaction of the solute are scaled and this makes the transition between different conformations of the solute faster.

Another consequence of the different scaling factors used for the E_{pw} term in REST1 and REST2 is the different acceptance ratio formulas in these two methods. It is easy to show by imposing detailed balance condition that the acceptance ratio for exchange between replicas m and n in REST2 is determined by:

$$\Delta_{mn}(REST2) = (\beta_m - \beta_n) \left[(E_{pp}(X_n) - E_{pp}(X_m)) + \frac{\sqrt{\beta_0}}{\sqrt{\beta_m} + \sqrt{\beta_n}} (E_{pw}(X_n) - E_{pw}(X_m)) \right]. \quad (6.4)$$

For replica m , the exchanges to neighboring replicas $m - 1$ and $m + 1$, are determined by the fluctuation of $E_{pp} + \sqrt{\beta_0}/(\sqrt{\beta_m} + \sqrt{\beta_{m-1}})E_{pw}$ and $E_{pp} + \sqrt{\beta_0}/(\sqrt{\beta_m} + \sqrt{\beta_{m+1}})E_{pw}$ respectively. Thus for discussion purposes, but not in the simulations, the fluctuation of $E_{pp} + (1/2)\sqrt{(\beta_0/\beta_m)}E_{pw}$ can be thought to determine the acceptance ratios for exchanges of the replica at temperature T_m to neighboring replicas because, to a good approximation, $\beta_{m-1} \approx \beta_m \approx \beta_{m+1}$. Note then that the difference in the acceptance ratio formulas between REST1 and REST2 lies in the replacement of the factor $1/2$ by the factor $1/2\sqrt{(\beta_0/\beta_m)}$ multiplying the term E_{pw} . This difference is also partly responsible for the improvement of

REST2 over REST1 due to an approximate cancellation of E_{pp} and the scaled E_{pw} in the acceptance probability or equivalently in Δ_{nm} of REST2 but not in REST1, as we shall see.

6.3 Results and Discussion

Using REST2, we simulated the trpcage system with DESMOND [31] using 10 replicas with effective temperatures of the solute at 300K, 322K, 345K, 368K, 394K, 423K, 455K, 491K, 529K, 572K. The OPLSAA force field[100] was used for the protein and the Tip4p model[32] was used for water. All the replicas were started from the ‘native’ NMR structure (PDB ID 1L2Y)[101] and the simulation lasted for 20ns. Conformations of the protein were saved every 0.5ps, and exchange of configurations between neighboring replicas are attempted every 2ps with an average acceptance ratio of about 30%.

Four representative temperature trajectories for the trpcage replicas started at 300K, 368K, 455K and 572K in the folded state are displayed in figure 6.1. It can be seen that the temperature trajectory for each replica visits all of the temperatures many times, even during the first 5ns of the simulation, and all of the replicas visit any given temperature many times during the simulation. This is a good indication of the efficiency of the sampling. By comparison, none of the temperature trajectories using REST1 were able to visit all of the temperatures during a 5ns simulation for the same system (see Fig. 6b in ref. [95]). In REST2 the time interval for attempted exchange was 2ps while in the REST1 simulation 0.4ps was used. We expect that even more rapid diffusion in temperature space could be achieved if shorter time intervals between attempted exchanges were used in REST2.

In the REST1 simulations, it was observed that only the folded structures were sampled at the lower temperatures while the folded structures were rarely sampled at higher temperatures like 572K after an initial equilibration phase (see Fig. 7a in ref.[95]). The REST2 simulation of the protein heavy atom deviation (RMSD) from the native structure is displayed in figure 6.2 for replicas with effective temperature of the protein at 300K, 423K, and 572K. It is clear that both the folded structure and the unfolded structures are sampled at 300K even in the first 5ns simulation (inset of figure 6.2). At the intermediate temperature (423K) the folded and unfolded structures are sampled with almost equal probability, and

the unfolded structures dominate at high temperature (572K). But unlike in REST1, the folded structures are also sampled at 572K after the initial equilibration phase.

The β -hairpin system is likewise more efficiently sampled by REST2. For the same number of replicas and the same temperature levels used in REST1[95], the temperature trajectories for three representative replicas, initially at low ($T = 310K$), intermediate ($T = 419K$), and high ($T = 684K$) temperatures, are shown in figure 6.3a. We also determined the protein heavy atom RMSD versus time at each of the above temperatures when replicas visited those temperatures which are shown in figure 6.3b. With a time interval of 2ps for attempted exchange, each replica is able to visit all the temperatures within 5ns and both the folded and unfolded structures are sampled at low and high temperatures. By comparison, using REST1, none of the replicas were able to visit all the temperatures (See Fig 3b in reference, [95]) whereas the low temperature replicas stayed folded and the high temperature replicas stayed unfolded after the initial stage (See Fig 4 in reference. [95]) Thus REST2 is clearly superior to REST1 for both the trpcage and the β -hairpin.

The different scaling factors used for the E_{pw} term in REST1 and REST2 are responsible for the improvement of REST2 over REST1 as expected from the discussion given in the previous section. Consider the constant temperature molecular dynamics trajectory between attempted replica exchanges. In REST1, the scaling factor for the E_{pw} term was $(\beta_0 + \beta_m)/2\beta_m$. In the limit when $T_m \rightarrow \infty$, REST1 will effectively sample the distribution $\exp(-\beta_0(E_{pw}/2 + E_{ww}))$. Since the unfolded structure has more favorable solute water interactions than the folded structure, replicas at higher temperature will sample the unfolded structure with dominating probability in REST1, and this was indeed observed in REST1 simulations for the trpcage as well as for the β -hairpin.[95] In REST1, the replicas at high temperatures can not sample the whole conformational space efficiently and replicas at high and low temperatures sample completely different regions of conformation space. This is one of the reasons for the observed inefficient sampling in REST1. By comparison, in REST2, we use a scaling factor $\sqrt{\beta_m/\beta_0}$ for the E_{pw} term. In the limit when $T_m \rightarrow \infty$, REST2 will effectively sample the distribution $\exp(-\beta_0 E_{ww})$. So both the folded and unfolded structures are sampled efficiently during the trajectories between attempted replica exchanges for the higher temperature replicas in REST2, and this is one of the reasons for

why REST2 is more efficient than REST1. This is shown in fig. 6.4, where it can be seen that in a constant temperature MD simulation using the scaled Hamiltonian of REST2 at high temperature the heavy atom RMSD for the β -hairpin fluctuates from the native structure from values close to 2.5Å to 5Å and back again, whereas in fig. 4 in ref[95] it stayed above 4Å after short times.

Fig. 6.5 displays the relation between the intra-molecular potential energy of the trpcage system, E_{pp} , and the interaction energy between the trpcage and water, E_{pw} , for replicas with different effective temperatures of the protein in REST2. At each temperature, there is a strong anti-correlation between those two terms. This is easy to understand: the more extended the structure is, the less favorable the intra-molecular potential energy of the protein, and the more favorable the interaction energy between the protein and water (which scales with the surface area of the protein). With increasing temperature, the probability for the unfolded structure gets larger, and the intra-molecular potential energy of the protein gets less negative. For the interaction energy between the protein and water, there are two counterbalancing effects. On the one hand, the higher the temperature, the more favorable the unfolded structure, and the more favorable the interaction between water and protein. On the other hand, every single component of the potential energy would increase with increasing temperature because of the generalized equi-partition theorem. This is exactly what we observe in fig. 6.5: with increasing temperature, the E_{pp} term get less negative while the E_{pw} term increases very slowly because of the compensation of the two effects mentioned above. This is clearly demonstrated in figure 6.6a and figure 6.6b, where the distribution of E_{pp} and E_{pw} are displayed for replicas at different temperatures. At 394K, both E_{pp} and E_{pw} are binomially distributed with the folded and unfolded structures almost equal probability. Below 394K, the folded structure dominates and above 394K the unfolded structure dominates. This is the reason why replicas running below 394K and above 394K were not able to exchange efficiently in REST1. [95] While E_{pp} increases monotonically with increasing temperature, the behavior of E_{pw} is more complicated. Below 394K, the center of the distribution for E_{pw} gets less negative, but the distribution gets boarder in the left tail of the distribution because of the increased probability of extended structures. Above 394K, E_{pw} increases with temperature because of the equipartition theorem.

The absence of a compensating term E_{ww} in the replica exchange probability of REST1 was suggested in our previous paper[95] to explain the observed better performance of TREM for the exchange between folded and unfolded structures where there is a big difference in the energies of these two states,[95] but in the AppendixB we show why we now think that this is not the reason for this difference. Actually the compensation or lack of compensation between E_{pp} and the scaled E_{pw} is more important than the loss of any compensation between E_{ww} and E_{pp} .

In REST1, it is the fluctuation of $(E_{pp} + (1/2)E_{pw})$ that determines the acceptance ratio. While the two terms can compensate each other to some extent, they both tend to increase with increasing temperature of the solute. By comparison, in REST2, it is the fluctuation of $E_{pp} + (1/2)\sqrt{(\beta_0/\beta_m)}E_{pw}$ that determines the acceptance ratio. Since $\sqrt{(\beta_0/\beta_m)}$ increases with increasing temperature of the solute, it will compensate the decrease of the magnitude of E_{pw} . (The E_{pw} term is negative, and the magnitude of it decreases with increasing temperature.) Fig. 6.6(c) displays the distribution of $(1/2)\sqrt{\beta_0/\beta_m}E_{pw}$ for replicas at different temperatures. It is quite clear that the factor $\sqrt{\beta_0/\beta_m}$ perfectly compensates the increase of E_{pw} . Below 394K, the distribution is centered at about $-380kcal/mol$ corresponding to the folded structure; above 394K, the distribution is centered at about $-495kcal/mol$ corresponding to the unfolded structure. At 394K, the folded and unfolded structures are almost equally distributed. With increasing temperature, the probability of the unfolded structure increases and the probability of folded structure decreases. The difference in the E_{pp} term between the folded and unfolded structures is compensated by the difference in $(1/2)\sqrt{\beta_0/\beta_m}E_{pw}$ term, which makes the distribution of $E_{pp} + (1/2)\sqrt{\beta_0/\beta_m}E_{pw}$ have sufficient overlap for neighboring replicas. (Figure 6.6d) The approximate cancellation of the contributions E_{pp} and E_{pw} in REST2 and their smaller cancellation in REST1 makes the acceptance ratio for replica exchange larger in REST2 than in REST1, and this is part of the reason for the more efficient sampling in REST2 than in REST1. For the trpcage system studied here, with the same number of replicas and the same temperature levels, we obtained an average acceptance ratio for REST2 of 30%, while in REST1 only 20% was obtained. In addition the more frequent barrier crossings in the the MD trajectories of REST2 than in REST1 contributes considerably to the better efficiency of REST2.

As mentioned, the rate of convergence of REST1 (relative to TREM) to the correct underlying distribution was shown to scale as $O(\sqrt{(f/f_p)})$ for small solutes like alanine dipeptide; however, for systems involving large conformation changes, REST1 fails to achieve this expected speed up.[94] The results presented in the above sections clearly demonstrate that REST2 is much more efficient than REST1 for sampling systems with large conformational change, but does REST2 do better compared to TREM for these problematic systems? To answer this question, we simulated the trpcage system starting from an almost fully extended configuration using both TREM and REST2. As before, 10 replicas were used for REST2, and 48 replicas were needed in TREM to maintain an appropriate acceptance ratio. It should be noted that the replica exchange ratio for TREM is 10% whereas for REST2 it is 30% so that we could have used fewer replicas in REST2 to get the same exchange ratio as in TREM. Within 2ns simulations, none of the replicas in TREM were able to visit all the temperatures while in REST2 all 10 replicas were able to visit all of the the temperatures, indicating that REST2 is much more efficient in diffusing through temperature space than TREM. The distribution of intramolecular energy of the trpcage for the lowest level replica calculated from TREM and REST2 are shown in fig. 6.7. For trajectories of the same length, REST2 samples a broader region in conformation space than TREM, and in addition the cpu cost of generating equal length trajectories is greater for TREM than for REST2 (see the appendix B).

6.4 Conclusion

We find that Replica Exchange with Solute Scaling (REST2) more efficiently samples the conformation space than REST1. We used a different scaling factor for the interaction energy between the protein and water, E_{pw} , than we used in REST1. Application of REST2 to the trpcage and β -hairpin systems results in an improvement over REST1 in sampling large systems involving large conformational energy changes. The better efficiency of REST2 over REST1 arises because there is a greater cancellation between the scaled terms E_{pp} and E_{pw} in REST2 than in REST1. This gives rise to REST2's larger replica exchange probability than REST1's, and also to its better sampling between replica exchanges at

high temperature, as we now discuss. For example, for $T_0 = 300K$ and $T_m = 600K$, the deformed potential for REST2 is $E_m^{(REST2)} = 0.5E_{pp} + 0.71E_{pw} + E_{ww}$, but it is run at $T_0 = 300K$; whereas for REST1 it is $E_m^{(REST1)} = E_{pp} + 1.5E_{pw} + 2E_{ww}$ but it is run at $T_m = 600K$. The exponents in the Boltzmann factors for these two cases, are

$$\beta_0 E_m^{REST2} = \beta_m [E_{pp} + 1.41E_{pw} + 2.0E_{ww}],$$

and

$$\beta_m E_m^{(REST1)} = \beta_m [E_{pp} + 1.5E_{pw} + 2.0E_{ww}].$$

The only difference between these is due to the different scaling factors of the E_{pw} term, which for $T_0 = 300K$ and $T_{max} = 600K$, is $(1/2)(\beta_0 + \beta_m)/\beta_m = 1.5$ vs $\sqrt{\beta_0/\beta_m} \approx 1.41$. We have seen in Fig. 6.5 that the E_{pw} term is usually much larger in magnitude than the E_{pp} term, so a small change in the scaling factor of the E_{pw} term leads to better sampling efficiencies for the high temperature MD between replica exchanges for REST2 than REST1. For the trpcage and β -hairpin systems studied here, in REST2 both folded and unfolded conformations are sampled at higher temperature replicas whereas in REST1 only the unfolded conformational space of the solute are sampled at higher temperature replicas. Because of the larger replica exchange probability and because of better constant temperature sampling these folded and unfolded conformations filter down to the replica at the temperature of interest, T_0 . In addition, since all the replicas are running at the same temperature in REST2, there is no need to rescale the velocity during exchange process which will save some computer time and makes it easier to implement in various MD programs. We also found REST2 to be more efficient in sampling the trpcage than TREM, because of the much smaller number of replicas and faster cpu times required to generate the MD trajectories in REST2 compared to TREM. In addition, the lowest level replica was found to explore a larger region of energy space for REST2 than for TREM for the same MD trajectory lengths for each replica. Thus, REST2 speeds up the sampling of the trpcage, by at least a factor of 9.6 over TREM.

We believe that REST2 should be used for investigating large protein-water systems especially when there are large conformation energy changes in the protein. The improvement comes from: (a) the larger replica exchange probabilities and concomitantly the smaller

number of replicas that can be used, and (b) the more efficient MD sampling of the conformational states between replica exchanges on the upper replica potential energy surfaces.

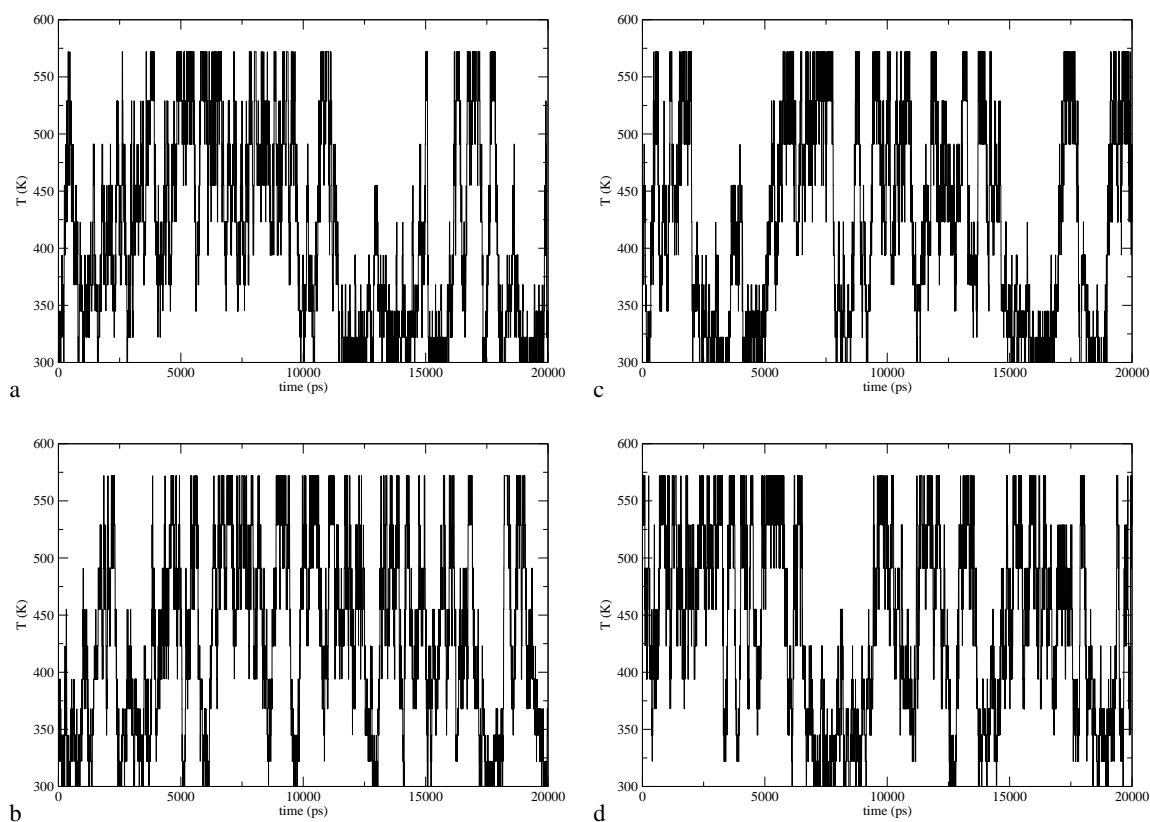


Figure 6.1: Temperature trajectories of four representative replicas with the effective temperature of the protein started at 300K (a), 368K (b), 455K (c), and 572K (d) for the trpcage system starting from the native structure. It should be noted that the temperatures referred to are the effective temperatures of the protein which arises from the scaling of the force field parameters of the protein, while the actual simulation is done at temperature T_0 .

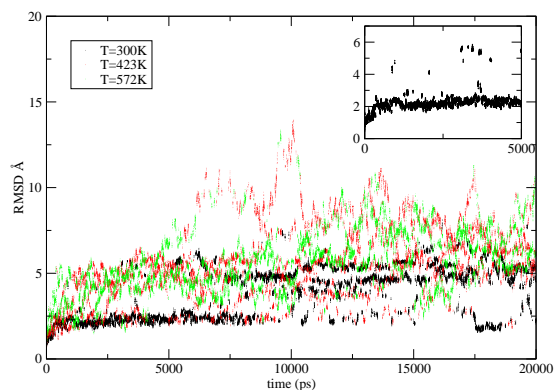


Figure 6.2: Protein heavy atom RMS deviation from the native structure as a function of simulation time for replicas with different effective temperatures of the protein for the trpcage system. Inset of the figure highlights the RMSD for replica at effective temperature 300K in the first 5ns simulation.

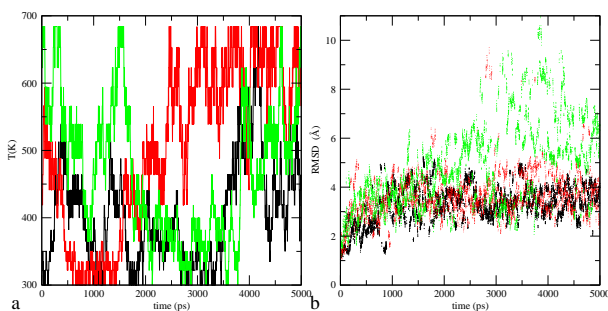


Figure 6.3: a. The temperature trajectories for three representative replicas with the effective temperature of the protein initially at low ($T = 310K$), intermediate ($T = 419K$), and high ($T = 684K$) temperatures for the β -hairpin system. b. The protein heavy atom RMSD versus time at each of the above temperatures when replicas visit those temperatures. (Black, $T = 310K$; Red, $T = 419K$; Green, $T = 684K$)

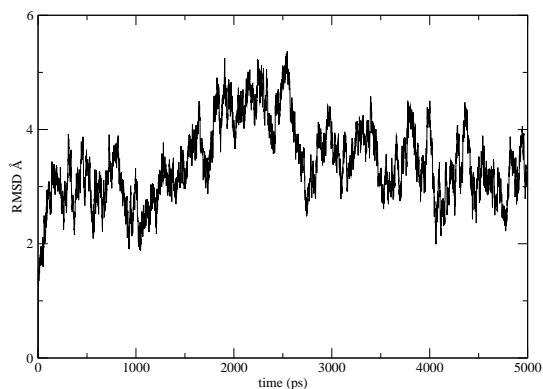


Figure 6.4: The heavy atom RMSD from the native structure of the β -hairpin as a function of simulation time with the effective temperature of the protein at 600K using the scaled Hamiltonian of REST2 without attempted replica exchanges. Both $\text{RMSD} > 4\text{\AA}$ and $< 4\text{\AA}$ are sampled, by comparison in REST1 only $\text{RMSD} > 4\text{\AA}$ are sampled at high temperatures.(Fig. 4 of reference [95])

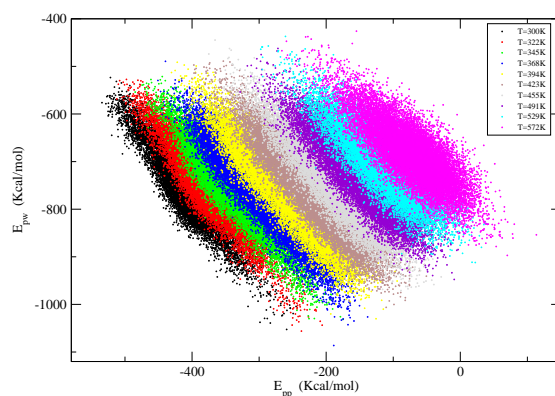


Figure 6.5: Anti-correlation between the intra-molecular potential energy of the protein and the interaction energy between the protein and water for replicas with different effective temperatures of the protein for the trpcage system.

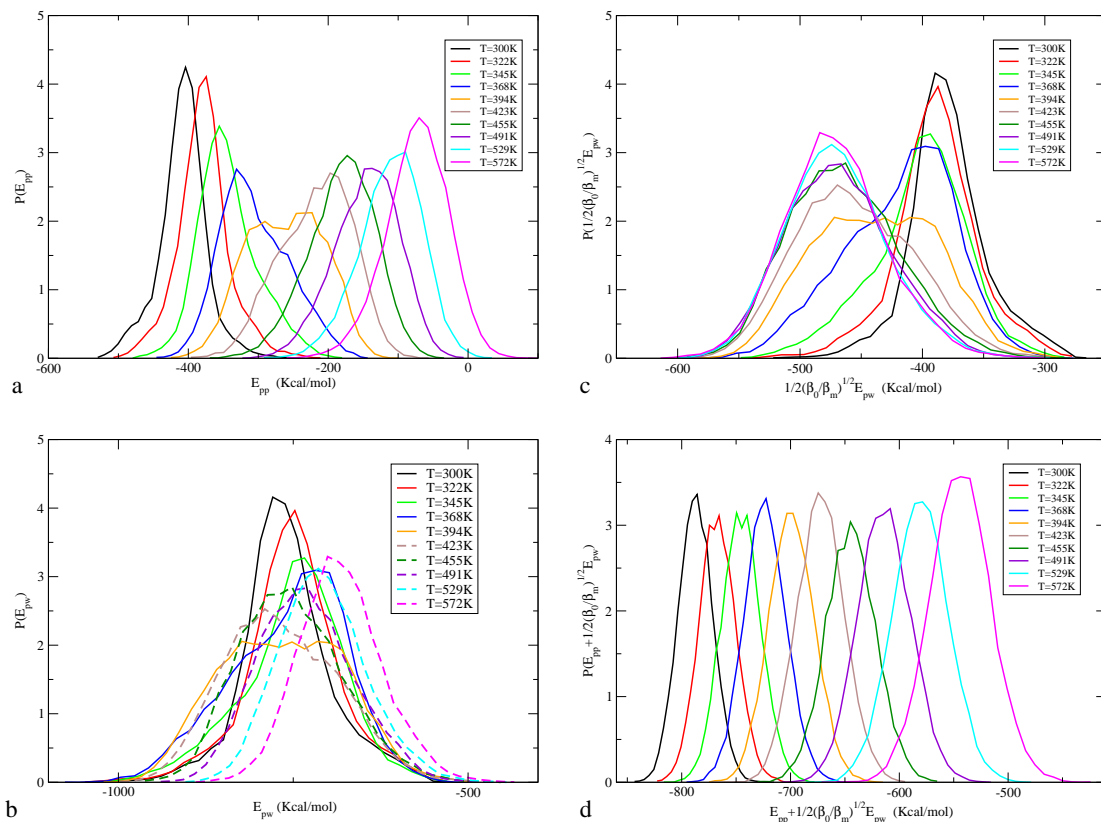


Figure 6.6: a. The distribution of intra-molecular potential energy of the protein for replicas with different effective temperatures of the protein. b. The distribution of interaction energy between protein and water for replicas with different effective temperatures of the protein. c. Distribution of $(1/2)\sqrt{\beta_0/\beta_m}E_{pw}$ for replicas with different effective temperatures of the protein. d. Distribution of $E_{pp} + (1/2)\sqrt{\beta_0/\beta_m}E_{pw}$ for replicas with different effective temperatures of the protein.

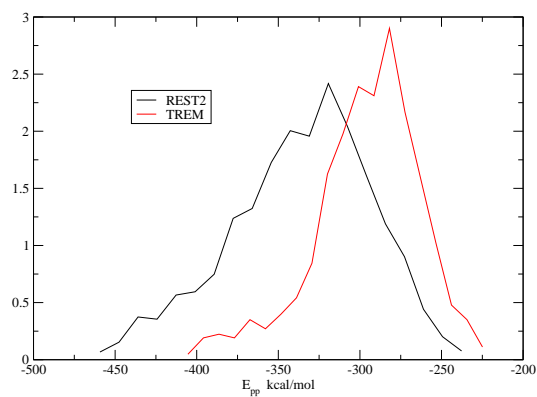


Figure 6.7: The distribution of intra-molecular potential energy of the protein at the lowest temperature replica using TREM and REST2 starting from an almost fully extended structure.

Chapter 7

On achieving high accuracy and reliability in the calculation of relative protein-ligand binding affinities

Abstract

We apply a new free energy perturbation simulation method, FEP/REST, to two modifications of protein-ligand complexes which lead to significant conformational changes, the first in the protein and the second in the ligand. The new approach is shown to facilitate sampling in these challenging cases where high free energy barriers separate the initial and final conformations, and leads to superior convergence of the free energy as demonstrated both by consistency of the results (independence from the starting conformation) and agreement with experimental binding affinity data. The second case, consisting of two neutral thrombin ligands which are taken from a recent medicinal chemistry program for this interesting pharmaceutical target, is of particular significance in that it demonstrates that good results can be obtained for large, complex ligands, as opposed to relatively simple model systems. To achieve quantitative agreement with experiment in the thrombin case, a next generation

force field, OPLS 2.0, is required, which provides superior charges and torsional parameters as compared to earlier alternatives.

7.1 Introduction

Biological processes often depend on protein-ligand binding so that accurate prediction of protein-ligand binding affinities is of central importance in structural based drug design.[49; 102; 103; 104] Among the existing methods used to calculate these binding affinities in explicit solvent, free energy perturbation (FEP) simulations provides one of the most rigorous simulation methods. Usually FEP is applied in the lead optimization stage of structure based drug design, and is used to rank-order a series of congeneric ligands in order to choose the most potent ones for further investigation.[49; 102; 103; 104]

Despite the potentially large impact that FEP could have on structure based drug design projects, practical applications in an industrial context have been limited over the past decade. High accuracy and reliability in the methodology are required to make productive decisions about compound modification during late stage lead optimization, but neither has yet been demonstrated by existing implementations. Two types of challenges stand in the way of developing FEP into a true engineering platform for drug candidate optimization. Firstly, converging explicit solvent simulations to the desired precision is far from trivial, even with the immense computing power that is currently available using low cost multiprocessor clusters or cloud computing platforms of various types. Secondly, errors in the potential energy models must be reduced to the point where they lead to errors in a converged calculation that are smaller than the desired errors in relative binding affinities compared to experiment, typically on the order of 0.5 kcal/mole. While the present article focuses primarily upon a new algorithm design to address the sampling challenge, we also provide an example, taken from the recent medicinal chemistry literature, illustrating that existing energy models, while substantially improved over the past 20 years via extensive effort in a number of research groups[105; 106; 100; 107; 108], require further refinement if the demanding target accuracy specified above is to be achieved.

FEP provides an in-principle rigorous method to calculate protein-ligand binding affinities within the limitations of the potential energy model as long as the simulation time is long enough that all the important regions in phase space are sampled. In practice, however, problems arise when there are large structural reorganizations in the protein or in the ligand upon the formation of the binding complex or upon the alchemical transformation from one ligand to another.[49; 103; 104] In these cases, there can be large energy barriers separating the different conformations and the ligand or the protein may remain kinetically trapped in the starting configuration for a very long time during brute-force FEP/MD simulations. The incomplete sampling of the configuration space results in the computed binding free energies being dependent on the starting protein or ligand configurations, thus giving rise to the well known quasi-nonergodicity problem in FEP. The slow structural reorganizations, even at a single side chain level,[58; 109] or some key solvent molecules in the binding pocket,[110; 3; 111] can affect the calculated binding affinities to a significant degree.

Recently, many groups have made efforts to reduce or eliminate the quasi-nonergodicity problem in FEP. In 2007, Mobley *et al.* proposed the “confine-and-release protocol”,[58] using umbrella sampling to calculate the potential of mean force (PMF) along the prior known slow degree of freedom. However, this method requires prior knowledge of the slow degrees of freedom, making it difficult to use for more complicated real systems. In 2010, the Roux group designed the 2-dimensional replica exchange method (REM) to compute absolute binding free energies of ligands,[109] with one REM on the Hamiltonian space for alchemical transformation, and the other REM on the sidechains surrounding the binding pocket which were assumed to include all the slow degrees of freedom without prior knowledge. However, the number of parallel replicas required in this method is very large, and it is nontrivial to apply this method to the case where the slow degrees of freedom are on the ligands.

In this article, we introduce a very efficient protocol called FEP/REST, which combines the recently developed enhanced sampling method REST (Replica Exchange with Solute Tempering)[94; 112; 96] into normal FEP to deal with the structural reorganization problem and use it to calculate relative protein-ligand binding affinities in some troublesome cases. The method assumes that the slow degrees of freedom are located

within a close neighborhood of the bound ligand without prior knowledge. The computational cost of this method is comparable with normal FEP, and it can be very easily generalized to more complicated systems of pharmaceutical interest. We apply this method on two systems; (a) the L99A mutant of the T4 Lysozyme (T4L/L99A),[113; 114] a popular model system with an engineered nonpolar binding pocket where the structural reorganization happens in the protein, and (b)Thrombin (Factor IIa),[115; 116; 117] an important drug target in the coagulation cascade where the structural reorganization happens in the ligand. (See Fig. 7.1) In both cases, the relative binding affinities calculated using FEP/REST agree with experiment within the error bars independent of starting conformation of the protein or the ligand, while normal FEP fails to characterize the effects of structural reorganization and thus gives incorrect free energies. In the latter case, we show that use of an upgraded force field model is essential in achieving the accuracy targets delineated above.

7.2 Results

Upon alchemical transformation from one ligand to another, structural reorganization might occur in the protein or in the ligand. In this article, we study two systems, the T4L/L99A and Thrombin, using both normal FEP method and the FEP/REST protocol as described in the methods section. Many aromatic molecules can bind to the nonpolar binding pocket of T4L/L99A and experimental binding affinity data are available for comparison.[114] Despite the rigidity of the protein and the simplicity of the nonpolar pocket, accurate prediction of the relative binding affinities for the ligands has proved challenging for methods ranging from rapid virtual screening and MM-GBSA to more rigorous FEP methods.[118; 119; 103; 104] The difficulty arises from the key residue Val111 surrounding the binding pocket: in the binding complex of small ligands like benzene and toluene, the Val111 stays in the “trans” conformation as in the apoprotein; in the binding complex of larger ligands like p-xylene and o-xylene, the Val111 changes its rotameric states from the “trans” conformation ($\chi \approx -180$) to the “gauche” conformation ($\chi \approx -60$),(Fig. 7.1a) which is usually called an induced fit effect.[113; 58; 109] Thrombin, a serine protease, is a very important drug target

in the coagulation cascade for many thromboembolic diseases such as deep vein thrombosis, myocardial infarction, and pulmonary embolism.[117; 115; 116] With the discovery of a neutral P1 substitute of the native substrates, a new generation of more potent inhibitors were designed with high levels of bio-availability and good pharmacokinetic properties, among which CDA and CDB are representative.[117; 115] In the binding complexes of CDA and CDB, the structures of the protein are essentially the same. However, with the addition of a methyl group on the P1 pyridine ring next to the fluorine atom, the ring flips.[117] This is shown in Fig. 7.1b where the two binding complexes are superimposed. While the fluorine atom on the P1 pyridine ring is pointing out of the S1 pocket in ligand CDA (denoted as “F-out” conformation), it is pointing into the S1 pocket in ligand CDB (denoted as “F-in” conformation). Both the reorienting of Val111 and the flipping of the pyridine ring are sufficiently slow that they are trapped in the initial conformation on the time scale of typical FEP simulation.

The estimated relative binding affinities of p-xylene with respect to benzene binding to T4L/L99A calculated using normal FEP, lambda hopping FEP (replica exchange between neighboring lambda windows),[109; 120] and FEP/REST starting from different conformations of the protein (“trans” vs. “gauche” of Val111) are given in Table 7.1. With a 2ns simulation, the normal FEP predicted relative binding affinities depend on the starting conformation and neither of them is within the error bars to the experimental result.[114] Starting from the “trans” conformation, the predicted binding affinity is more positive than experimental result (0.95 vs. 0.52 kcal/mol); starting from the “gauche” conformation, the predicted binding affinity is less positive than experimental result (0.30 vs. 0.52 kcal/mol). Using lambda hopping, the predicted binding affinities are a little closer to the experimental value than normal FEP, but a similar discrepancy as with normal FEP was found. By comparison, the estimated binding affinities determined by FEP/REST for the same 2ns simulation time are independent of the starting conformations, and are very close to the experimental result.

The side chain dihedral angle of Val111 (N-CA-CB-CG1) for the initial lambda window (binding complex of benzene) and the final lambda window (binding complex of p-xylene) as a function of simulation time starting from the “trans” conformation using normal FEP and

FEP/REST are given in Fig. 7.2. It is clear that, starting from the “trans” conformation, the Val111 was trapped in that conformation during a 2ns simulation in normal FEP. By comparison, using FEP/REST, for the same 2ns simulation time the Val111 was able to make many transitions between the different rotameric states, and while the initial state favors the “trans” conformation, the final state favors the “gauche” conformation after a short equilibration time, in agreement with experimental results.[113] Similar kinetic trapping in normal FEP and enhanced sampling in FEP/REST are observed starting from the “gauche” conformation of Val111.

We determined the probabilities for the initial and final states being in the “trans,” “gauche+,” and “gauche-” conformations, and calculated the free energy to confine the binding complex in each of these conformations using FEP/REST. For the binding complex of benzene (initial state), the probability of the “trans” conformation is 0.6, of the “gauche+” conformation is 0.4, but because the free energy of the remaining “gauche-” conformation is very high, its probability is very close to 0. For the binding complex of p-xylene (final state) the probability of the “gauche” conformation is 0.75 and of the “trans” conformation is 0.24, in agreement with previous results using umbrella sampling (0.76, 0.23, 0.002) or 2 dimensional replica exchange with a boosting potential (0.73, 0.16, 0.11).[58; 109] In normal FEP calculations, the protein was found to be “virtually” confined in the starting “trans” or “gauche” conformation, and we can correct their free energies by adding the “confine and release” free energies for each conformation according to the confine and release protocol proposed by Mobley *et al.*[58]. We thus add $(0.90+0.30-0.85=0.35$ kcal/mol) for the “trans” conformation, and $(0.30+0.54-0.17=0.67$ kcal/mol) for the “gauche” conformation, finding that the corrected results fall within the error bars of experimental value. This validates that the error of normal FEP is due to incomplete sampling of conformational space.

We used the FEP and FEP/REST protocols to calculate the relative binding affinity of ligands CDB and CDA to Thrombin, using the OPLS 2005 force field for the ligands,[100; 108] starting from different conformations of the ligand (denoted by “F-in” or “F-out” respectively). The results from 3ns simulations are given in table 7.2. The structures of the ligands are much more complicated than the T4L/L99A case, and the error bars

for these free energy results are larger. Similar to the T4L/L99A system, the calculated binding affinities using normal FEP depend on the starting conformation of the ligands, and neither of them comes close to the experimental value.[117] Using FEP/REST, the calculated binding affinities are within error bars of each other, independent of the starting conformation. From the simulated trajectories, we observed that the P1 pyridine ring was trapped in the starting conformation using normal FEP, while it flipped many times using FEP/REST, indicating the efficiency of enhanced sampling. However, none of these calculated free energies are within error bars of the experimental value.

Upon closer investigation of the FEP/REST simulated trajectories using the OPLS 2005 force field for the ligands, we found another important conformation of the ligand different from the two conformations identified in the crystal structures. The correct binding pose of ligand CDA from the crystal structure and the erroneous conformation from the simulation are given in Fig 7.3. In the correct binding pose, the P3 pyridine ring of the ligand is in the S3 pocket of the protein while in the erroneous conformation the P3 pyridine ring moves out of the S3 pocket pointing into solvent. The S3 pocket is a hydrophobic pocket and the P3 pyridine ring binds to the S3 pocket through hydrophobic interaction and edge-to-face $\sigma - \pi$ interaction between P3 aryl group and Trp215.[117] However, in the OPLS 2005 force field, there are large partial charges on the atoms of the pyridine ring (as large as -0.68 on the nitrogen atom), so the P3 pyridine ring incorrectly prefers to point into solvent. (See Table C.1 in AppendixC) In addition, the distribution of the dihedral angle involved in the flipping of P1 pyridine ring (N-C-C-C labeled in Fig. 7.1) also has an erroneous state which might be due to the incorrect dihedral angle terms in the OPLS 2005 force field (See Fig. C.2 in AppendixC). These investigations point out the deficiency of the OPLS 2005 force field and lead us to use an improved version of force field for the ligands, OPLS 2.0, which assigns the partial charges and the bonded interaction terms through high accuracy quantum mechanics calculation. The major differences between the OPLS 2005 and OPLS 2.0 force fields are the different partial charges on the atoms of the pyridine ring and the different torsional angle terms. (See detailed comparison in AppendixC)

The calculated relative binding affinities using the OPLS 2.0 force field for the ligands from normal FEP and FEP/REST are given in table 7.3. Significantly improved results are

obtained compared with those obtained from the OPLS 2005 force field. Using normal FEP, the calculated binding affinities depend on the starting conformation, with an error of about 0.6 kcal/mol compared with experimental result starting from the “F-out” conformation. By comparison, the FEP/REST predicted results are within the error bar of the experimental value independent of the starting conformation of the ligand. The dihedral angle involved in the flipping of the pyridine ring (N-C-C-C labeled in Fig. 7.1) as a function of simulation time for the initial and final states using normal FEP and FEP/REST starting from the “F-out” conformation ($\chi \approx -100$) are given in Fig. 7.4 a and 7.4 b. It is clear that the ligand was trapped in that conformation using normal FEP while it flipped between the “F-out” ($\chi \approx -100$) and “F-in” ($\chi \approx 90$) conformations many times after an initial equilibration time using FEP/REST. A similar enhanced sampling effect was observed using FEP/REST starting from the “F-in” conformation.

The flipping of the pyridine ring in the Thrombin system occurs more slowly than the transitions between rotameric states in the T4L/L99A system, and more intermediate lambda windows were needed to help converge its free energy, thus it takes a much longer time to equilibrate the two “F-in” and “F-out” conformations. To shorten the simulation time to get close to equilibrium, we performed two additional FEP/REST simulations; (a) one with the first half of the lambda windows starting from “F-in” conformation and the last half of the lambda windows starting from “F-out” conformation (denoted as “F-in/out” in table 7.3), and (b) the other with an inverted starting conformation for each lambda window (denoted as “F-out/in”). The calculated relative binding affinities from these two simulations (table 7.3) are within the error bar of the experimental result independent of whether starting conformations (a) or (b) are used. The dihedral angle involved in the flipping of the pyridine ring is given as a function of simulation time for the initial and final lambda windows in Fig. 7.4c. Indeed, the time required to get close to equilibrium was much shorter than what was found from a single conformation for each lambda window and, importantly, higher precision results were obtained. Thus when the binding poses for the two ligands are known *a priori*, it will be more efficient to start the FEP/REST simulation with each lambda window starting from different conformations.

It should be pointed out that the final equilibrium distribution and the free energy are

independent of the starting conformation for each replica as long as there are a sufficient number of conformational transitions in the middle lambda window in FEP/REST. Using different starting conformations for different replicas, as opposed to the same starting conformations, can shorten the simulation time for getting close to equilibrium within the same error bars. This is because the time scale for a transition from one conformation to another in MD is much longer than the time scale for the exchange of two conformations between neighboring replicas. We note that the time required to truly equilibrate is the same for any starting configuration except if we started with the equilibrium distribution. The fact that the calculated free energies using different replica starting conformations (“F-in”, “F-out”, “F-in/out”, “F-out/in”) are within the error bars of each other indicates that a 3ns simulation time is sufficiently long enough to equilibrate the generalized ensemble in this case.¹

In the FEP/REST simulations, we also calculated the probabilities for the initial and final states being in the two conformations (“F-in” vs “F-out”) which is displayed in Fig. 7.4d. For the final state (binding complex of CDB), the “F-in” conformation is the major conformation, in agreement with the experimental crystal structure; however, for the initial state (binding complex of CDA), the two conformations have almost equal probability in contrast to the experimental crystal structure where it was found to be in the “F-out” conformation. This discrepancy might be due to the different physical conditions in experiment (crystal) and in simulation (in solution). To confirm this argument, we performed another two FEP/REST simulations with the protein heavy atoms harmonically restrained to the initial position (corresponding to the crystal structure) starting from different ligand conformations for each lambda window. The trajectories from these simulations confirm that the “F-out” conformation is a major conformation for the initial state and the “F-in” conformation is a major conformation for the final state when the protein heavy atoms are restrained (see Fig. 7.5), validating the hypothesis that the solution environment may shift

¹ For example, in two state kinetics, the deviation of the concentration of reactant (or product) from its equilibrium concentration decays as $\delta c(t) = \delta c(0) \exp(-t/\tau)$, where τ is the relaxation time. Thus all choices of the initial deviation decay on the same time scale, but the smaller $\delta c(0)$ is, the less time it will take to reach $\delta c = 0$ within the specified error bar.

the relative population of the two conformations from what is found in the solid. Interestingly, the calculated relative binding affinities from the two simulations with protein heavy atoms restrained (Table 7.3) converge to the same value but different from experimental result by about 0.8 kcal/mol, indicating a 0.8 kcal/mol difference in protein restrain free energy for the two binding complexes.

7.3 Discussions and Conclusions

The results reported comprise only a few test cases. However, the performance of the algorithm is encouraging with regard to overcoming problems due to significant configurational changes, in either the protein or ligand, upon ligand modification. The REST methodology in both examples facilitates rapid interconversion between the phase space region separated by barriers in normal FEP, at a relatively low computational cost and without the requirement of prior knowledge of the slow degrees of freedom. If these properties are shown to hold for a larger, diverse set of test cases, this will represent a significant advance in the convergence of FEP simulations. Other groups have succeeded in the T4L/L99A case, but as pointed out above, at a substantially higher computational cost. The thrombin example is no longer a toy problem, but represents the sort of modification made on a routine basis on complex ligands in late stage drug discovery projects. The striking success of FEP/REST in this case offers hope that it will be applicable, in its current form, to real world problems as well as model systems. The efficient sampling of FEP/REST allows us to observe that one is free to play with the Hamiltonian of intermediate states in FEP as long as the correct physical states are achieved at the end-points. It is also worth noting that both the 2-dimensional replica exchange method[109] and FEP/REST in their current forms enhance the sampling only of the localized region around the ligands, which might not be sufficient for treating delocalized conformational changes (allosteric regulation). A possible procedure to treat this problem is the following: (1) include a larger “hot” region in a first round FEP/REST simulation, and find those key residues responsible for the allosteric regulation; (2) run a second round FEP/REST just including those key residues in the “hot” region.

Three other points are worthy of discussion. Firstly, the improved results obtained with OPLS 2.0, as opposed to OPLS 2005, were achieved without any specific parameter adjustment based on the experimental FEP data. Rather, much more extensive fitting to basic quantum mechanical data for charges and torsional parameters yields a superior force field which can be expected to display similarly enhanced results for other ligands and receptors. Preliminary results in which statistical measures of errors in conformational energies for OPLS 2.0 as compared not only to OPLS 2005, but also to alternative force fields like MMFF, support this suggestion.[121] For ligands relevant to medicinal chemistry efforts, it is likely that improvements in both sampling and the potential energy model are needed to approach agreement with high quality experimental data on a routine basis.

Secondly, our results suggest that FEP/REST methods can be substantially improved in efficiency and reliability if the endpoints of the calculation (i.e, the co-crystallized structures that would be obtained experimentally for the two ligands) are known. Often, one endpoint is available from experiment (the lead compound which is being modified in the lead optimization process). The other endpoint can then be generated via conformational search calculations using induced fit docking (IFD) algorithms[122], which are typically much less expensive than the FEP simulation itself. In some challenging cases the IFD calculation will generate a small number (typically 2-3) alternatives for the endpoint; here, FEP/REST can be used to select between these alternatives with improved accuracy, while at the same time using the truncated list of alternatives to reduce FEP/REST calculation time, and focus FEP/REST sampling on relevant phase space regions.

Finally, the differences between results obtained with crystal packing as compared to free solution, to our knowledge the first to rigorously explore this issue, are of significant interest, although a large data set will have to be investigated to draw firm conclusions. One would not expect crystal packing to lead to very large changes in structure or binding affinity in an active site cavity (which typically is recessed and hence has few direct contacts with neighboring protein molecules of the crystal) except in unusual cases, and our results are consistent with this intuition. However, a nontrivial effect, big enough to be relevant to the potency targets in drug discovery projects, is observable, and the better agreement of the solution calculation with experiment (performed in solution) confirms that

the computational estimation of the effect is likely to be a good estimate.

7.4 Methods

7.4.1 FEP/REST

The incomplete sampling of configurational space in normal FEP results from the large energy barriers separating the relevant conformational states. Our strategy to solve the quasi-nonergodicity problem is to combine enhanced sampling techniques into FEP.

Recently, our group proposed Replica Exchange with Solute Tempering (REST) in which, through Hamiltonian scaling, only a small region of interest of the system is effectively “heated up” while the rest of the system stays “cold.”[94] In this way, a small number of replicas are sufficient to maintain the sampling efficiency, in contrast to the large number of replicas needed in the usual temperature replica exchange. Here, we are using a more recently developed version of REST (called REST2) where the effective temperature of the hot region is achieved at the Hamiltonian level through scaling the potential energy terms of the hot region.[112]

In FEP/REST, along the alchemical transformation from the initial lambda window to the final lambda window, the effective temperature of the “hot region” (the region we are interested in, usually including the ligand and the protein residues surrounding the binding pocket) is gradually increased from T_0 for the initial lambda window to T_h for the middle lambda window, and then gradually decreased from T_h for the middle lambda window back to T_0 for the final lambda window. The effective temperature of the hot region is achieved by scaling the Hamiltonian, and exchange of configurations between neighboring lambda windows is attempted using the Hamiltonian Replica Exchange Method (HREM).[98] (All of the replicas are run at the same temperature, and the velocities and kinetic energies of all of the atoms whose interactions are scaled remain always in contact with a single heat bath at this same temperature.) In this way, enhanced sampling is achieved through the increased effective temperature of the hot region at intermediate lambda windows,[112] and through replica exchange the initial and final lambda windows can sample the different conformations. The effective temperature for the initial and final states is at T_0 , which

is the temperature we are interested in, and the sum of free energy difference between all neighboring lambda windows gives the relative binding affinity between the two ligands. The intermediate accessory states not only help to bridge the different phase space regions for the initial and final states as in normal FEP, but also helps the sampling of different conformational states through the increased temperature of the hot region. This method does not require prior knowledge of the slow degrees of freedom and can be easily applied to complicated real systems of medicinal interest.

7.4.2 Details of the simulations

In FEP/REST, in addition to the different energy terms introduced in alchemical transformation in normal FEP, the different effective temperatures of the “hot region” in REST will make the free energy difference between neighboring lambda windows larger, and the precision of the free energy results might be reduced. This is the price paid to get enhanced sampling. The larger the hot region, and the higher the effective temperature of the hot region, the stronger the enhanced sampling effect, but the error bars in the resulting calculated free energy energies between neighboring lambda windows are also increased. So a proper choice of the hot region and effective temperature profile reflects a trade off between the precision of free energy results and the efficiency of the enhanced sampling; consequently the “hot region” should be as small as possible but still be able to sample structural reorganization effects. In the two systems studied in this article, we know the slow degrees of freedom, so only the residue Val111 or the P1 pyridine ring was included in the hot region. In general, if there is no prior knowledge about the slow degrees of freedom, a proper choice of hot region would include the ligand and the protein residues surrounding the ligand because usually the structural reorganization involves the ligand and the protein residues surrounding the binding pocket.

The free energy difference, ΔF , between neighboring lambda windows depends on the distribution functions $P_0(\Delta E)$ and $P_1(\Delta E)$ of energy differences (ΔE) in forward and backward sampling respectively,[123] through,

$$P_1(\Delta E) = P_0(\Delta E) \exp^{-\beta(\Delta E - \Delta F)}. \quad (7.1)$$

The two distributions are equal for the specific energy difference $\Delta E = \Delta F$, and the

accuracy of the free energy ΔF depends on the overlap of the two distributions.[123] At the same time, it is easy to show by imposing detailed balance condition that the acceptance ratio for attempted replica exchanges between neighboring lambda windows also depends on the energy difference from forward and backward sampling:[95]

$$\Delta_{01} = \beta(E_1(X_0) + E_0(X_1) - E_1(X_1) - E_0(X_0)) \quad (7.2)$$

$$= \beta(\Delta E(X_0) - \Delta E(X_1)). \quad (7.3)$$

Here, X_0 and X_1 are the configurations sampled in the forward and backward directions, and E_0 and E_1 are the potential energy functions for the two states. So both the accuracy of the free energy result and the efficiency of the enhanced sampling are maximized when neighboring lambda windows have regions of overlap in potential energy distribution, and an optimal alchemical lambda schedule and effective temperature schedule will generate equal acceptance ratios for all neighboring lambda windows.

All simulations were done using the Desmond program.[31] The starting structures for the simulations were taken from crystal structures with PDBIDs 181L (Val111 “trans”) and 187L (Val111 “gauche”) for T4L/L99A,[113] and with PDBIDs 1MU6 (“F-out” conformation) and 1MU8 (“F-in” conformation) for Thrombin.[117] The structures of the proteins were modified using protein preparation wizard [86] and the protonation states were assigned assuming the systems are at pH 7.0. The OPLS 2005 force field[100; 108] was used for the protein and the Tip4p water model[32] was used for the solvent. Both the OPLS 2005 and the OPLS 2.0 force fields were used for ligands CDA and CDB, and the OPLS 2005 force field was used for benzene and p-xylene. Simulations lasted for 2ns for the T4L/L99A complexes, 3ns for the Thrombin complexes and 5ns for ligands in pure solvent.

A dual topology ideal gas molecule end state method was used to define the mutation path, which facilitates the sampling through the double tunneling mechanism.[120] The electrostatic interactions unique to the initial ligand were turned off before the Lennard Jones (LJ) interactions, and the LJ interactions unique to the final ligand were turned on followed by the electrostatic interactions. The core of the LJ interactions is made softer to avoid the singularities and instabilities in the simulation.[67] The mutation path is symmetric, so mutation from either direction will give identical free energy result. To get more efficient

enhanced sampling, the fluorine atom on the P1 pyridine ring was mutated to an identical atom, so that the effective volume of P1 pyridine ring was made smaller in the middle lambda window and the transition between the two conformations was faster. The lambda values, the scaling factors for the “hot” region, and the free energy difference between all neighboring lambda windows for the two systems are given in Table. 7.4 and 7.5.

In the two systems we studied, with a total number of 16 lambda windows and highest effective temperature of 1200K for the T4L/L99A system, and a total number of 23 lambda windows and highest effective temperature of 1784K for the Thrombin system, and a time interval of 1ps between attempted exchanges among neighboring replicas, we obtained an average acceptance ratio of 0.54 for T4L/L99A system and 0.59 for the Thrombin system. The energy difference between neighboring λ windows for each configuration was calculated, and only data generated after the equilibration stage were used to calculate the free energy through the Bennett acceptance ratio method[30]. The error was calculated using block averages with a 500ps bin width.

The bonded interactions involving the dummy atoms are treated differently in this article to avoid singularities and instabilities, with the details given in the Appendix D. This problem is not appreciated in the literature on FEP.

Table 7.1: Predicted relative binding affinities of p-xylene to T4L/L99A compared with benzene using various methods

Starting conformation	method	ΔG in complex	$\Delta\Delta G$
trans	FEP	-3.31 ± 0.10	0.95 ± 0.15
	λ -hopping	-3.36 ± 0.10	0.90 ± 0.15
	FEP/REST	-3.78 ± 0.10	0.48 ± 0.15
gauche	FEP	-3.96 ± 0.10	0.30 ± 0.15
	λ -hopping	-3.83 ± 0.10	0.43 ± 0.15
	FEP/REST	-3.77 ± 0.1	0.49 ± 0.15
exp			0.52 ± 0.09

Free energies in kcal/mol; ΔG in solvent is -4.26 ± 0.05 kcal/mol.

Table 7.2: Predicted relative binding affinities of ligand CDB to Thrombin compared with ligand CDA using OPLS 2005 force field for the ligands

Starting conformation	method	ΔG in complex	$\Delta\Delta G$
F-out	FEP	2.04 ± 0.20	-0.14 ± 0.30
	FEP/REST	0.50 ± 0.20	-1.68 ± 0.30
F-in	FEP	0.32 ± 0.20	-1.86 ± 0.30
	FEP/REST	0.70 ± 0.20	-1.48 ± 0.30
exp			-0.85

Free energies in kcal/mol; ΔG in solvent is 2.18 ± 0.10 kcal/mol. “F-in”/“F-out” means the fluorine atoms on the P1 pyridine ring pointing into or out of the P1 pocket of Thrombin.

Table 7.3: Predicted relative binding affinities of ligand CDB to Thrombin compared with ligand CDA using OPLS 2.0 force field for the ligands

Starting conformation	method	ΔG in complex	$\Delta\Delta G$
F-out	FEP	1.09 ± 0.20	-0.21 ± 0.30
	FEP/REST	0.18 ± 0.20	-1.12 ± 0.30
F-in	FEP	0.07 ± 0.20	-1.23 ± 0.30
	FEP/REST	0.27 ± 0.20	-1.03 ± 0.30
F-in/out	FEP/REST	0.30 ± 0.15	-1.00 ± 0.25
F-out/in	FEP/REST	0.52 ± 0.15	-0.78 ± 0.25
F-in/out	FEP/REST(res)	1.22 ± 0.10	-0.08 ± 0.20
F-out/in	FEP/REST(res)	1.44 ± 0.10	0.14 ± 0.20
exp			-0.85

Free energies in kcal/mol; ΔG in solvent is 1.30 ± 0.10 kcal/mol. “F-in/out” means the first half lambda windows start from the conformation with the fluorine atoms on the P1 pyridine ring pointing into the P1 pocket of Thrombin and the last half lambda windows start from the conformation with the fluorine atoms on the P1 pyridine ring pointing out of the P1 pocket. The reversed starting conformations were used for “F-out/in.” “FEP/REST(res)” means FEP/REST simulation with the protein heavy atoms harmonically restrained to the initial position (corresponding to the crystal structure).

Table 7.4: Lambda values, scaling factors and free energy difference between neighboring lambda windows for T4L/L99A system

λ	0	1	2	3	4	5	6	7
bondedA	1.0	0.933	0.867	0.8	0.733	0.667	0.6	0.533
bondedB	0.0	0.067	0.133	0.2	0.267	0.333	0.4	0.467
chargeA	1.0	0.75	0.5	0.25	0.0	0.0	0.0	0.0
chargeB	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
vdwA	1.0	1.0	1.0	1.0	1.0	0.857	0.714	0.571
vdwB	0.0	0.0	0.0	0.0	0.0	0.143	0.286	0.429
scaling	1.00	0.8464	0.7056	0.5776	0.4624	0.3721	0.3025	0.25
ΔG_{FEP}	-2.4997	-2.5577	-2.6721	-2.7345	-0.7043	0.3626	0.8722	0.2075
$\Delta G_{FEP/REST}$	1.6669	1.6307	1.5967	1.5675	3.0860	3.6672	3.4928	-0.0463
λ	8	9	10	11	12	13	14	15
bondedA	0.467	0.4	0.333	0.267	0.2	0.133	0.067	0.0
bondedB	0.533	0.6	0.667	0.733	0.8	0.867	0.933	1.0
chargeA	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
chargeB	0.0	0.0	0.0	0.0	0.25	0.5	0.75	1.0
vdwA	0.429	0.286	0.143	0.0	0.0	0.0	0.0	0.0
vdwB	0.571	0.714	0.857	1.0	1.0	1.0	1.0	1.0
scaling	0.25	0.3025	0.3721	0.4624	0.5776	0.7056	0.8464	1.00
ΔG_{FEP}	-0.5328	-0.9243	-1.1963	2.4196	2.3307	2.2174	2.0848	
$\Delta G_{FEP/REST}$	-3.3633	-4.2287	-4.9902	-1.9003	-1.9433	-1.9780	-2.0394	

Table 7.5: Lambda values, scaling factors and free energy difference between neighbor
 lambda windows for Thrombin system

λ	0	1	2	3	4	5	6	7	8	9	10	11
bondedA	1.00	0.95	0.91	0.86	0.82	0.77	0.73	0.68	0.63	0.59	0.54	0.50
bondedB	0.00	0.05	0.09	0.14	0.18	0.23	0.27	0.32	0.37	0.41	0.46	0.50
chargeA	1.00	0.75	0.50	0.25	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
chargeB	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
vdwA	1.00	1.00	1.00	1.00	1.00	0.68	0.46	0.33	0.25	0.19	0.12	0.00
vdwB	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
scaling	1.00	0.9216	0.8464	0.7569	0.6724	0.5776	0.49	0.4096	0.3364	0.2704	0.2116	0.1681
ΔG_{FEP}	2.5092	2.2140	1.9203	1.6569	0.1517	-0.3940	-0.3947	-0.2190	-0.1103	-0.0743	0.0912	-0.3519
$\Delta G_{FEP/REST}$	1.3099	0.7003	-1.1062	-0.5334	-1.7920	-1.9750	-1.9259	-1.7625	-1.6065	-1.4625	-1.0676	0.9690
λ	12	13	14	15	16	17	18	19	20	21	22	
bondedA	0.46	0.41	0.37	0.32	0.27	0.23	0.18	0.14	0.09	0.05	0.00	
bondedB	0.54	0.59	0.63	0.68	0.73	0.77	0.82	0.86	0.91	0.95	1.00	
chargeA	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	
chargeB	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.25	0.50	0.75	1.00	
vdwA	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	
vdwB	0.12	0.19	0.25	0.33	0.46	0.68	1.00	1.00	1.00	1.00	1.00	
scaling	0.2116	0.2704	0.3364	0.4096	0.49	0.5776	0.6724	0.7569	0.8464	0.9216	1.00	
ΔG_{FEP}	-0.0413	0.0475	0.0895	0.1379	-0.0907	-1.2740	-0.0423	-0.8294	-1.5706	-2.4305		
$\Delta G_{FEP/REST}$	1.4394	1.6101	1.7749	1.8948	1.6985	0.8331	1.5105	0.8022	-0.3918	-1.8040		

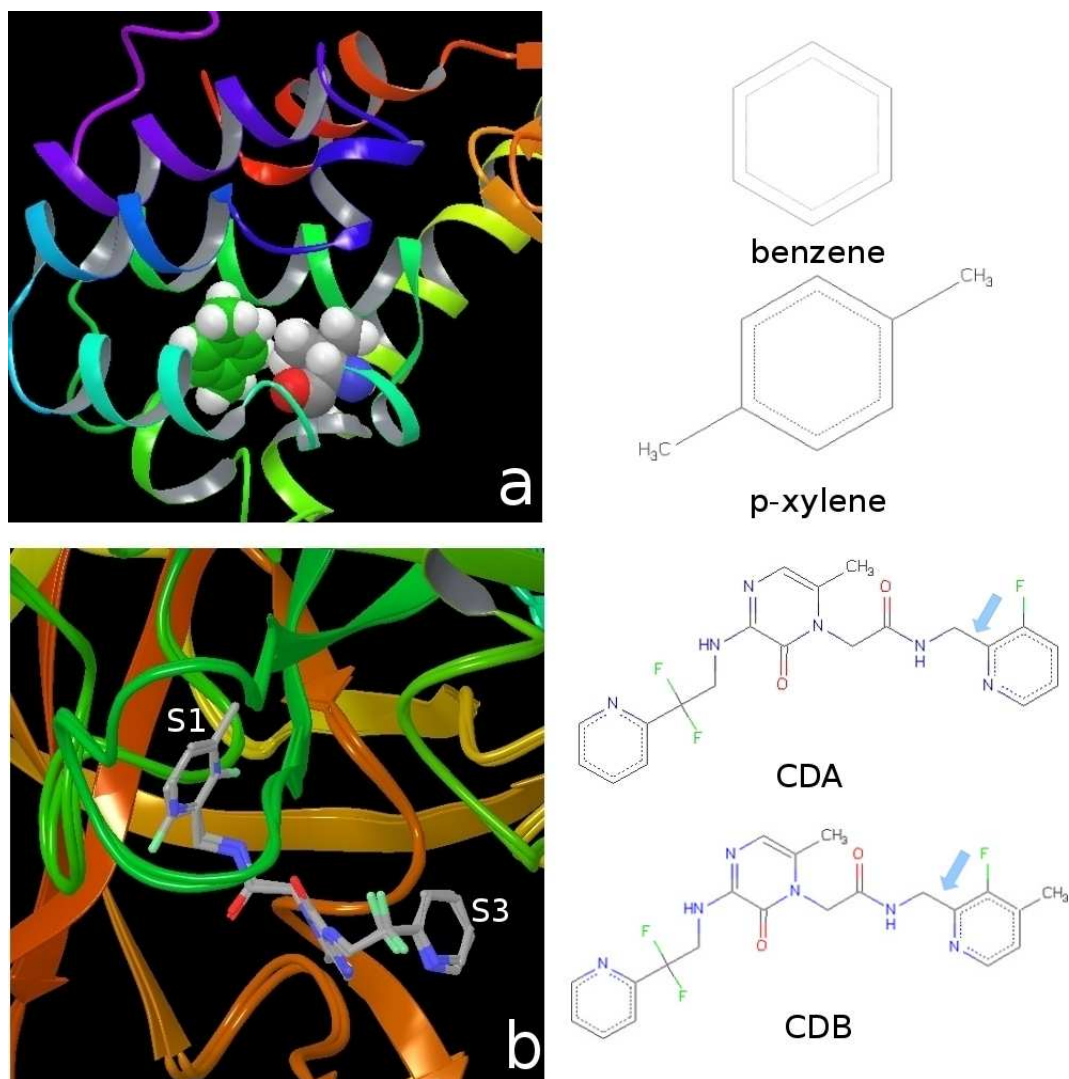


Figure 7.1: a. The nonpolar binding pocket of T4L/L99A with p-xylene bound. The key residue Val111 and p-xylene are displayed in VDW mode. The structures of the two ligands, benzene and p-xylene, for the relative binding affinity calculation are given on the right. b. The binding pocket of Thrombin with the ligands CDA and CDB superimposed. With the addition of the methyl group on the P1 pyridine ring of ligand CDB, the ring flips. In the binding complex of Thrombin/CDA, the Fluorine atom on the P1 pyridine points out of the S1 pocket (“F-out” conformation), while the fluorine atom points into the S1 pocket (“F-in” conformation) in the Thrombin/CDB binding complex. The structures of the two ligands CDA and CDB for relative binding affinity calculation are given on the right with the dihedral involved in the flipping of the P1 pyridine ring (N-C-C-C) indicated by an arrow.

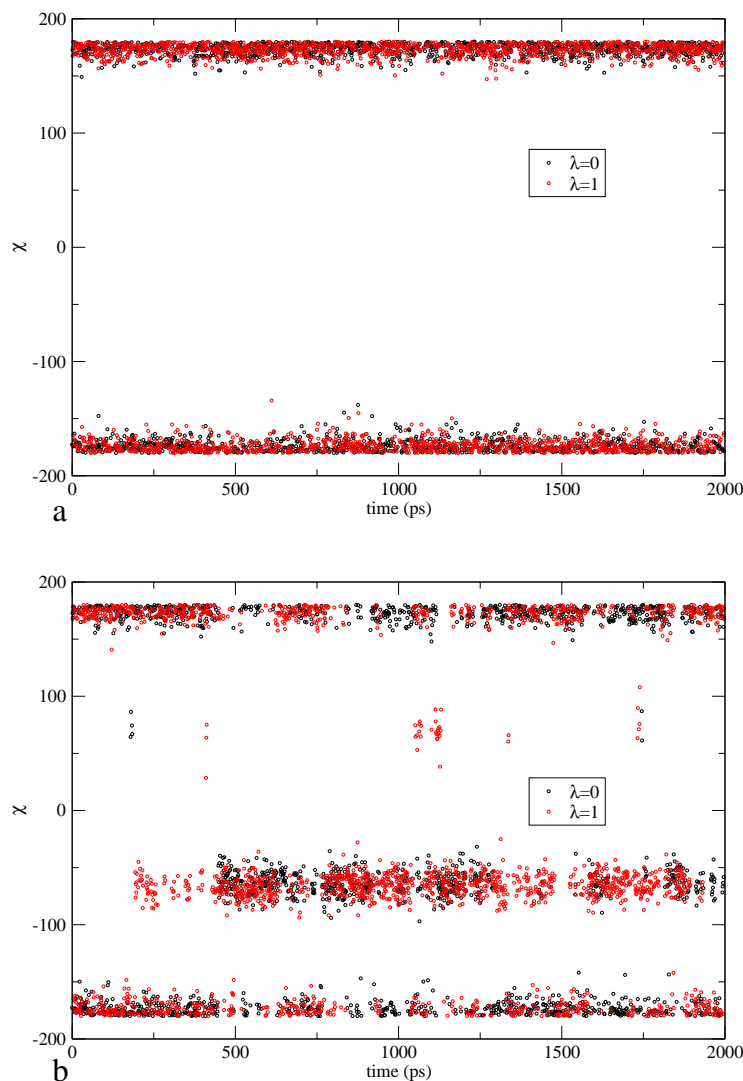


Figure 7.2: The Val111 side chain dihedral angle (N-CA-CB-CG1) as a function of simulation time for the initial and final lambda windows. Initial lambda window corresponds to the T4L/L99A/benzene binding complex, and the final lambda window corresponds to the T4L/L99A/p-xylene binding complex. a. Results from normal FEP simulation starting from the “trans” conformation. The Val111 was trapped in the “trans” conformation through the 2ns simulation time. b. Results from FEP/REST simulation starting from the “trans” conformation. After a short equilibration time, the Val111 transits between the “trans” and “gauche” conformation with a dominating “gauche” conformation for the final state and a dominating “trans” conformation for the initial state, in agreement with experiment. Similar enhanced sampling was observed using FEP/REST starting from the “gauche” conformation.

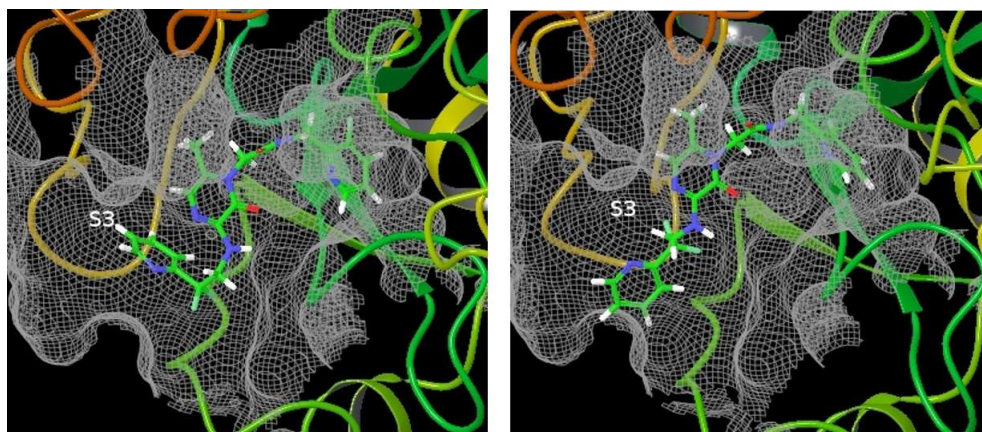


Figure 7.3: The correct binding pose from the crystal structure (left) and the erroneous conformation (right) observed in simulation using the OPLS 2005 force field for the ligands. In the correct binding pose, the P3 pyridine ring points into the S3 pocket of Thrombin while the P3 pyridine ring moves out of the S3 pocket and points into solvent in the erroneous conformation.

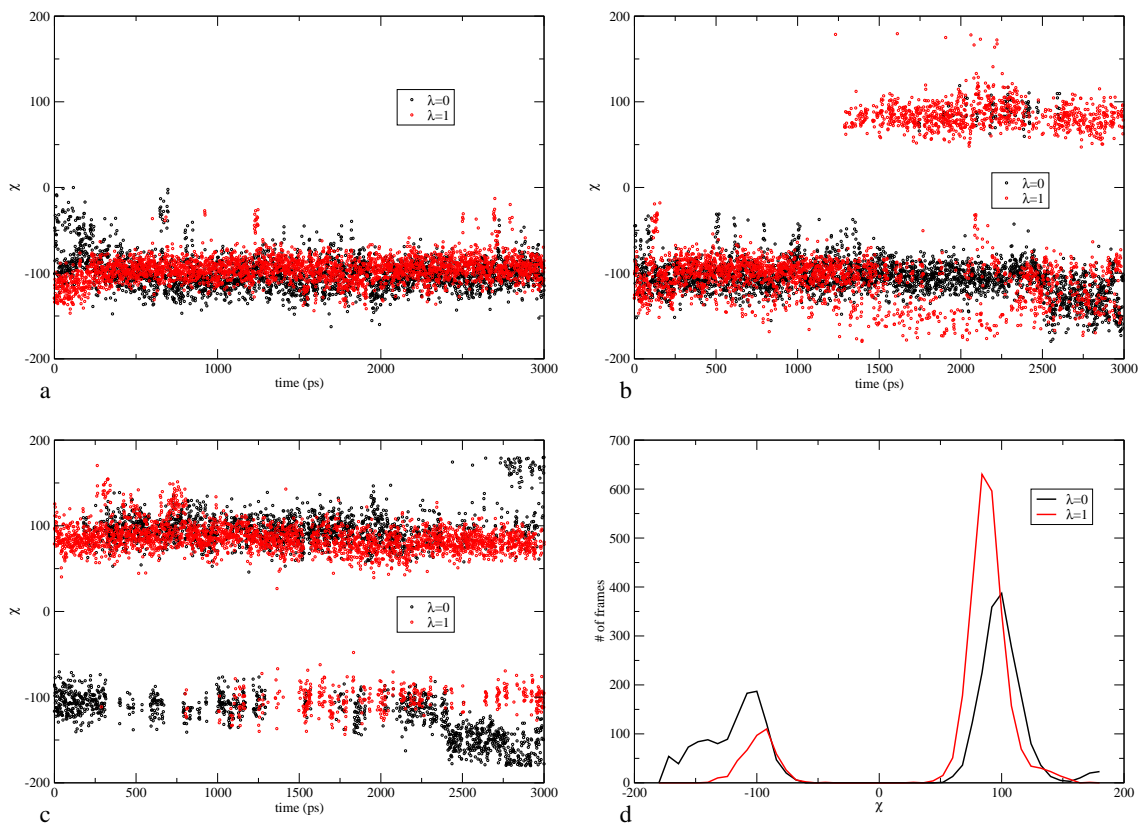


Figure 7.4: The distribution of the dihedral angle involved in the flipping of P1 pyridine ring (N-C-C-C labeled in Fig. 7.1) for the initial and final lambda windows using OPLS 2.0 force field for the ligands. a. The dihedral angle as a function of simulation time using normal FEP starting from the “F-out” conformation ($\chi \approx -100$). The ligands were trapped in that conformation through the 3ns simulation time. b. The dihedral angle as a function of simulation time using FEP/REST starting from the “F-out” conformation ($\chi \approx -100$). After the equilibration stage, the pyridine ring transits between the “F-in” ($\chi \approx 90$) and “F-out” ($\chi \approx -100$) conformations. c. The dihedral angle as a function of simulation time using FEP/REST with the first half lambda windows starting from “F-out” conformation and the last half lambda windows starting from “F-in” conformation. The equilibration time was much shorter compared with b. d. The distribution of the two conformations for the initial and final states. The binding complex of Thrombin/CDB ($\lambda = 1$) favors the “F-in” ($\chi \approx 90$) conformation in agreement with crystal structure, while the binding complex of Thrombin/CDA ($\lambda = 0$) has almost equal probability for the two conformations. This slight discrepancy with experimental crystal structure might be due to the different physical conditions in simulation and in experiment (in solution vs. in crystal).

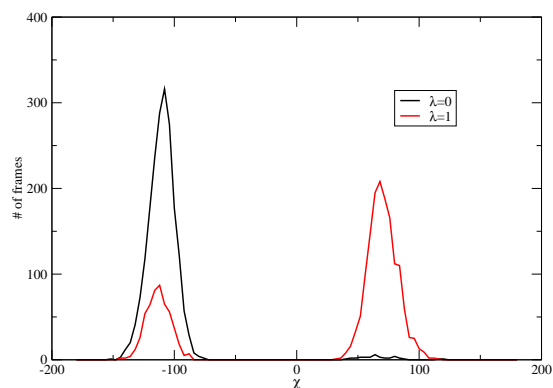


Figure 7.5: The distribution of the dihedral involved in the flipping of P1 pyridine ring from a FEP/REST simulation with the protein heavy atoms harmonically restrained to the initial position.

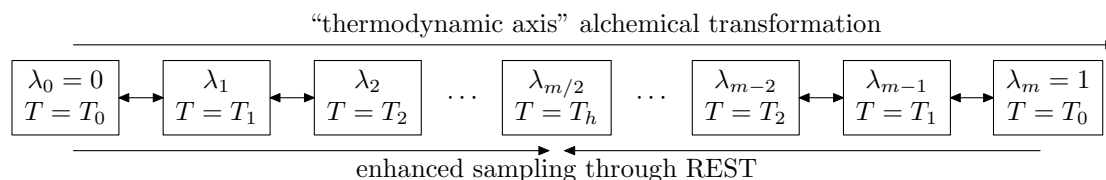


Figure 7.6: 1-dimensional replica exchange protocol combining REST into FEP. Each box represents a lambda window with the input parameters given by λ , the thermodynamic coupling parameter, and T , the effective temperature of the hot region. The double arrow symbols indicate attempts to exchange configurations between neighboring replicas.

Part III

Investigations about hydrophobic interaction and electrostatic interaction

Chapter 8

Introduction of hydrophobic interactions and electrostatic interactions

In the previous two sections, I have introduced several techniques to calculate protein-ligand binding affinities. In the protein-ligand binding process, usually there is no covalent bond forming or breaking, and all that is involved is the hydrophobic interaction and electrostatic interaction. In this section, I will present some investigations towards the foundational understanding of hydrophobic interaction and electrostatic interaction. To be specific, the nonadditivity effect in hydrophobic interactions and the competition of hydrophobic interaction and electrostatic interaction between a hydrophobic particle and model enclosures are discussed.

In the WaterMap method, when a ligand displaces two hydration site water molecules, the binding affinity contribution of displacing the two water molecules is assumed to be the sum of the contribution from displacing each water molecule separately. In other words, we assume that the effect of displacing multiple water molecules is pairwise additive. Similarly, in the cavity contribution term, if several atoms of the ligand are located in the dry region of the binding pocket, we also assume the effect is pairwise additive. In general, pairwise additivity assumes that the potential of the mean force (PMF) holding a cluster of N particles

together is equal to the sum of pairwise interaction free energies (or pair pmfs). In Chapter 9, The binding affinities between a united-atom methane and various model hydrophobic enclosures were studied through high accuracy free energy perturbation methods (FEP) and the nonadditivity of the hydrophobic interaction in these systems, measured by the deviation of its binding affinity from that predicted by the pairwise additivity approximation, is investigated. Many of the implicit solvent models attempt to account for hydrophobicity in terms of nonpolar surface exposure to water, among which solvent accessible surface area (SASA) and molecular-surface-area (or Connolly surface area, MSA) models are the most popular. We investigated how well these implicit solvent models can characterize the nonadditivity effect and found that implicit solvent models based on the molecular surface area (MSA) performed much better, not only in predicting binding affinities, but also in predicting the non-additivity effects, compared with models based on the solvent accessible surface area (SASA), suggesting that MSA is a better descriptor of the curvature of the solutes.

In popular implicit solvent models, like MMPBSA or MMGBSA, the solvent is treated as a dielectric media, and the free energy to turn on the charge on solute atoms is quadratically dependent on the magnitude of the charge of the solute. In Chapter 10, we studied the binding between a united atom methane and model hydrophobic plates with different charge densities and different charge patterns on the plates. From this study, we observed that the binding affinity is reduced when the plates are charged, and with increased charge density, the plates can change from “hydrophobic like” (pulling the particle into the interplate region) to “hydrophilic like” (ejecting the particle out of the interplate region), demonstrating the competition between hydrophobic and electrostatic interactions. In addition, the electrostatic contribution to the binding affinity is quadratically dependent on the magnitude of the charge for symmetric systems, but linear and cubic terms also make a contribution for asymmetric systems. We explain these results by statistical perturbation theory and show when and why implicit solvent models fail.

Chapter 9

Hydrophobic interactions in model enclosures from small to large length scales: nonadditivity in explicit and implicit solvent models.

Abstract

The binding affinities between a united-atom methane and various model hydrophobic enclosures were studied through high accuracy free energy perturbation methods (FEP). We investigated the non-additivity of the hydrophobic interaction in these systems, measured by the deviation of its binding affinity from that predicted by the pairwise additivity approximation. While only small non-additivity effects were previously reported in the interactions in methane trimers, we found large cooperative effects (as large as $-1.14 \text{ kcal mol}^{-1}$ or approximately a 25% increase in the binding affinity) and anti-cooperative effects (as large as $0.45 \text{ kcal mol}^{-1}$) for these model enclosed systems. Decomposition of the total potential of mean force (PMF) into increasing orders of multi-body interactions indicates that the

contributions of the higher order multi-body interactions can be either positive or negative in different systems, and increasing the order of multi-body interactions considered did not necessarily improve the accuracy. A general correlation between the sign of the non-additivity effect and the curvature of the solute molecular surface was observed. We found that implicit solvent models based on the molecular surface area (MSA) performed much better, not only in predicting binding affinities, but also in predicting the non-additivity effects, compared with models based on the solvent accessible surface area (SASA), suggesting that MSA is a better descriptor of the curvature of the solutes. We also show how the non-additivity contribution changes as the hydrophobicity of the plate is decreased from the dewetting regime to the wetting regime.

9.1 Introduction

The hydrophobic interaction (HI) plays a very important role in the formation and stability of many self-assembled aggregates and biological structures,[1] and is considered to be the driving force for protein folding.[124; 125] A fundamental understanding of HI is crucial to the study of many important biological phenomena, such as protein folding, micelle formation, protein-ligand binding.

An important question concerning the protein folding problem is whether hydrophobic associations are pairwise additive, cooperative, or anti-cooperative.[126; 127] In other words, is the potential of mean force (PMF) holding a cluster of n hydrophobic particles together equal to the sum of pairwise interaction free energies (or pair pmfs), or is it more negative (cooperative) or more positive (anti-cooperative) than what is predicted by the pairwise additivity approximation.

Nemethy and Scheraga[128] found that the hydrophobic interactions between more than two solute particles can not be expressed as a sum of pairwise solute-solute interactions. Palma also observed the non-additivity of solvent induced potential of mean force for hydrophobic particle solutions by molecular dynamics simulation.[129] In 1997, Rank and Baker[130] studied the three methane molecules in an isosceles triangle geometry with two methane molecules at contact distance forming a fixed base, and they found that the three-

body PMF was anti-cooperative for distance up to 6.5 Å. The conclusions of this work were not reliable, however, because the error associated with the baseline of the PMF was of the same order as the non-additive term itself. After that, Chan [131; 126] and Scheraga[132; 133; 134] performed studies of methane trimers using the weighted histogram analysis method (WHAM) and the test particle insertion method respectively, and found contradictory results. After exchange of comments between these two groups,[135; 136; 137; 138] they found that the disagreement came from the assignment of the baseline for the PMFs. Comparing these two methods Scheraga[139] stated that methane dimer seemed to be the largest system that can be treated by the test particle insertion method. The effect of size, pressure, temperature, and salts on the non-additivity of the three particle system was also investigated.[140; 141; 142; 143] Recently, cooperative effects on the association of four methane molecules were also investigated by Scheraga's group.[144]

In this paper, we used the FEP method to study the binding affinities between a united-atom methane and various model hydrophobic enclosures. The binding affinities were compared with the predictions of the pairwise additivity approximation to access the non-additive contributions to the hydrophobic interactions in these systems. Two different empirical models were also used to predict the binding affinity and non-additivity effect, and comparisons were made with the FEP reference data. We found that the molecular surface area (MSA, or Connolly surface area) model performed much better than solvent accessible surface area (SASA) model both for the prediction of binding affinities and the nonadditivity effects, which is consistent with previous findings[130; 132; 133; 126]. Detailed analysis of these two models indicates that there is problem intrinsic to the SASA model, which can not predict the cooperative effects. Decomposition of the total PMF into increasing orders of multi-body interactions indicated that the higher order multi-body interactions can be either positive or negative, and increasing the order of the multi-body interactions considered did not necessarily improve the accuracy.

9.2 Definition of cooperative and anti-cooperative effects

In general, for a solution with n solute particles forming a cluster, the potential of mean force (PMF), $W(1, 2, \dots, n)$, is related to the n -particle correlation function $g^{(n)}(1, 2, \dots, n)$ by the following definition:[145]

$$g^{(n)}(1, 2, \dots, n) = e^{-\beta W(1, 2, \dots, n)}, \quad (9.1)$$

where $\beta^{-1} = k_B T$, k_B is Boltzmann's constant, and T is the temperature. $W(1, 2, \dots, n)$ is a short-hand notation for $W(r_1, r_2, \dots, r_n)$, where r_i denotes the position of the i -th particle, and it corresponds to the free energy to bring the n particles from infinitely far apart to the current configuration.

The n -particle PMF can be decomposed into single-body, pairwise and multi-body contributions:

$$\begin{aligned} W(1, 2, \dots, n) &= F(1, 2, \dots, n) - \sum_{i=1}^n F(i) & (9.2) \\ &= \sum_{i < j} \delta F(i, j) + \sum_{i < j < k} \delta F(i, j, k) + \dots + \delta F(1, 2, \dots, n) \\ &= W_2 + W_3 + \dots + W_n & (9.3) \end{aligned}$$

Where $F(1, 2, \dots, n)$ is the hydration free energy for the specified configuration of the solute particles, $F(i)$ is the hydration free energy of solute article i in an infinitely dilute solution, $\delta F(i, j)$ is the same as the normalized two-body PMF $W(i, j)$, and $\delta F(i, j, \dots)$ corresponds to subsequent higher order multi-body interactions. W_2 is the sum of pairwise contributions, and W_m is the sum of m -body interactions. Truncating the series to n -body term leads to the Generalized Kirkwood Superposition Approximation (GKSA) to the n -th order.[19; 20] Specifically, the hypothesized pairwise additivity of PMF is obtained by truncating the series at the pairwise term. In clusters this approximation neglects shielding effects where a third particle could shield a pair from other particles and from the solvent.

For the three-particle case, pairwise additivity is equivalent to the Kirkwood superposition approximation,[146; 147; 148] and $\delta F(i, j, k)$ measures the non-additive part of the three body interactions. Cooperativity is defined when $\delta F(i, j, k)$ is negative, meaning the free energy between the third particle and the remaining two particles is more negative

and the configuration is more favorable than pairwise additivity predicts, and similarly, anti-cooperativity is defined when $\delta F(i, j, k)$ is positive, meaning the free energy between the third particle and the remaining two particles is more positive and the configuration is less favorable than pairwise additivity predicts. For the n -particle case, if the interaction between a specific particle and the remaining $n-1$ particles is equal to the sum of the pairwise interactions between that specific particle with each of the other $n-1$ particles, and if this condition holds for each of the n particles, the total n -body PMF will be equal to the sum of pairwise interactions. So we further generalize the cooperative and anti-cooperative concepts to the n -particle case: when the interaction energy between one specific particle and the remaining $n - 1$ particles is more negative than the sum of pairwise free energies between the specific particle and each particle in the remaining $n - 1$ cluster, we term it “cooperative”; “anti-cooperative” is similarly defined. In other words, if we label the specific particle as n , and introduce the notation $\delta W(1, 2, \dots, n - 1; n)$, which is the sum of all higher than two body interactions involving particle n ,

$$\begin{aligned} \delta W(1, 2, \dots, n - 1; n) &= F(1, 2, \dots, n) - F(1, 2, \dots, n - 1) - F(n) \\ &\quad - \sum_{i=1}^{n-1} \delta F(i, n) \end{aligned} \tag{9.4}$$

$$\begin{aligned} &= \sum_{i < j \leq n-1} \delta F(i, j, n) + \sum_{i < j < k \leq n-1} \delta F(i, j, k, n) \\ &\quad + \dots + \delta F(1, 2, \dots, n - 1, n) \end{aligned} \tag{9.5}$$

then, cooperativity or anti-cooperativity is defined when $\delta W(1, 2, \dots, n - 1; n)$ is negative or positive, respectively. In Eq. (9.4), the term, $F(1, 2, \dots, n) - F(1, 2, \dots, n - 1) - F(n)$, gives the interaction energy between particle n and the remaining $n - 1$ particles, and the term, $\sum_{i=1}^{n-1} \delta F(i, n)$, gives the sum of pairwise interactions between particle n and each of the other $n - 1$ particles. The difference between these two terms provides a measure of the non-additive contribution to the interaction.

9.3 Simulation details

In this paper, molecular dynamics simulations were performed using the DESMOND program[31] to study the binding affinities between a united-atom methane and 13 model hydrophobic

enclosures depicted in figure 9.1. The geometry of the model hydrophobic plate in these systems is displayed in the bottom right of figure 9.1. It consists of 19 single-layer atoms arranged in a triangular lattice with a bond length of 3.2 Å. For systems consisting of two plates, the two plates were parallel and in-registry with separation distance of $D = 7.46\text{Å}$. The LJ atoms forming the enclosures were uniformly represented with Lennard Jones parameters $\sigma = 3.73\text{ Å}$ and $\epsilon = 0.294\text{ kcal/mol}$, which are the same as the united-atom methane parameters used in these simulations.[149] The inserted methane particle (displayed in green in figure 9.1) was placed at the contact distance ($d = 3.73\text{ Å}$) with the other atoms and plate(s) forming the enclosures. In order to study the non-additivity of hydrophobic interactions between the insertion methane and the enclosures, another 10 systems corresponding to the subsystems of the 13 model enclosures (depicted in figure 9.2, and named after the corresponding system in figure 9.1 by adding a “prime”) were also studied. The binding affinities for the other four systems in figure 9.2 [G'', H'', I'', K''] were calculated through a thermodynamic cycle by combining the binding affinities calculated for related systems. For example, the binding affinity for system G'' was calculated by combination of binding affinities calculated for systems G', G and E.

The free energy perturbation (FEP) method was used to determine the binding affinities between the inserted methane and each of the enclosures. The Maestro System Builder utility [65] was used to insert each enclosure into a cubic water box with a 10 Å buffer. The SPC water model[38] was used to describe the solvent. The atoms of the enclosures were constrained to their initial positions throughout the dynamics, and only the solvent degrees of freedom were sampled. The united-atom methane was “turned on” inside the model enclosures over 9 lambda windows with $\lambda=[0, 0.125, 0.25, 0.375, 0.50, 0.625, 0.75, 0.875, 1]$, where λ is the coupling parameter to turn on/off the LJ interaction between the methane and the rest of the system with initial state and final state correspond to $\lambda = 0$ and $\lambda = 1$ respectively. In these simulations, the core of the LJ potential is made softer[67] as $\lambda \rightarrow 0$ to avoid singularities and numerical instabilities. For each of the λ windows, molecular dynamics simulations were performed. The energy of the system was minimized, and then equilibrated to 298 K and 1 atm with Nose-Hoover[33; 34] temperature and Martyna-Tobias-Klein[35] pressure controls over 100 ps of molecular

dynamics. A cutoff distance of 9 Å was used to model the Lennard Jones interactions, and the particle-mesh Ewald method[37] was used to model the electrostatic interactions. Following the equilibration, a 20 ns production molecular dynamics simulation was performed and configurations of the system were collected every 1.002 ps. The energy difference between neighboring λ windows for each configuration saved was calculated and the Bennett acceptance ratio method[30] was used to calculate the free energy difference between neighboring states. The sum of the free energy difference between neighboring states gave the solvation free energy of methane in the enclosures. The same procedure was followed to calculate the solvation free energy of methane in bulk water. The difference between the two solvation free energies gave the binding free energy to bring a methane from infinitely far to inside the hydrophobic enclosure, which is the potential of mean force (PMF) between the methane and the enclosure. The error associated with these binding affinities was of the order of $\pm 0.02 \text{ kcal.mol}^{-1}$.

9.3.1 Implicit solvent model calculations

Due to the large computational cost of running explicit solvent model simulations, in protein folding or protein-ligand binding problems, one is often forced to use implicit solvent models to reduce the computational cost. Many of the implicit solvent models attempt to account for hydrophobicity in terms of nonpolar surface exposure to water,[150] among which solvent accessible surface area(SASA)[151] and molecular-surface-area (or Connolly surface area, MSA)[152] models are the most popular. The solvent accessible surface (SAS) is traced out by the probe sphere center as it rolls over the solute, and the molecular surface (MS) is the surface traced by the inward-facing surface of the probe sphere.

In this paper, the SASA and MSA of each enclosure, both with and without the bound methane, were computed with the Connolly molecular surface package,[68] as was the SASA and MSA of the methane particle by itself. From this data the buried surface area upon methane-enclosure complexation was determined. The direct Lennard Jones interaction energy upon the binding of methane to each enclosure was similarly computed. The buried surface area times the surface tension is often used to approximate the solvent induced potential of mean force. Together with the direct Lennard Jones interaction energy, the

total binding affinity between the methane and each enclosure can then be calculated, as has been done in many empirical methods for calculating binding affinities.

In addition, to investigate whether these two implicit solvent models can predict the non-additivity of the hydrophobic effect, the predicted pairwise additive buried surface area upon methane binding to each enclosure was also calculated. The deviation of the actual buried surface area from the pairwise additivity predicted allowed us to estimate the non-additive part of hydrophobic interactions based on these models.

9.4 Results and discussion

The binding free energies between methane and the model hydrophobic enclosures, as measured by FEP, are reported in table 9.1. It is found that the range of binding free energies of the methane for the model enclosures is nearly 5 kcal.mol^{-1} . Also reported in table 9.1 are the pairwise additivity predicted binding affinity upon complexation, the buried surface area upon complexation, (both SASA and MSA), the pairwise additivity predicted buried surface area upon complexation, and the deviation between corresponding terms which gives the non-additive contributions.

9.4.1 Comparison for different implicit solvent models in predicting the binding affinity

From the data presented in table 9.1, we can determine how well the buried surface area/molecular mechanics model predicts the binding affinity. Tuning the surface tension coefficient to minimize the mean-average-error (MAE) of fit with FEP reference data, we obtained an optimal surface tension coefficient of $\gamma = 0.00763 \text{ kcal mol}^{-1}\text{\AA}^{-2}$ for the SASA and $\gamma = 0.03767 \text{ kcal mol}^{-1}\text{\AA}^{-2}$ for the MSA models of these enclosures. For comparison, the surface tension for SASA determined by fitting to experimental solvation free energy for linear or branched alkanes by Honig[153] was $0.005 \text{ kcal mol}^{-1}\text{\AA}^{-2}$, whereas the macroscopic water-alkanes surface tension was $0.070 \text{ kcal mol}^{-1} \text{\AA}^{-2}$. Since there are no overlaps between the inserted methane and the enclosures, the buried van der Waals area is zero for all these systems. So the van der Waals surface model would just predict the binding affinity to be

the direct LJ interaction energy. The predicted binding affinities versus the FEP reference data were reported in figure 9.3. From this figure we see that, for the most part, the MSA model performed better than the SASA model. This is indicated by a higher R^2 value (0.89 vs 0.76) and smaller MAE (0.40 vs 0.57), which is consistent with previous findings.[130; 126; 132; 133] Both of these models performed much better than the van der Waals surface area model, which has $R^2 = 0.70$ and MAE= 0.94. The SASA based model cannot differentiate the hydrophobicity between systems [J, K, K', K'', L, L', M, and M'] (predicting a similar binding affinities of about -4.1 kcal/mol) nor between systems [D, E, F, F', G, G', G'', H, H', H'', I, I', and I''] (predicting a similar binding affinities of about -2.2 kcal/mol), while the MSA model performed much better for these systems and predicted the right order of hydrophobicity among these systems to some extent.

The SASA model found enclosures J and K' in which a methane molecule is bound between two hydrophobic plates to be the most hydrophobic of all the systems [J, K, K', K'', L, L', M, M']. The buried SASAs for these two systems, upon methane complexation, were the largest because large swaths of formerly accessible surface area on the faces of the plates are buried by the presence of the binding methane. For the other enclosures, several methane molecules already lie between the plates in the absence of the binding methane so that part of the surface area buried by the binding methane was already buried by the other particles. For the extreme cases of systems L and M, which are most hydrophobic, the SASA model strongly underestimates the binding affinity, because the buried SASAs are smaller for these two systems than in systems J and K'. On the other hand, the MSA model to some extent predicts the right order of binding affinity among these systems. A similar analysis of systems [D, E, F, F', G, G', G'', H, H', H'', I, I', I''], yields similar conclusions.

9.4.2 Comparison of MSA and SASA model predictions of the non-additivity effect

The non-additive contributions to the binding affinities of methane to all of the enclosures are listed in table 9.1. We see from the table that while only small non-additive effects (± 0.2 kcal/mol) were observed for the methane trimers by the Chan and Scheraga groups,[144; 134; 126; 131] large cooperative effects (as large as -1.14 kcal/mol for system L') and

anti-cooperative effects (as large as 0.45 kcal/mol for system J) were observed for these systems.

The data shown in table 9.1 allows us to phenomenologically identify a connection between the sign of the non-additivity effect and the curvature of molecular surface. To wit, for all the enclosures that exhibit an anti-cooperative effect, the molecular surfaces of the enclosures were convex without any concave or saddle parts, whereas for all enclosures that exhibit a cooperative effect, the molecular surfaces had concave or saddle parts, and the methane binds close to the saddle or concave part of the surface.

Consider systems E and G for example. The molecular surface of enclosure E has no concave part, and the binding affinity of methane is anti-cooperative, (less favored than pairwise additivity predicted). In the presence of another methane in enclosure G, the molecular surface has a saddle part, and when the methane binds to this saddle part, it shows a cooperative effect. Similar analyses can be done on systems J and K, K' and M', K' and K''. Systems I', H', and F have the same enclosure whose molecular surface exhibits a saddle region, but the inserted methane binds at different locations of the enclosure. Comparing the methane binding affinities for these systems, we found that as the methane moves closer and closer to the saddle part (from systems I' to H' to F), the binding affinity gets more and more favorable (from -1.74 kcal/mol to -1.97 kcal/mol to -2.63 kcal/mol), and the cooperative effect changes from relatively weak (-0.05 kcal/mol for system I') to strong (-0.51 kcal/mol and -0.37 kcal/mol for systems H' and F, respectively). Similar analysis can be done for systems M', L', and K.

The cooperative effects in systems H' and L' were found to be stronger than for systems F and K respectively. This might seem to be at odds with intuition because systems F and K are more compact than systems H' and L' respectively. This is not surprising though: the distance between the two methane molecules in systems H' and L' corresponds to the de-solvation barrier on the PMF curve of two methanes molecules in bulk water. In the presence of the plate(s), the de-solvation barrier between the two methanes may not exist, so that the binding affinities might be much more favorable than would be predicted by the pairwise additivity approximation.

In addition to the data for the non-additive contribution to the changes of SASA and

MSA, we can investigate whether either of these implicit solvent models can predict the non-additivity in the hydrophobic interactions. Multiplying the non-additive part of the changes in MSA and SASA in methane complexation by the surface tension coefficient obtained in the previous section, we can determine the MSA and SASA predicted non-additive contributions to the hydrophobic interactions. Figure 9.4 depicts the relationship between the SASA/MSA predicted non-additive contributions and the FEP results of the non-additive hydrophobic effects. It can be seen that the MSA predicted results have a strong correlation with the FEP results while SASA predicted results anti-correlate with the FEP results.

At first glance at figure 9.4, it might seem that MSA anti-correlates with SASA, in contradiction to our common understanding of these models. To better understand this “strange” behavior, the non-additive contribution of SASA and MSA for a model methane trimer system was investigated. With the two methanes kept at their contact distance, the position of the third methane can be specified by two coordinates (θ, d) . (See the top right corner of figure 9.5). Figure 9.5 shows the predictions of the non-additive part of hydrophobic interactions given by calculating the changes in the MSA and the SASA each multiplied respectively by their corresponding surface tension coefficients (as in the previous section) as a function of the distance d when $\theta = 0$. Actually, this figure is representative of what we found for all of the angles θ . (See figure 9 of ref. [126]) We see from this figure that for $d = 3.23 \text{ \AA}$, which corresponds to the configuration where the third methane is in contact with the other two methanes, the predictions of MSA and SASA have opposite signs. For all of the systems studied in this paper, the binding methane is in contact with the other methane(s) and/or plate(s) in the enclosures, thus it is not surprising to find the anti-correlation between MSA and SASA predictions observed in figure 9.4.

In addition, it can be seen from figure 9.5 that while the MSA model can predict additive, anti-cooperative, and cooperative effects, the SASA model only predict additive and anti-cooperative effects. This conclusion seems to be generally valid. In the SASA model, the surface considered is formed from overlapping spheres, each of whose radii is the sum of the corresponding atomic van der Waals radii plus the radius of water.[151] The buried surface area between a sphere and a cluster of spheres, when the cluster of spheres has no overlaps,

is always equal to the sum of buried surface areas between the sphere and each individual sphere in the cluster, and smaller than this when the cluster of spheres has overlaps. This means that the SASA model can never predict cooperative effects. However, in the MSA model, the surface is traced by the inward-facing surface of the probe sphere,[152] and it has the potential to predict all of the non-additivity effects. This partially explains why the MSA based model performs better than the SASA based model both for the binding affinity and for the non-additivity effects.(See Figure 9.6)

It should be noted from figure 9.4 that while the MSA model successfully predicts how the non-additivity effects correlate with the FEP results, there were a few outliers. We will interpret these outliers case by case.

In system M', the region between the plates is enclosed on four different sides by hydrophobic moieties. This causes the density of water in the enclosed region to be much smaller than in bulk (one fourth of bulk value), close to a hydrophobic dewetting condition. Part of the time there was one water molecule in the enclosed region and part of time it was empty, indicating an interface different from the normal nonpolar solute/liquid water interface. It is well known that in hydrophobic dewetting there is a strong driving force to bring the hydrophobic particles together,[1] which explains the strong cooperative effect observed in system M'.

For systems H' and L', as previously mentioned, the distance between the two methane molecules corresponds to the de-solvation barrier on the PMF curve of two methane molecules in bulk water, and in the presence of the plate(s), the de-solvation barrier may not exist at all. Because of the barriers in the pair potential of mean force, the pairwise additivity approximation predicts that this configuration will be very unfavorable, but the full calculation shows that in fact they are not unfavorable. In fact, the binding affinities in these particular configurations are much more favorable than would be predicted by assuming pairwise additivity. This may be the reason for the strong cooperative effects observed for these systems. In addition, the dewetting argument discussed in connection with systems M' also applies to system L', which further validates the strong cooperative effect observed for system L'.

To investigate the reason for the strong anti-cooperative effects observed in systems J

and K' , we analyzed the structure of water between the two plates and found that, at the surface of one hydrophobic plate (system D), water breaks one hydrogen bond on average, which caused the average interaction energy between water at the surface and the rest of the system to be higher than that in bulk water by 1.12 kcal/mol. However, between two hydrophobic plates (system J), water breaks less than two hydrogen bonds on average because of its flexibility in making hydrogen bonds. This is supported by the fact that the average interaction energy between water molecules located between the two plates and the rest of the system is higher than in bulk water by only 2.05 kcal/mol, less than twice 1.12 kcal/mol ($2 \times 1.12 = 2.24$) expected from doubling the effect of one plate. It is well known that the hydrophobic effect for large scale systems is enthalpy driven,[1; 154] because of broken of hydrogen bonds at the surface of hydrophobic plates. A large contribution of the methane-enclosure binding affinity comes from de-solvation of solvent between the plates.[155] The anti-cooperative effect observed for systems J and K' may be due to the the fact that the excess of the interaction energy of water located between two hydrophobic plates (system J) over that of bulk water is found to be smaller than twice the value of the water at one hydrophobic plate (system D).

9.4.3 Non-additivity effect at wetting-dewetting transition

In the previous sections, we have shown how the non-additivity effects of methane binding affinities in enclosures with different topologies correlate with the MSA measurements. It is well known that the hydrophobicity of the enclosures depends not only on the topologies but also on the LJ parameters for atoms making up the enclosures.[1] So it will be interesting to study how the non-additivity effects depend on the LJ parameters for particles making up the enclosure. In this section, we will explore this effect by changing the LJ ϵ parameter for particles making up the plates for one representative enclosure, enclosure J.

Figure 9.7 depicts how the binding affinities of methane in enclosure J (ΔG_J) and in enclosure D (ΔG_D) changes as a function of the LJ ϵ parameter for particles making up the plate(s) from FEP simulations. The pairwise additivity predicted binding affinities for enclosure J, which is two times that for enclosure D, are also depicted in the figure. As can be seen from this figure, while the binding affinities for enclosure D decreases (or alternatively

the free energy becomes more positive) monotonically with increasing value of ϵ , the binding affinities for enclosure J first increases (the free energy becomes more negative) and then decreases, and the non-additivity effect goes from slightly anti-cooperative at very low ϵ region (smaller than 0.06 kcal/mol), to cooperative at intermediate ϵ region (between 0.06 and 0.23 kcal/mol), back to anti-cooperative at high ϵ region (higher than 0.23 kcal/mol).

The solvation free energy of methane in the enclosure can be decomposed into two components: the free energy to create a cavity with the size of methane in the enclosure and the free energy to turn on the attractive part of interactions between methane and the rest of the system. For enclosure D, when increasing the ϵ parameter, the free energy to create the cavity will become more unfavorable because the solvent will become denser at the surface of the plate; however, the free energy to turn on the attractive part of interactions between methane and the plate will become more favorable with increasing value of ϵ . These two factors having opposite effects, but the first component dominates, so the overall binding affinity will decrease slightly. For enclosure J, at low value of ϵ , the region between the plates dewets, so the free energy to create the cavity is almost zero and changes slightly with increasing value of ϵ , but the free energy to turn on the attractive interactions between the methane and the plates becomes more negative, so the binding affinity will increase in this dewetting region. The critical value of ϵ corresponding to the wetting-dewetting transition is $\epsilon \approx 0.15\text{kcal/mol}$. At this point, the probability for observing a dry inter-plate region is 50%. For ϵ larger than this, the free energy to create the cavity grows rapidly with increasing value of ϵ , becoming the dominant effect, so that the overall binding affinity decreases rapidly with increasing value of ϵ . At the critical value of ϵ , that is at the wetting-dewetting transition, there is a large cooperative non-additive effect on the binding of the methane between the plates. With increasing values of ϵ , the two plates affect the density fluctuation of solvent by more than twice what one plate does, and the slope of the binding affinity versus ϵ is therefore much larger for enclosure J than that for enclosure D. For sufficiently large ϵ there is a large anti-cooperative deviation from additivity. At $\epsilon = 0.23$, these two effects balance each other, so the free energy is additive at this point. The ϵ value for system J studied in the previous section is $\epsilon = 0.294\text{kcal/mol}$, which is higher than 0.23, so we observed an anti-cooperative effect there. As ϵ is decreased

below the critical value of ϵ , the binding affinity increases for enclosure D and decreases for enclosure J and becomes anti-cooperative in this very low ϵ region.

9.4.4 Higher order multi-body interactions

In the previous section, we investigated the non-additive effects manifested when methane is inserted into different model enclosures. For systems consisting of three components, the non-additive effect corresponds to three-body interactions; for systems consisting of more than three components, the non-additive effect is the summation of all higher-than two-body interactions involving the insertion methane, corresponding to $\delta W(1, 2, \dots, n-1; n)$, defined in Eq.(9.5). We now investigate the contributions beyond three body interactions defined recursively in Eq.(9.3).

Table 9.2 lists the total PMF, the two-body contribution to the PMF, W_2 , the three-body contribution to the PMF, W_3 , and the subsequent higher order contributions for all the systems consisting of more than three components. The deviations found by truncating the series up to order n , ΔW_n , equivalent to the GKSA approximation to n -th order, is also listed in the table. (For example, $\Delta W_2 = W(1, 2, \dots, n) - W_2$.) From this table, we see that the error arising from truncation of the total PMF at the pairwise term can be as large as 1.43 kcal/mol (system L), indicating the importance of non-additivity effects in these systems. In addition, the higher order multi-body contributions can be either positive or negative, suggesting that hydrophobic interactions can be quite complex.

Figure 9.8 and 9.9 gives the deviations of the PMF predicted by the GKSA from the total PMF, when the latter is truncated to order n , as a function of n . For most of the cases, the higher the order of multi-body interactions considered, the smaller the magnitude of the error (systems G, H, I, L', L, M). However, this is not generally true since for systems K and M', inclusion of the three-body interactions increases the error.

9.5 Conclusion

In this Chapter, the binding affinities between a united-atom methane and various model hydrophobic enclosures were determined using high accuracy FEP molecular dynamics sim-

ulations. Comparisons were made between the binding affinities from FEP and predictions based on assuming pairwise additivity, and through this it was possible to investigate the non-additive contributions of the hydrophobic interactions. Small non-additivity effects were found in the methane trimer systems by the Chan and Scheraga groups,[144; 134; 126; 131] but we find large cooperative effects (as large as -1.14 kcal/mol) and anti-cooperative effects (as large as 0.45 kcal/mol) for our relatively larger systems. Although approximations based on pairwise additivity of PMF have been used to study the transition state of protein folding[156] and the force-extended behavior of protein,[157] simulations done in this Chapter indicate that higher order correlations may be very important in real biological systems such as protein folding, protein-ligand binding and other relevant fields. This should not be surprising since the Kirkwood superposition approximation fails in dense simple fluids.

Phenomenologically, the sign of the non-additive contributions to the binding affinities of methane in the enclosures was found to be correlated with the curvature of the enclosures. To be specific, anti-cooperative effects were observed only in enclosures whose molecular surfaces were convex without any concave or saddle part, and cooperative effect were observed in enclosures whose molecular surfaces having concave or saddle parts, in which case the methane was found to bind close to the saddle or concave part of the surface. Such observations might be useful for further development of models to incorporate the non-additivity effect.

We also investigated whether two kinds of implicit solvent models are consistent with the observed binding affinities and non-additivity effects. In these models the solvent induced free energy for particle insertion was computed from the product of a “surface tension” and the area of the buried surface upon methane complexation. The area of the buried surface was computed by using the solvent accessible surface area (SASA) and by using the molecular surface area (MSA) as described in the text. We found that the MSA based model performed much better than the SASA based model in predicting the binding affinities, an observation consistent with previous findings.[130; 126; 132; 133] In addition, the SASA based model always predicts non-additive effects which anti-correlate with the FEP reference data, whereas the MSA based model performed reasonably well in predicting the non-

additive effects, except for a few outliers which were explained in the text. Further analysis indicated that the MSA based model predicts all cooperative and anti-cooperative non-additivity effects, and the SASA based model exhibits an intrinsic problem in that it can never predict cooperative effects. Furthermore, because of the correlation we observed between the non-additive contributions and the curvature of molecular surface in the MSA based model, we believe that the MSA based model is a far better descriptor of the curvature of simple solutes than the SASA based model.

The non-additivity effect depends not only on the topology but also on the LJ parameters for atoms making up the enclosure. By changing the LJ ϵ parameter for atoms making up the plates for enclosure J, we observed a wetting-dewetting transition in the inter-plate region, and the non-additivity effect changes from anti-cooperative in the low ϵ region, to cooperative in the intermediate ϵ region, and then to anti-cooperative again in the higher ϵ region. This complicated re-entrant behavior of the non-additivity effect results from the competition between the two factors contributing to the solvation free energy: the free energy to create the cavity and the the free energy to turn on the attractive interactions between the solute and the rest of the system. While the first factor dominates in the higher ϵ region (wetting region), the second factor dominates in the lower ϵ region (dewetting region).

The decomposition of the PMF for cluster formation can be expressed as a sum of increasing orders of multi-body interactions. We found that the multi-body correlations can be either negative or positive, implying that the hydrophobic interaction will depend on the topology of the surfaces enclosing the particle. In addition, increasing the order of multi-body interactions included can usually improve the accuracy, but this was not generally true.

Table 9.1: The binding thermodynamics of methane for the various model hydrophobic enclosures.

	ΔG_{FEP}	ΔG_2	$\Delta\Delta G$	$\Delta SASA$	$\Delta SASA_2$	$\Delta\Delta SASA$	ΔMSA	ΔMSA_2	$\Delta\Delta MSA$	sign
A	-0.60	-	-	-57.44	-	-	-4.15	-	-	-
B'	-0.06	-	-	0.00	-	-	0.00	-	-	-
C'	0.18	-	-	-25.75	-	-	3.54	-	-	-
D	-1.66	-	-	-89.09	-	-	-14.44	-	-	-
E'	-0.06	-	-	0.00	-	-	0.00	-	-	-
F'	-1.64	-	-	-88.54	-	-	-14.16	-	-	-
B	-1.15	-1.21	0.06	-114.88	-114.88	0.00	-8.29	-8.29	0.00	anti
E	-2.17	-2.26	0.09	-146.53	-146.53	0.00	-18.58	-18.58	0.00	anti
G'	-1.44	-1.46	0.02	-114.29	-114.29	0.00	-10.61	-10.61	0.00	anti
J	-2.86	-3.31	0.45	-178.17	-178.17	0.00	-28.87	-28.87	0.00	anti
K'	-2.83	-3.27	0.44	-177.08	-177.08	0.00	-28.31	-28.31	0.00	anti
C	-1.41	-1.21	-0.20	-95.61	-114.88	19.27	-11.46	-8.29	-3.16	coop
F	-2.63	-2.26	-0.37	-115.03	-146.53	31.50	-25.14	-18.58	-6.56	coop
G	-3.41	-2.86	-0.54	-153.20	-203.97	50.77	-24.48	-14.76	-9.72	coop
G''	-2.68	-2.06	-0.62	-120.96	-171.73	50.77	-24.48	-14.76	-9.72	coop
H	-3.44	-2.86	-0.57	-129.81	-203.97	74.15	-38.00	-22.73	-15.27	coop
H'	-1.97	-1.46	-0.51	-101.36	-114.29	14.16	-13.01	-10.61	-2.40	coop
H''	-2.77	-2.06	-0.71	-116.14	-171.73	55.59	-25.88	-14.76	-11.12	coop
I	-3.47	-2.86	-0.60	-140.97	-203.97	63.00	-35.84	-22.73	-13.11	coop
I'	-1.74	-1.69	-0.05	-88.54	-88.54	0.00	-14.16	-14.16	0.00	coop
I''	-2.57	-2.29	-0.28	-114.49	-145.98	31.49	-24.86	-18.31	-6.55	coop
K	-4.59	-3.92	-0.67	-172.62	-235.61	62.99	-46.13	-33.46	-13.11	coop
K''	-4.56	-3.77	-0.68	-171.53	-234.52	62.99	-45.57	-32.46	-13.11	coop
L	-5.24	-4.52	-0.72	-164.01	-293.05	129.04	-64.10	-37.17	-26.93	coop
L'	-4.23	-3.09	-1.14	-176.97	-202.83	25.86	-29.57	-24.77	-4.80	coop
M	-5.45	-4.52	-0.94	-167.06	-293.05	125.99	-63.39	-37.17	-26.22	coop
M'	-3.80	-3.33	-0.48	-177.09	-177.09	0.00	-28.32	-28.32	0.00	coop

ΔG_{FEP} denotes the binding free energies from FEP, (free energy perturbation) ΔG_2 denotes the predicted binding affinities by assuming pairwise additivity, and $\Delta\Delta G$ denotes the deviation of pairwise additivity predicted binding affinities from the corresponding FEP results (Which corresponds to non-additivity of hydrophobic interactions defined by formula 9.5). $\Delta SASA$ denotes the change of SASA (solvent accessible surface area) upon methane binding. the $\Delta SASA_2$ denotes the pairwise-additive contribution to the change

Table 9.2: Multi-body PMF (potential of mean force) calculated from FEP, (free energy perturbation) and the contribution from two-, three-, four-, five-body interactions.

Systems	W	W_2	ΔW_2	W_3	ΔW_3	W_4	ΔW_4	W_5	ΔW_5
G	-4.90	-4.37	-0.52	-0.47	-0.06	-0.06	0.00	-	-
H	-7.04	-5.95	-1.08	-1.46	0.38	0.38	0.00	-	-
I	-6.85	-6.19	-0.65	-0.74	0.09	0.09	0.00	-	-
K	-7.42	-7.19	-0.24	0.14	-0.37	-0.37	0.00	-	-
L'	-7.07	-6.36	-0.71	-0.14	-0.56	-0.56	0.00	-	-
M'	-6.64	-6.60	-0.04	0.78	-0.82	-0.82	0.00	-	-
L	-12.31	-10.88	-1.43	-1.40	-0.03	-0.53	0.51	0.51	0.00
M	-12.09	-11.12	-0.98	-0.21	-0.73	-1.39	0.63	0.63	0.00

The PMF between two plates was set to zero when $D = 7.46\text{\AA}$, which corresponded to the configuration for all the systems including two plates. The choice of base line for the PMF between the two plates does not affect the multi-body contributions.

Free energies in $kcal.mol^{-1}$.

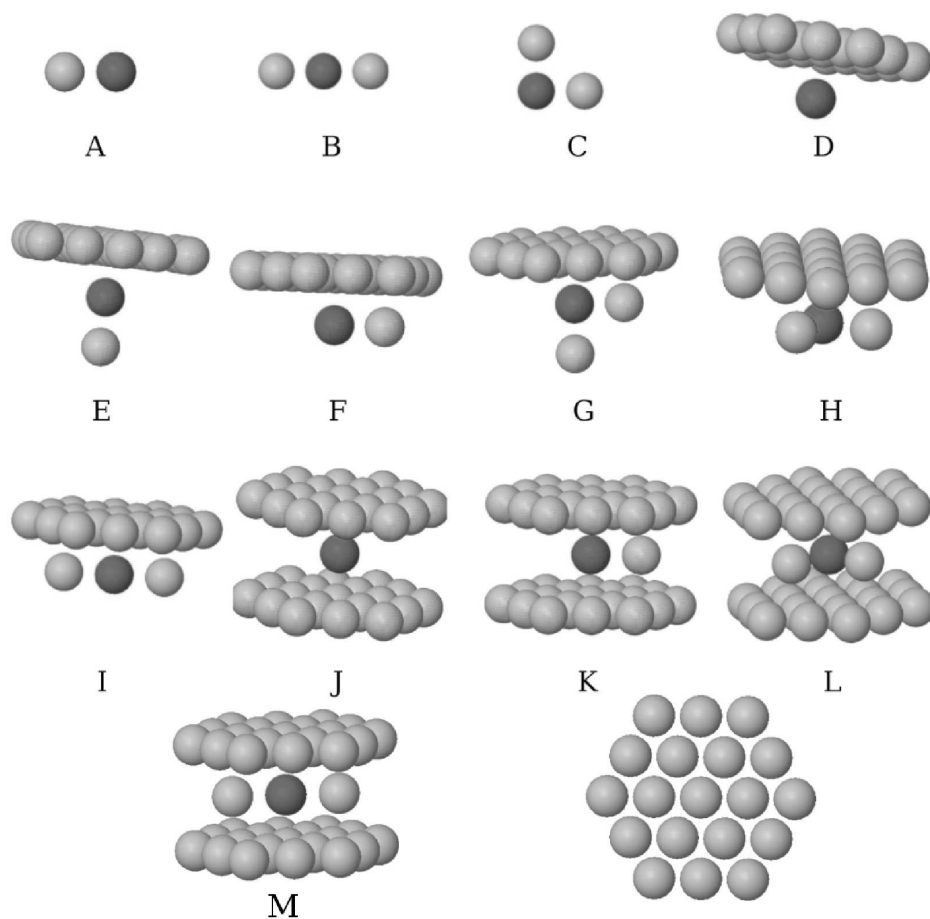


Figure 9.1: The 13 model systems studied. The hydrophobic enclosures were depicted in gray. The location of the methane molecule when bound to the respective hydrophobic enclosure was depicted in black. The geometry of the hydrophobic plate was depicted in the right bottom of the figure.

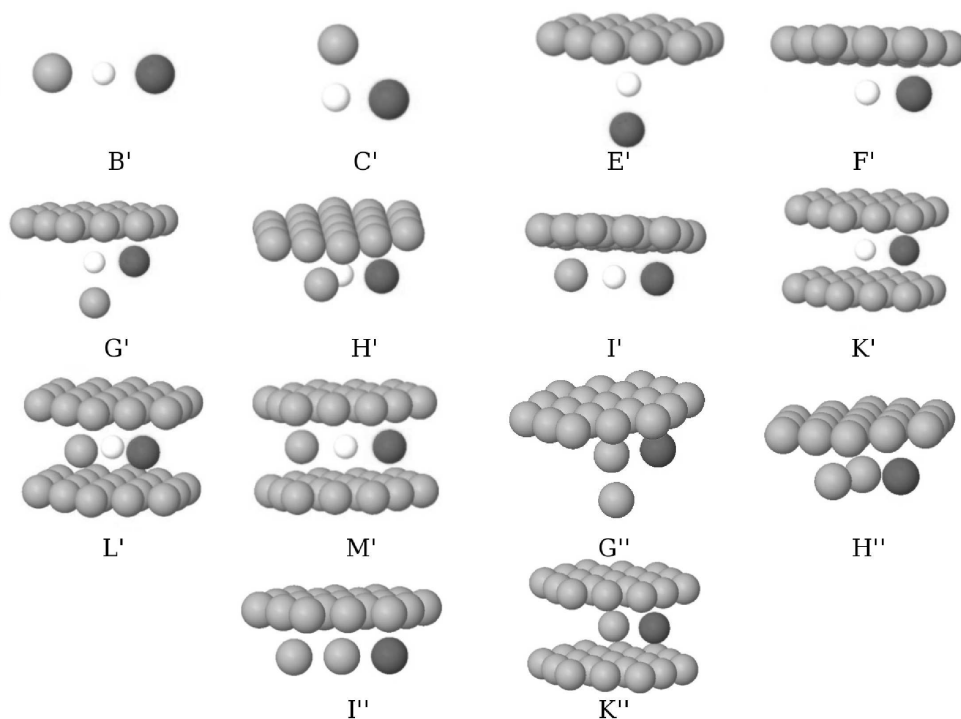


Figure 9.2: The 14 model systems corresponding to subsystems depicted in figure 7.1. The black particle in each system denote the methane that will bind to the enclosure which was depicted in gray, and the small white particle denote a pseudo-particle that specified the position of the binding methane for the corresponding system in figure 7.1. (The binding affinity for system G'', H'', I'' and K'' were calculated through thermodynamic cycles by combination of the binding affinities calculated for related systems.)

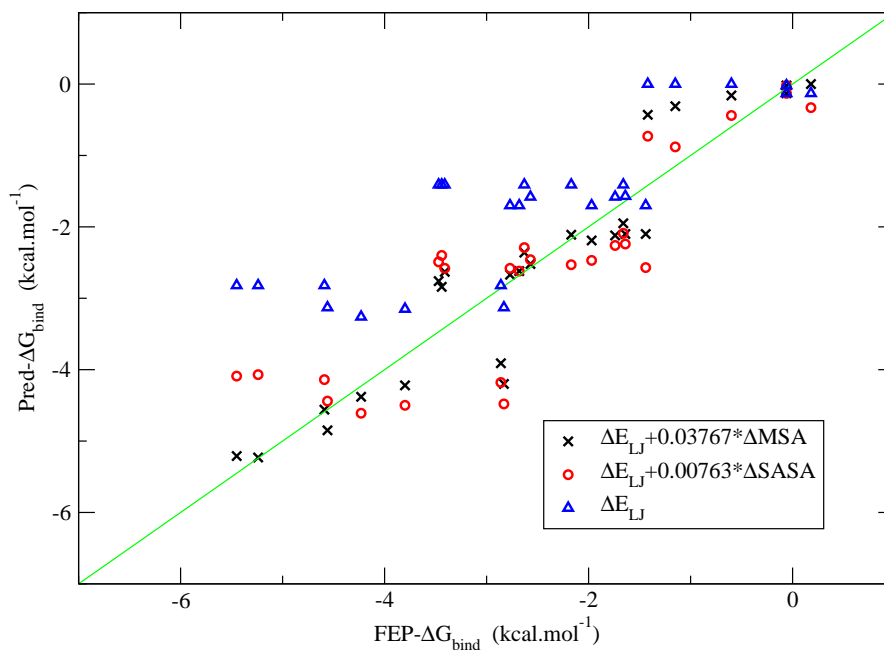


Figure 9.3: Buried surface area/molecular mechanics prediction of methane-enclosure binding affinities. The surface tension coefficients were chosen to minimize the MAE (mean average error) between the predicted and FEP results for the binding affinities. Predictions based on MSA (molecular surface area) model performed better than those based on SASA (solvent accessible surface area) model, indicated by a higher R^2 value (0.89 vs 0.76) and smaller MAE (0.40 vs 0.57). Both of these models performed much better than predictions based on the van der Waals surface model which takes into account the direct LJ interaction (giving $R^2 = 0.70$, and MEA = 0.94 kcal/mol). SASA can not differentiate the hydrophobicity among systems [J, K, K', K'', L, L', M, M'] (predicting a similar binding affinity of about -4.1 kcal/mol) nor among systems [D, E, F, F', G, G', G'', H, H', H'', I, I', I''] (predicting a similar binding affinity of about -2.2 kcal/mol), while MSA based model performed much better for these systems and predicted the right order of hydrophobicity among these systems to some extent.

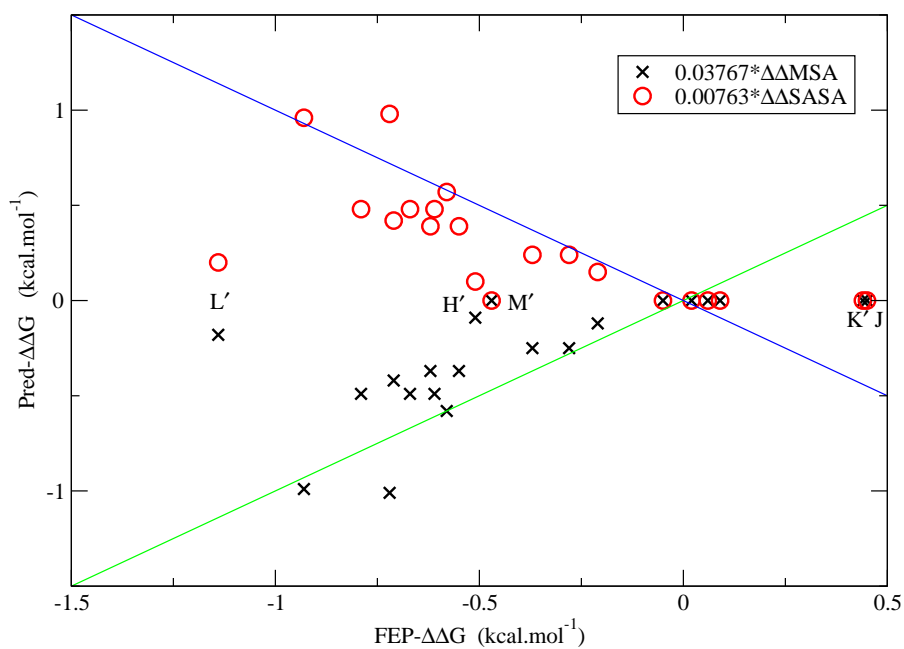


Figure 9.4: Relationship between the surface area models predicted non-additivity of hydrophobic effects and the corresponding FEP results. There is a strong correlation between MSA (molecular surface area) predictions and FEP results, but the SASA (solvent accessible surface area) predicted results anti-correlate with the FEP results.

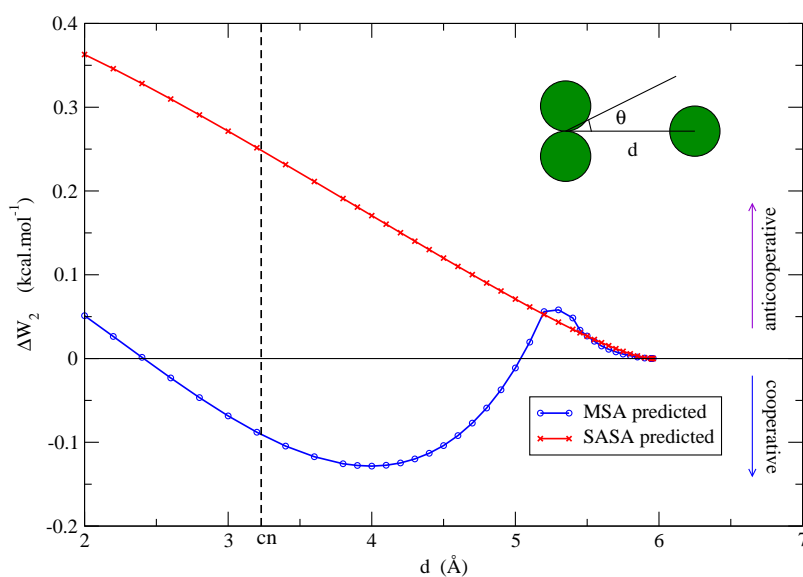


Figure 9.5: The MSA/SASA (molecular surface area/solvent accessible surface area) predicted non-additivity effects of hydrophobic interactions (non-additive part of the change in MSA and SASA multiplied by the corresponding surface tension coefficient) for methane trimer system as a function of the distance d for the configuration depicted in the top right corner when $\theta = 0$. The distance corresponding to the configuration where the third methane is in contact with the remaining two were indicated by the dotted line labeled “cn”.

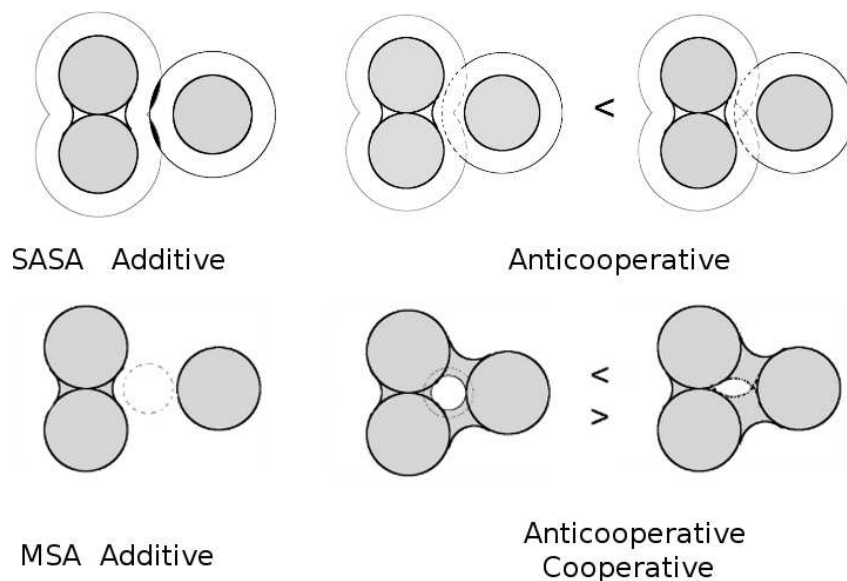


Figure 9.6: In the SASA model (upper), the surface considered is formed from overlapping spheres, each of whose radii is the sum of the corresponding atomic van der Waals radii plus the radius of water.[151] The buried surface area between a sphere and a cluster of spheres, when the cluster of spheres has no overlaps, is always equal to the sum of buried surface areas between the sphere and each individual sphere in the cluster (left), and smaller than this when the cluster of spheres has overlaps (right). This means that the SASA model can never predict cooperative effects. However, in the MSA model (lower), the surface is traced by the inward-facing surface of the probe sphere. The buried surface area between a third particle and the remaining two particles in MSA model, is additive when the third particle has no overlap with the remaining two particles (left), and can be larger or smaller than pairwise additivity predicted buried surfaces (right) depending on the geometry of the three particles.

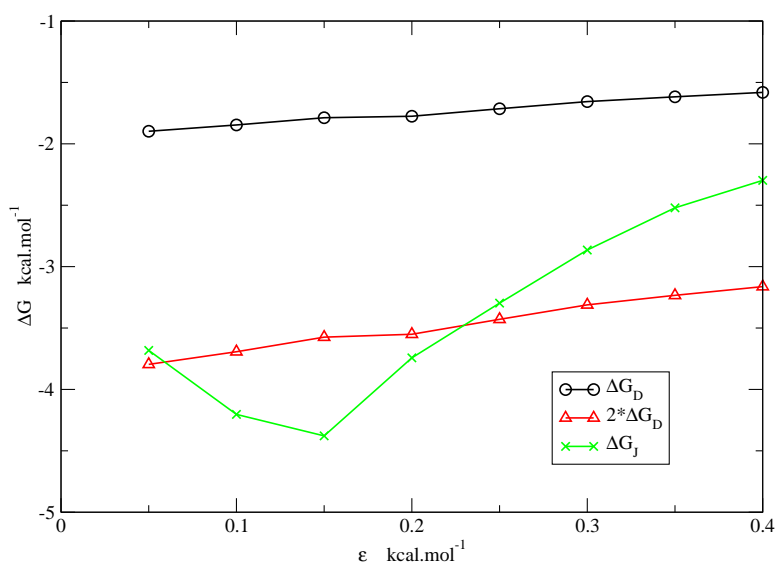


Figure 9.7: The methane enclosure binding affinities for enclosure D and J as a function of the LJ ϵ parameter for atoms making up the plate(s). The binding affinity increases monotonically for enclosure D with increasing value of ϵ , while it decreases at the lower ϵ region and increases at higher ϵ region for enclosure J. The non-additivity effect for enclosure J goes from anti-cooperative in the lower ϵ region, (smaller than 0.06 kcal/mol), to cooperative in the intermediate ϵ region, (between 0.06 and 0.23 kcal/mol), and then to anti-cooperative again in the higher ϵ region (larger than 0.23 kcal/mol).

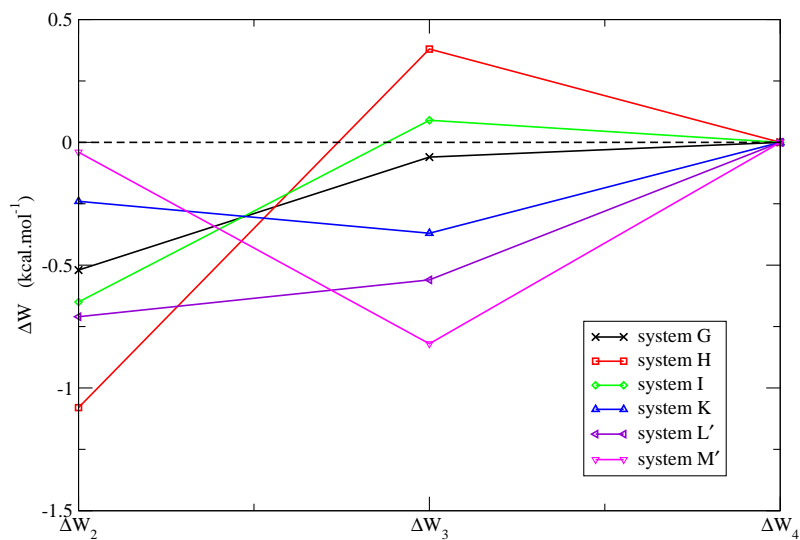


Figure 9.8: The deviation of GKSA (generalized Kirkwood superposition approximation) predicted PMF (potential of mean force) from the total PMF by truncating the total PMF to the second-, third-, and fourth- order as a function of the order.

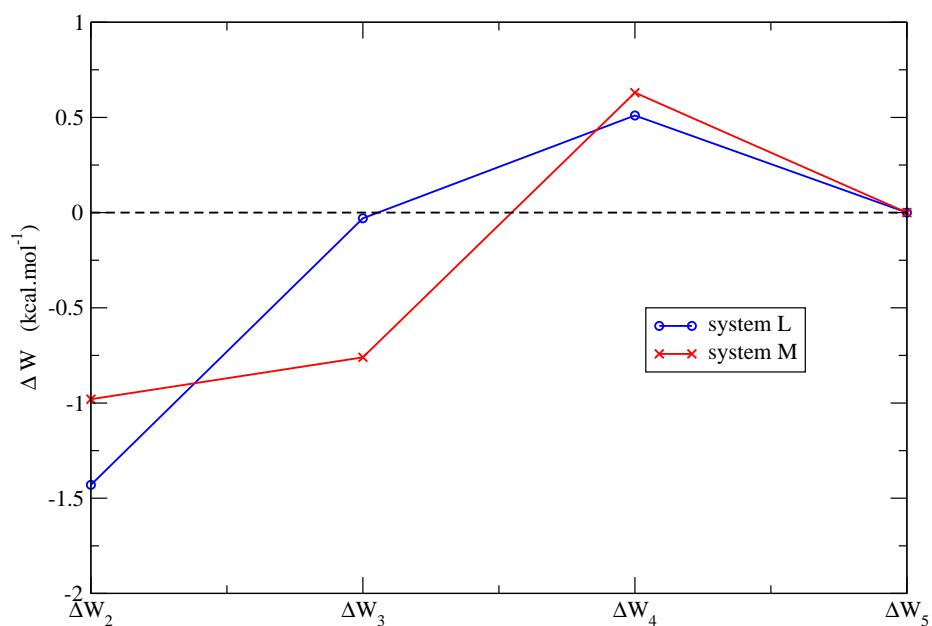


Figure 9.9: The deviation of GKSA (generalized Kirkwood superposition approximation) predicted PMF (potential of mean force) from the total PMF by truncating the total PMF to the second-, third-, fourth- and fifth- order as a function of the order.

Chapter 10

Competition of electrostatic and hydrophobic interactions between small hydrophobes and model enclosures

Abstract

The binding affinity between a probe hydrophobic particle and model hydrophobic plates with different charge (or dipole) densities in water was investigated through molecular dynamics simulations free-energy perturbation calculations. We observed a reduced binding affinity when the plates are charged, in agreement with previous findings. With increased charge density, the plates can change from “hydrophobic like” (pulling the particle into the interplate region) to “hydrophilic like” (ejecting the particle out of the interplate region), demonstrating the competition between hydrophobic and electrostatic interactions. The reduction of the binding affinity is quadratically dependent on the magnitude of the charge for symmetric systems, but linear and cubic terms also make a contribution for asymmetric systems. Statistical perturbation theory explains these results and shows when and why implicit solvent models fail.

10.1 Introduction

Hydrophobic interactions give rise to solvent induced attractions between nonpolar particles when solvated in water. They play an important role in protein folding, protein ligand binding, and micelle formation.[125; 124; 1] While great efforts have been made by many groups to study the interactions between pure hydrophobic particles or plates, from small to large length scales,[154] relatively less effort has been made to understand the effect of electric charge on the hydrophobic interactions. Yet, most bio-molecular solutes, such as proteins, carry partial charge. It is of interest to further study how the solute-solvent electrostatic interactions affect the binding free energies of nonpolar particles in charged hydrophobic enclosures.

There has been recent work looking at the structure and compressibility of water at hydrophobic/hydrophilic interfaces[158], connecting the hydrophobicity of the surface with the solute binding affinity; heterogeneous surfaces with mixed hydrophobic and hydrophilic patches were also studied.[159] However, all these studies were concerned with one surface, and the structure and dynamics of water in enclosed systems, where water are surrounded on multiple sides by hydrophobic or hydrophilic moieties, a key motif in many important protein receptors for its molecular recognition, have not been studied. In this study we investigate in quantitative detail the effects of enclosure on a model system containing both hydrophobic and hydrophilic components. The model system work is complementary to our investigation of protein active sites which has had significant impact on the drug discovery community.[3]

It is well known that when two sufficiently large hydrophobic plates are closer than a critical distance, the interplate region dewets.[160; 161; 162; 154] And in such heterogeneous environments, there is a sensitive coupling of hydrophobicity to the changes in local geometry, dispersion, and electrostatic interactions.[1] Recently, Hansen *et al.*[163] observed a strong reduction of the critical distance for dewetting between two nanoscale solutes when they were charged, and the effective hydrophobic interactions between the solutes were also reduced. In addition, the reduction of the interactions is sensitive to the charge pattern on the solutes, and there is a significant asymmetry between anionic and cationic solute pairs.[164] The asymmetry between cationic and anionic solvation free

energy is a well known fact which has been investigated by many groups.[165; 166; 167; 168; 169; 170; 171] Recently, by studying the electric field dependence of the density and polarization density of water between two graphite-like plates,[172] Rasaiah and coworkers found that applying the electric field decreases the density of the water between the plates, contrary to Hansen's conclusions, and to bulk fluid electrostriction. Rossky *et. al.* also observed an enhanced hydrophobicity of silica surfaces when the charges on Si and O are inverted compared to that of a fictitious neutral silica surface.[173] Thus, surface polarity is important and sometimes acts in unexpected ways. In addition, Zangi and coworkers [174] have studied the effect of cosolute ions on the potential of mean force (PMF) between two hydrophobic plates, and they found that, for cosolute ions with charge density higher than 0.90, the PMF between the plates will increase; and for cosolute ions with charge density lower than 0.90, the PMF will decrease.

In this paper, we study the binding affinities between a probe hydrophobic particle and model hydrophobic plates through molecular dynamics simulations, and by placing charges or dipoles on the plates, we investigate electrostatically induced interactions between the probe particle and the plate. The plate-water interaction is such that there is no dewetting between the two plates as in the above studies. We find that, for small charges, the binding free energy is negative, indicating the plates remain hydrophobic; however, for large charges, the binding free energy is positive. Thus, as expected, the electrostatic interaction between the charges on the plates and the solvent can drive the plates from being hydrophobic to being hydrophilic.

We also find that the binding affinity of the small particle depends quadratically on the magnitude of the charge (or dipole) on parallel symmetric plates, that is plates with the same sized ions (or dipoles). This is not surprising. The electrostatic contribution to binding affinity between the probe particle and the plates is the difference of electrostatic contribution to the solvation free energy for systems with and without the probe particle. Thus, implicit solvent models such as GB or PB also predict a quadratic dependence on the magnitude of charge.[175; 176; 177; 178] However, for plates with different sized ions, the linear and cubic charge (or dipole) dependent terms make small contributions to the solvation free energy, which is contrary to the implicit solvent model predictions.[177] All of

the observed effects can be explained by statistical perturbation theory using results from explicit solvent models, but not implicit solvent models.

10.2 Details of simulation

We performed molecular dynamics simulations using the DESMOND program [31] to study the binding affinities between a united-atom methane and two hydrophobic plates. The geometry of the model hydrophobic plate is displayed in figure 10.1a. It consists of 19 single-layer “atoms” arranged in a triangular lattice with a bond length of 3.2 Å. In two plate systems, the plates are parallel and in-registry with a separation distance of $D = 7.46$ Å. (which is two times the LJ σ parameter of methane, so the methane can just fit in between the plates.) The plate atoms forming the enclosures all have Lennard Jones parameters $\sigma = 3.73$ Å and $\epsilon = 0.294$ kcal/mol, which are the same as the united-atom methane parameters used in these simulations.[149] The inserted methane particle (displayed in green in figure 10.1) is placed at the center of the two plates. Then we place opposite charges on the two center atoms of the two plates, or two dipoles pointing in opposite directions, to see how electrostatic perturbation of water affects the binding affinities. The two oppositely charged atoms can be the same size or of different sizes. The plates with the same sized ions (or dipoles) are designated a symmetric system, whereas the plates with different sized ions is designated an asymmetric system.

The free energy perturbation (FEP) method was used to determine the binding affinities between the inserted methane and the two plates. We used the Maestro System Builder utility [65] to insert each system into a cubic water box with a 10 Å buffer. The water molecules interact through the SPC model.[38] In these simulations, the atoms of the plates were constrained to their initial positions, and only the solvent degrees of freedom were sampled. The united-atom methane was “turned on” inside the two plates over 9 lambda windows with $\lambda = [0, 0.125, 0.25, 0.375, 0.50, 0.625, 0.75, 0.875, 1]$, where λ is the coupling parameter to turn on/off the LJ interaction between the methane and the rest of the system with initial state and final state correspond to $\lambda = 0$ and $\lambda = 1$ respectively. The core of the LJ potential for methane is made softer[67] as $\lambda \rightarrow 0$ to avoid singularities and

numerical instabilities for FEP simulation. For each of the λ windows, molecular dynamics simulations were performed. The energy of the system was minimized, and then equilibrated to 298 K and 1 atm with Nose-Hoover[33; 34] temperature and Martyna-Tobias-Klein[35] pressure controls over 100 ps of molecular dynamics. A cutoff distance of 9 Å was used to model the Lennard Jones interactions, and the particle-mesh Ewald method[37] was used to model the electrostatic interactions. Following equilibration, a 20 ns production molecular dynamics simulation was performed and configurations of the system were collected every 1.002 ps. The energy difference between neighboring λ windows for each configuration saved was calculated and the Bennett acceptance ratio method[30] was used to calculate the free energy difference between neighboring states. The sum of the free energy differences between neighboring states gives the solvation free energy of methane in the enclosure between the plates. The same procedure was followed to calculate the solvation free energy of methane in bulk water. The difference between the two solvation free energies gives the binding affinity between the methane and the two plates. The error associated with these binding affinities is of order ± 0.02 kcal/mol.

As indicated in the thermodynamic cycle in figure 10.1, the electrostatic contribution to the binding affinity $\Delta F_2 - \Delta F_1$, is equal to $\Delta F_4 - \Delta F_3$, which is the free energy difference of charging the plates in water with and without the inserted methane. In order to investigate the electrostatic contribution to the binding affinities as a function of charge, we did additional FEP simulations to turn on the electrostatic interaction between the plates and the rest of the system for systems with and without the inserted methane. The FEP protocols were similar to that used for the calculation of the solvation free energy of methane, but here we used 16 lambda windows with $\lambda = [0, 0.02, 0.04, 0.06, 0.08, 0.10, 0.15, 0.20, 0.30, 0.40, 0.50, 0.60, 0.70, 0.80, 0.90, 1.00]$, and 6ns of data collection for each lambda window, where λ is the coupling parameter to turn on the electrostatic interaction between the charges on the plates and the rest of the system.

10.3 Results and discussion

10.3.1 Binding affinity results

The free energy results for each process with unit charge on corresponding atoms are depicted on the thermodynamic cycles in figure 10.1. The free energy changes along different paths of each half cycle only differed by 0.1 kcal/mol, indicating the high accuracy and precision of these free energy results. We see clearly that without charges or dipoles on the hydrophobic plates, there is a strong thermodynamic driving force to pull the methane into the region between the plates ($\Delta F = -2.865$ kcal/mol); however, if we put charges or dipoles on the plates, the methane is ejected from the enclosed region ($\Delta F > 10$ kcal/mol). This agrees with previous findings for the reduction of the hydrophobic interaction between two hydrophobic particles when they are charged.[163; 164] By putting charges or dipoles on the hydrophobic plates, the plates change from “hydrophobic like” (methane absorption) to “hydrophilic like” (methane ejection). It also indicates that even small hydrophilic patches on hydrophobic surface can have a strong effect on the hydrophobicity of the surface, which was observed in previous studies.[159] This behavior is expected: without charges or dipoles on the plates, water molecules can not make hydrogen bonds with the plates, so they would prefer to be away from the region between the plates to make hydrogen bonds with other water molecules, and methane would be driven into the region between the plates because it can neither make hydrogen bonds with water nor with the plates. However, if there are sufficiently large charges or dipoles on the plates, water molecules can make hydrogen bonds with the plates, or at least have an attractive polar interaction with the plates, so that it would be favorable for them to be there over the methane, and methane would be ejected from that region.

The binding affinity difference ($\Delta F_2 - \Delta F_1$), as mentioned before, arises from the free energy difference of charging the plates with and without the inserted methane ($\Delta F_4 - \Delta F_3$). This is also equal to the difference of the electrostatic contributions to the solvation free energy for the two systems, because the direct electrostatic interactions in solutes for the two systems are the same. It is well known that the more the ions are exposed to water, the more the electrostatic interaction contributes to the solvation free energy.[176] The ions

on the plates are more exposed to water without the inserted methane, so ΔF_3 is more negative than ΔF_4 , which provides another perspective for understanding the transition from “methane absorption” to “methane ejection” due to putting charges or dipoles on the hydrophobic plates.

10.3.2 Dependence of the binding affinity (Solvation free energy) on the magnitude of charge

To investigate quantitatively how the magnitude of charges or dipoles affects the hydrophobic or hydrophilic properties of the plates, we calculated the binding affinities of systems with different charge densities on the charged or polar atoms of the plates. (Here the radius of the atoms are fixed and only the magnitudes of the charges are varied.) Figure 10.2 depicts the relationship between the free energies of charging the plates for the two systems (corresponding to ΔF_3 and ΔF_4 in part b of figure 10.1) and the magnitude of the charge on the atoms, q , and also q^2 (Inset of the figure). We see clearly from this figure that the free energy of charging the electrostatic interactions for these two systems are proportional to the square of the magnitude of the charge on the atoms (or the charge density). From these data, we determine the methane-plates binding affinities as a function of the magnitude of the charge. Clearly, the binding affinity should also have a quadratic dependence on the magnitude of the charge. Figure 10.3 shows the methane-plates binding affinity as a function of q^2 , and can be perfectly fit by a straight line. If we define the plates to be hydrophobic or hydrophilic by the sign of binding affinity, negative binding affinity corresponding to hydrophobic and positive binding affinity corresponding to hydrophilic, then in the low charge region ($q < 0.37$), the plates are hydrophobic, and in the high charge region ($q > 0.37$), the plates are hydrophilic. The crossover point occurs at about $q \approx 0.37$ ($q^2 \approx 0.137$), where the plates change from being hydrophobic to being hydrophilic. Interestingly, Zangi et al.[174] have studied the effect of cosolute ions on the PMF between two hydrophobic plates, and they found that for cosolute with a charge density of 0.90, the PMF was the same as that in pure water, and lower charge density cosolute will decrease the hydrophobic interaction between the plates, and higher charge density cosolute will increase the hydrophobic interaction. However, both the trend and the crossover point charge density observed here

is different. This is not surprising: in their studies many ions were dissolved in solvent, while here one ion was placed on one hydrophobic plate and one oppositely charged ion is placed on the other hydrophobic plate. Also the size of the plates and LJ parameters for atoms making up the plates are different, so dewetting occurred in their systems but not here. For hydrophobicity as defined here, the crossover point charge density will depend on the LJ parameters for atoms making up the plates, and the size of the plates, but the trend should be the same. Similar results were observed for systems with dipoles on the two plates,(results not shown) only the slope was different.

People familiar with implicit solvent models such as PB or GB,[177; 178; 175; 176] would not be surprised by the quadratic dependence of the solvation free energy on the magnitude of charge. In these models, the electrostatic potential or the induced surface charge is proportional to the magnitude of the charge on the solute, so the electrostatic contribution to the solvation free energy is proportional the square of the magnitude of the charge.[177; 178; 175; 176] The direct electrostatic interactions are trivially proportional to the square of the magnitude of the charge. So the free energy of charging the plates, which is the sum of the two terms, should also have a quadratic dependence on the magnitude of the charge. However, the constant dielectric approximation in such models is clearly not a good approximation for these systems. Figure 10.4 depicts the projection of the orientation of water molecules in the region between the two plates from 10000 frames with unit charge on the corresponding atoms of plates. Clearly, water molecules are highly structured in this region, and they tend to make hydrogen bonds with the two charged atoms on the plates, so the constant dielectric approximation does not apply here. Although different dielectric constants can be assigned for different regions of the solution when solving the PB equation, it is generally difficult to assign these parameters without prior knowledge of the structure of solvent, and this technique is usually used only for the solute region. In the next section, we will explain this effect by a theory based on explicit solvent models, and the quadratic dependence of the solvation free energy on the magnitude of charge for such systems comes naturally from this theory.

10.4 Theoretical derivation for electrostatic contribution to the solvation free energy

For a solute molecule composed of N_A atoms solvated in N_s solvent molecules, the total interaction energy of the systems is:

$$U(\mathbf{r}_A, \mathbf{r}_s, \epsilon) = U_A(\mathbf{r}_A) + U_s(\mathbf{r}_s) + \sum_i^{N_A} \sum_s^{N_s} u_{np}(\mathbf{r}_i, \mathbf{r}_s) + \epsilon \sum_i^{N_A} \sum_s^{N_s} u_p(\mathbf{r}_i, \mathbf{r}_s) \quad (10.1)$$

$$= U_A(\mathbf{r}_A) + U_s(\mathbf{r}_s) + U_{As}^{np} + \epsilon U_{As}^p \quad (10.2)$$

where U_A is the intramolecular interactions of the solute and U_s is the intra- and inter-molecular interactions between the N_s solvents molecules, the first summation term on the right hand side is the nonpolar interactions between the solute and the solvent, and the last term on the right hand side is the polar (or electrostatic) interactions between the solute with charge scaled by a scaling parameter ϵ and the solvent. Through thermodynamic perturbation theory, the electrostatic contribution to the solvation free energy of the solute with charge scaled by ϵ can be expressed as:

$$\beta \Delta F_p = -\ln \langle e^{-\beta \epsilon U_{As}^p} \rangle_0 \quad (10.3)$$

where $\beta^{-1} = k_B T$, k_B is the Boltzmann constant, and $\langle \dots \rangle_0$ means the ensemble average of the mechanical properties over unperturbed state where there is no electrostatic interaction between the solute and the solvent. Here, to make the derivation neater, the solute is kept fixed, and only the solvent degrees of freedom are integrated over. Expanding Eq. 10.3 in powers of ϵ we get the electrostatic solvation free energy in powers of the magnitude of charge on the solute,

$$\Delta F_p = \epsilon \langle U_{As}^p \rangle_0 - \frac{\beta}{2} \epsilon^2 \langle (U_{As}^p - \langle U_{As}^p \rangle_0)^2 \rangle_0 + \frac{\beta^2}{6} \epsilon^3 \langle (U_{As}^p - \langle U_{As}^p \rangle_0)^3 \rangle_0 + \dots \quad (10.4)$$

This result is similar to what Hummer *et al.*[170; 171] get in their studies of ion hydration, except that in their studies the overall charge of the system is not zero, so the finite size effect had to be included explicitly.

Let us now analyze the coefficients of the linear and quadratic terms. The linear term is the average of the electrostatic interaction between the fixed solute and the solvent over the

unperturbed configurations of the solvent ($\epsilon = 0$). For neutral solute molecules, there exists excellent cancellation between interactions from positively charged atoms on solute and interactions from negatively charged atoms on solutes, so the coefficient of the linear term should be small. For symmetric systems like the symmetric parallel plates we studied, this linear term should be exactly zero. For systems with only a single charge, the linear term will be non-negligible and of opposite signs for cations and anions. This explains the asymmetry between cations and anions both for the solvation free energy [165; 166; 167; 168; 169; 170; 171] and for the reduction of the PMF.[163; 164] The quadratic term is proportional to the variance of distribution of U_{As}^p , which is nonzero, so the coefficient should be a large negative number, which makes sense because the electrostatic solvation free energy is negative for almost all systems studied up till now, and for implicit solvent models such as PB or GB.[177; 178; 175; 176] In addition, the coefficient for the second order term is symmetric with respect to charge inversion, which also is consistent with PB or GB predictions. (In other words, if the sign of the charge on the solute was reversed, the coefficient of this term does not change.) The coefficient of the cubic term also depends on the symmetry of the system: for symmetric systems, the distribution of U_{As}^p should also be symmetric, so the cubic term is exactly zero; however, this term is nonzero for asymmetric systems. Again there exists excellent cancellation in the cubic dependence term, so it should also be small.

Furthermore, since U_{As}^p , the electrostatic interaction between the solute and the solvent, is long ranged and is the sum of many terms, the distribution function of U_{As}^p is expected to be approximately a Gaussian distribution function according to the central limit theorem. So only the first few lower order terms in Eq.10.4 make non negligible contributions to the electrostatic solvation free energy. In addition, according to our analysis, the coefficients of the linear and cubic dependence terms are small, so the quadratic dependence term is the dominating contribution.

Comparing the final results of this theory and the implicit solvent models, it is clear that PB or GB models only predict the quadratic dependent term, which is the most important term as predicted from the theory above. This may be the reason why PB or GB models generally give good results for electrostatic solvation free energies, even though the constant dielectric picture is clearly not true for these systems. In the next section, we will present

some further evidence that the PB or GB models does not give even qualitatively correct predictions for asymmetric systems.

10.5 Further evidence to validate the theory

The four systems studied are all symmetric systems, so the linear and cubic terms should be exactly zero as predicted by theory and indeed FEP gives quadratic dependence of the solvation free energy on the magnitude of charge. In addition, there should also be nonzero linear and cubic terms, if the system is asymmetric, although the magnitude of these terms may be small. For this reason we simulated two plates one with a sodium ion and the other with a chloride ion (parameters for these ions are from ref.[179]) placed on the each center atom on the plates respectively. (Part d of figure 10.1) Now the solute is asymmetric with respect to the size of the ions because the sodium and chloride ions have different LJ parameters.

The free energies for each step of the thermodynamic cycle are given in part d of figure 10.1, and the electrostatic contributions to the solvation free energy in the absence of the inserted methane is given as a function of the magnitude of charge in the left side of figure 10.5. Overall, the quadratic functional still characterizes the trend, but not as well as those for the equal sized ions in the above four systems, and deviations of the fitted curve from FEP data are observed for medium and large charges. If the linear term is included in the fit, the overall performance of the fitting gets better, but there are still large deviations in the small charge region.(See inset of figure 10.5) Only if both linear and cubic terms are included does the fit become excellent over the whole charge range. This observation agrees with the theory: both the linear and cubic terms depend on the symmetry of the system, so they both make contributions to the solvation free energy for asymmetric systems. But overall, the quadratic term is still most important, the coefficient of this term being much larger from those of the linear and cubic terms.

The theory shows that if the charge on the solute were reversed, the sign of the coefficient for the linear term should also be reversed. So another model system was studied where the charge on the sodium and chloride ions were reversed. (reversed sodium chloride ion

system) The electrostatic contribution to the solvation free energy as a function of the magnitude of the charge is shown on the right side of figure 10.5 for this system. Similar to the sodium chloride ion system, both linear and cubic terms contribute to the solvation free energy. More importantly, the coefficients of both the linear and the quadratic terms were of similar magnitude as the sodium chloride system, but the sign of the linear term was reversed, which agrees well with what the theory predicts. The exact coefficient of the cubic term, should also be of the same magnitude but opposite sign upon charge reversal, just like the linear terms. However, because the cubic term only makes small contributions and we can often ignore the terms of higher order, $O(q^4)$, in the fitting, but then the coefficients of the cubic terms we obtain from the fitting will have the same sign and will be different in magnitude for charge reversal. This discrepancy points to a deficiency of fitting with polynomials. In addition for asymmetric systems if the fitting is done without the cubic term the observed deviations of the fitted curve from the FEP data in the small charge region is also caused by similar deficiencies of this approach to curve fitting to polynomials in the charge.

Interestingly, implicit solvent models such as PB or GB incorrectly predict identical electrostatic contributions to the solvation free energy for systems with reversed charge distributions, whereas the perturbation theory correctly predicts it. In our situation, the electrostatic contributions to the solvation free energy for the two systems studied with reversed charge distribution were found to be different (-106 kcal/mol vs. -129 kcal/mol for unit charge), which is in agreement with previous findings of the asymmetry between anionic and cationic solutes.[165; 166; 167; 168; 169; 163; 164; 170; 171] In addition, the sign of the linear term for these two systems is correctly predicted by the perturbation theory. It is well known that water will break one hydrogen bond at the surface of large hydrophobic plates pointing one of its hydrogen atoms towards the plates.[1; 154] Since the sodium ion is smaller than the chloride ion, hydrogens pointing to the uncharged sodium atom get closer to it than to the uncharged chloride atom in the uncharged state ($\epsilon = 0$ state). So the interactions between the positive charge on the sodium ion and the unperturbed solvent ($\epsilon = 0$ state) is larger in magnitude than that for chloride ion, which will result in an overall positive linear term for the sodium chloride ion system and negative linear term for

the reversed sodium chloride ion systems. In contrast, the PB or GB models will always predict a negative electrostatic solvation free energy. However, perturbation theory and FEP simulations show that if the coefficient of the linear term is positive, the electrostatic solvation free energy will be positive in the low charge region. To test whether this is true, additional simulations were performed for the two systems with a reversed charge distribution at small charge [0-0.1]. The electrostatic contribution to the solvation free energy as a function of the magnitude of charge is shown in figure 10.6. From this figure, we can see clearly that the linear term is important for this region, and the electrostatic solvation free energy is positive in the small charge region for the sodium chloride ion system, which further validates the theory presented here.

10.6 conclusions

We have studied the binding affinity between a probe hydrophobic particle and model hydrophobic plates with different charge (or dipole) densities. We found that the binding affinity of the probe particle is strongly decreased by putting charges (or dipoles) on the plates, which agrees with previous observations of the reduction in hydrophobic interaction between two solutes when they were charged.[163] The plates can be either hydrophobic or hydrophilic depending on the charge density of the ions on the plates: in the low charge density regime, the effective free energy of binding of the probe particle in the plate enclosure is negative, and the plates manifest hydrophobic property by pulling the hydrophobic particle into the enclosure; in the high charge density regime, the effective binding free energy is positive, and the plates manifest a hydrophilic property by ejecting the hydrophobic particle out of the enclosure between the plates. The effect of charge on the hydrophobicity of the plates is opposed to the effect of cosolute ions on the PMF between hydrophobic plates studied by Zangi *et al.*[174] because in the latter case the low charge ions can form a double layer around the plates and act as a surfactant.

Quantitatively, the observed reduction of binding affinity is quadratically dependent on the magnitude of charge (or dipole) on the plates. Although implicit solvent models such as PB or GB can predict the quadratic dependence, the constant dielectric approximation

in such implicit solvent models is clearly not valid in the simulated systems. However, from perturbation theory, which does not assume a constant dielectric approximation, the quadratic charge dependence of the solvation free energy for symmetric systems can easily be explained. The quadratic charge dependence of the solvation free energy results from the cancellation of the interactions of the positively and negatively charged atoms on the plates with the solvent molecules. However, for asymmetric plates, the two interactions mentioned above do not cancel exactly, so the theory predicts small linear and cubic terms with charge, which we confirmed by explicit solvent FEP simulations. But implicit solvent models can not predict such effects.

In addition, we found that the electrostatic contribution to the solvation free energy is different for asymmetric systems with reversed charge distribution, in agreement with previous observations of the asymmetry between anion and cation pairs,[165; 166; 167; 168; 169; 163; 164; 170; 171] also not predicted by implicit solvent models. This reversed charge effect is easily explained and predicted by perturbation theory. In addition, we observed a small positive value of the electrostatic contribution to the solvation free energy in the low charge density regime for the sodium chloride plates, as predicted by perturbation theory but not by the implicit solvent models. All of these observations give evidences that perturbation theory provide a guide for understanding the electrostatic contributions to solvation free energy of complicated solutes.

The inability of current implicit solvent models to predict linear and cubic in charge terms in the solvation free energy, the asymmetry between positive and negative ions, and the possible positive electrostatic solvation free energies at low charge, indicates some deficiencies of these models. It has also been shown that the effective solute-solvent interface in these implicit solvent models can vary according to the local electrostatic and dispersion potentials.[180; 181] Recently, there have been some attempts to couple nonpolar and polar solvation free energies into implicit solvent models.[182; 183] The theory and observations in this paper might be helpful for further development of implicit solvent models to incorporate such effects.

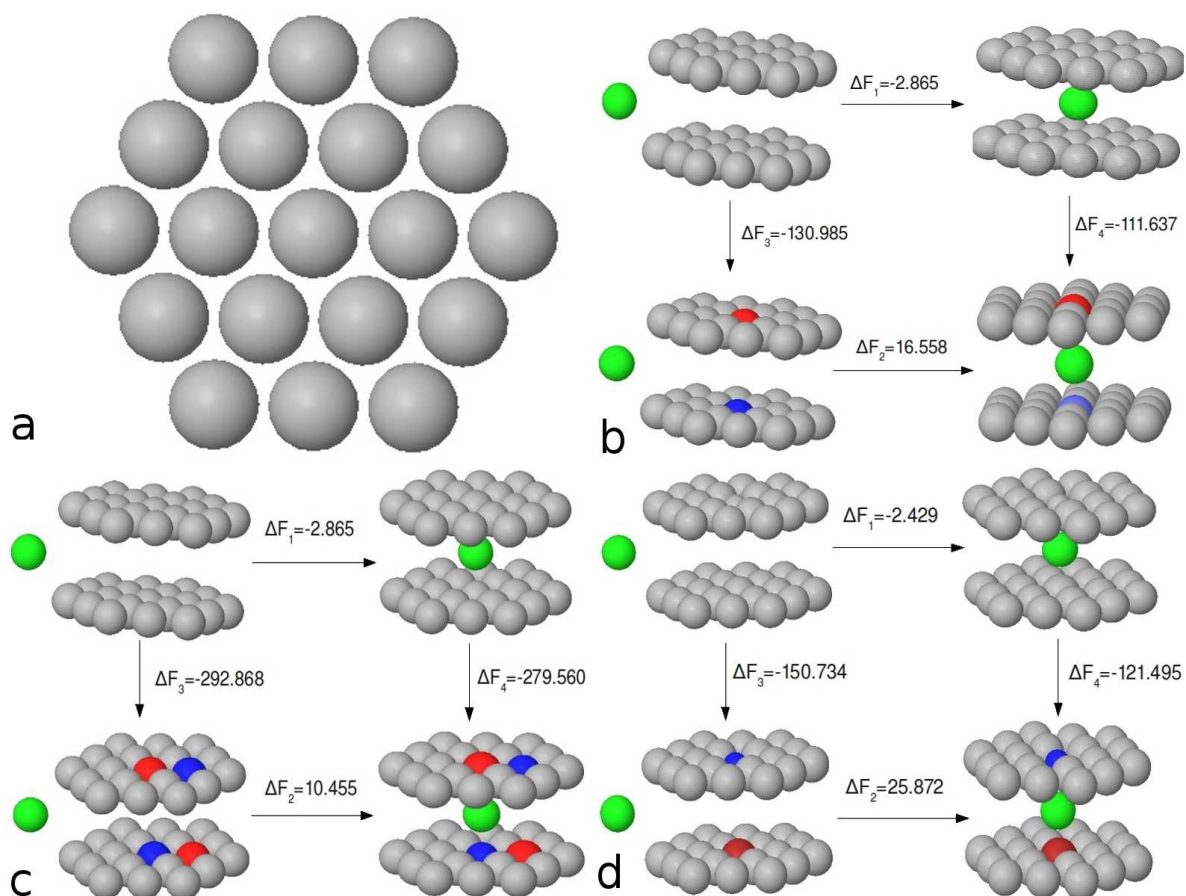


Figure 10.1: Thermodynamic cycles connecting methane-plates binding affinities and the free energies of charging the plates in water. The gray particles represent the LJ atoms forming the enclosure, the red particles represent negatively charged ions, blue particles represent positively charged ions and green particles represent united-atom methane which will bind to the enclosures. a) the configuration of the plate; b) thermodynamic cycle depicting the effect of charges on the methane-plates binding affinity; c) thermodynamic cycle depicting the effect of dipoles on the methane-plates binding affinity; d) the same process as that in part b, but the center ions were replaced by sodium and chloride ions respectively. The free energy changes for each step of the thermodynamic cycle were given in units of kcal/mol.

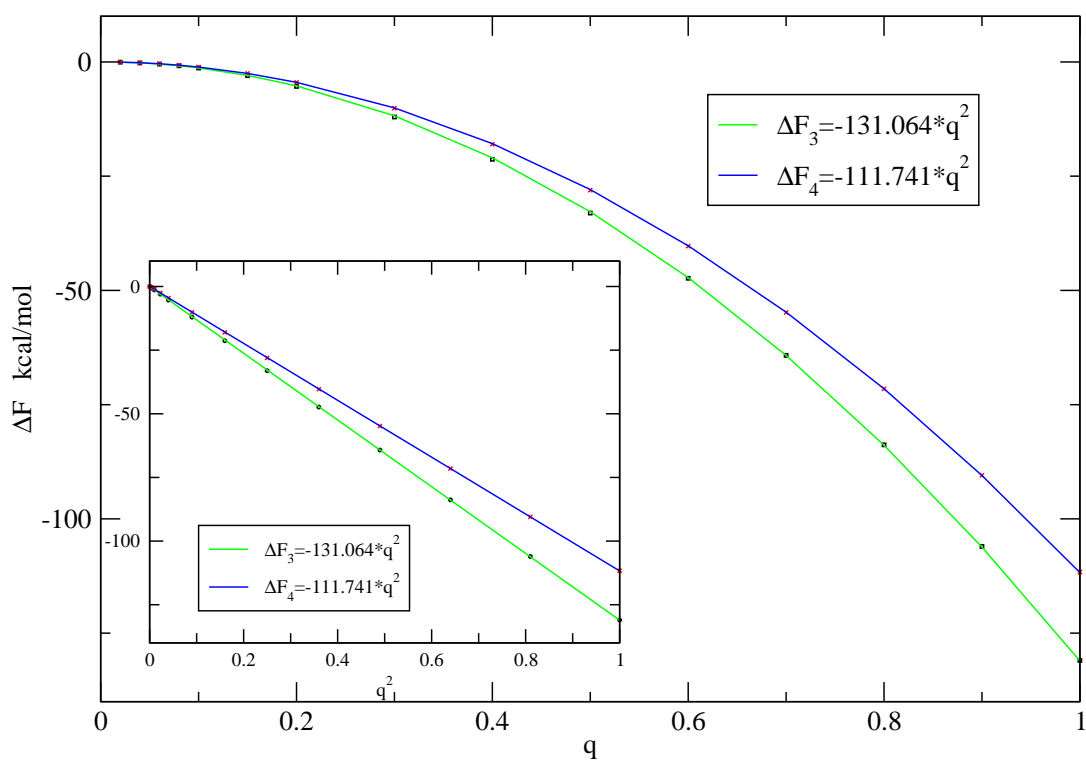


Figure 10.2: Free energy of charging the plates in water with and without the inserted methane as a function of the magnitude of charge, q , and the square of the magnitude of charge, q^2 (inset of the figure). Perfect quadratic dependence of the free energy on the magnitude of charge were displayed by these systems.

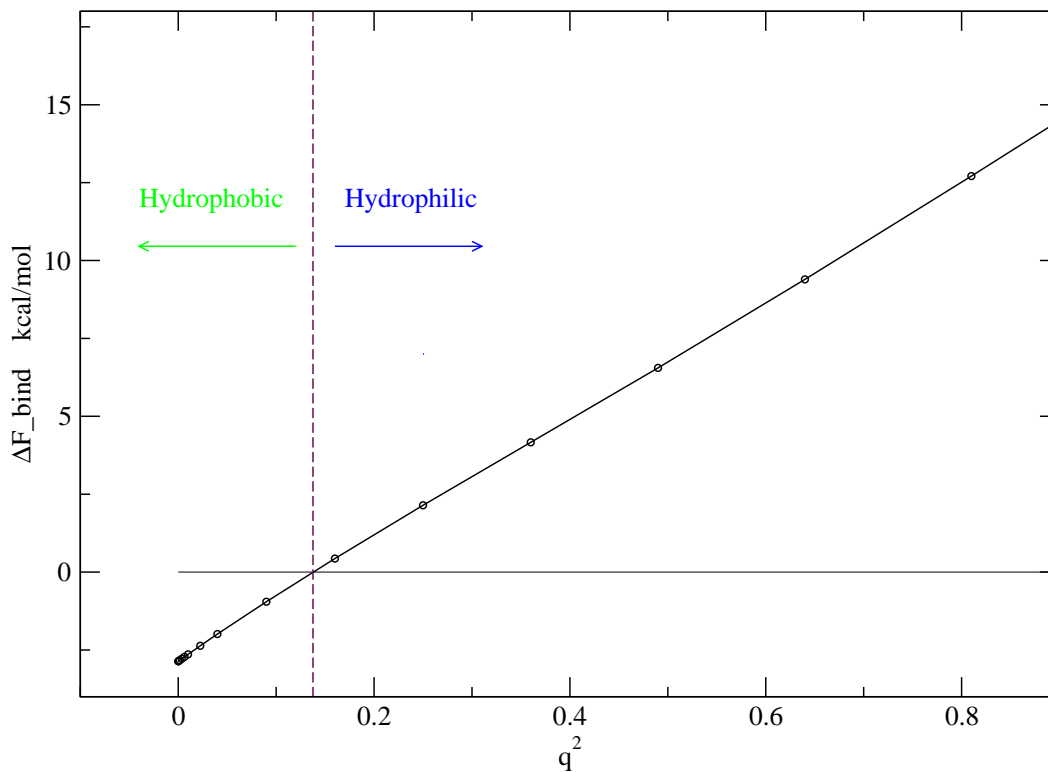


Figure 10.3: Methane-plates binding affinity as a function of the square of the magnitude of charge, q^2 . At low charge density, the binding affinity is negative, displaying hydrophobic property of the plates; however, at high charge density, the binding affinity is positive, displaying hydrophilic property of the plates.

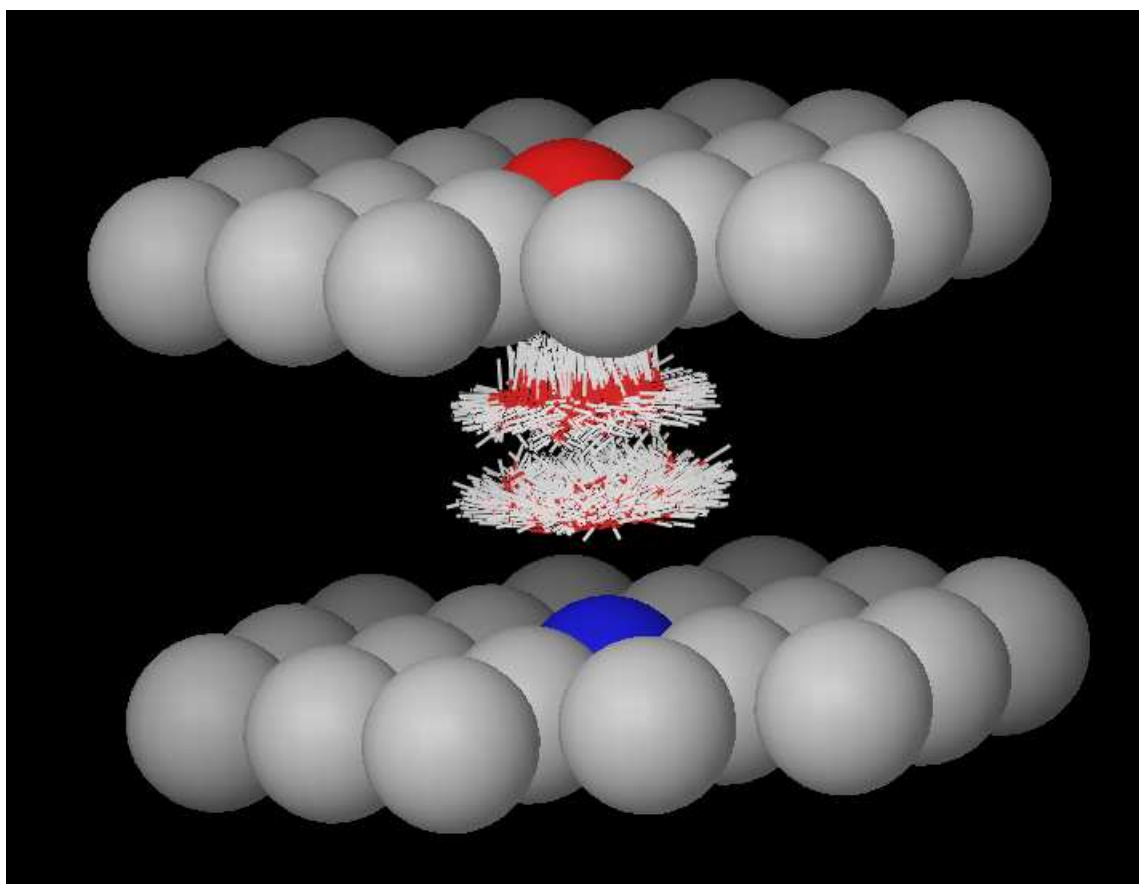


Figure 10.4: Projection of configurations of water between the two plates with two opposite unit charges on the center atoms of the plates from 16000 frames. Water is highly structured in this region, which clearly breaks the constant dielectric assumption of implicit solvent models.

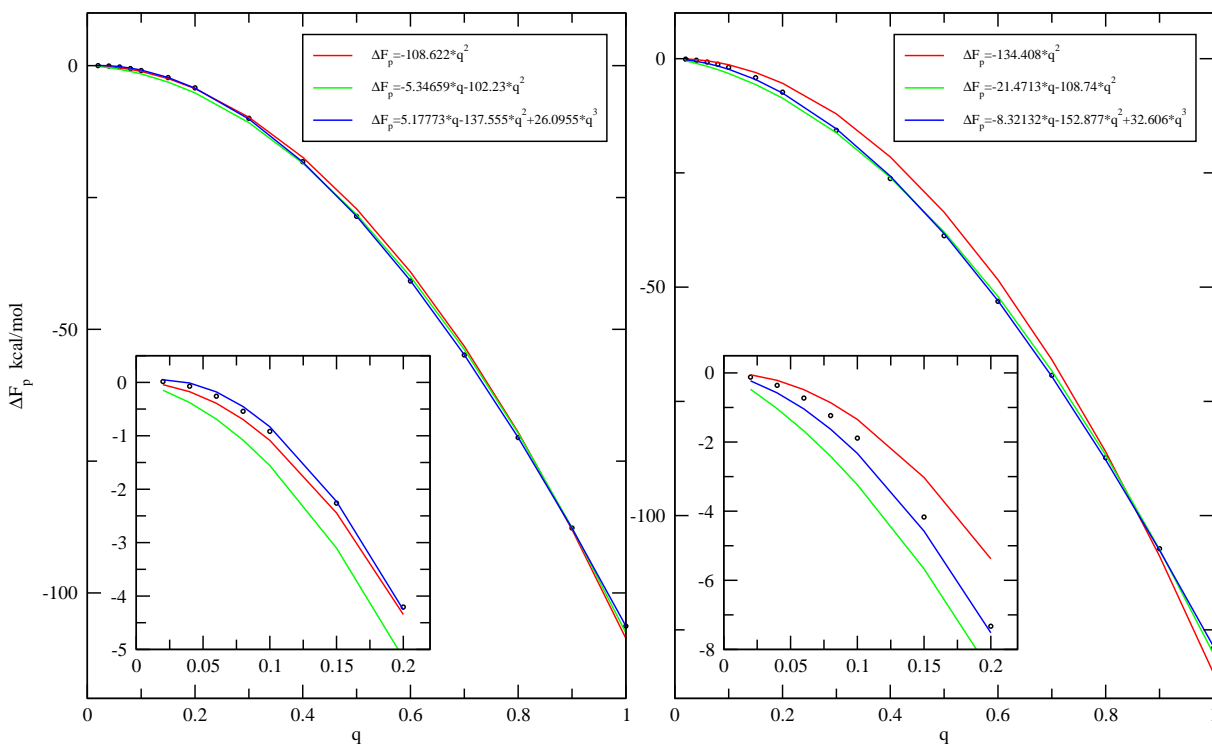


Figure 10.5: Electrostatic contributions to the solvation free energies as a function of the magnitude of charge for sodium chloride ions system (left) and the reversed sodium chloride ion systems (right). Insets of the figures depicts the same curves in the small charge region. Deviations from quadratic dependence appear for these systems. Linear and cubic terms also contribute to the electrostatic solvation free energy. The linear term coefficients for these two systems are approximately of the same magnitude but opposite sign.

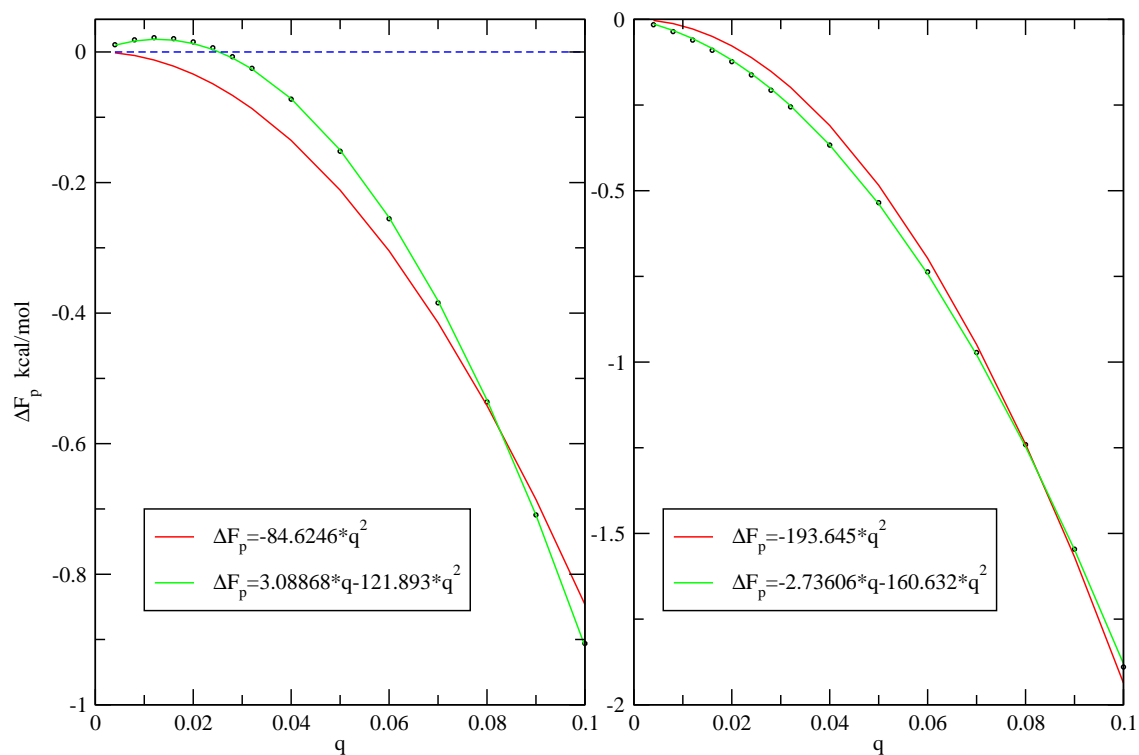


Figure 10.6: Electrostatic contributions to the solvation free energies as a function of the magnitude of charge for sodium chloride ions systems (left) and the reversed sodium chloride ion systems (right) in the small charge region. It is quite clear that the linear terms are important in this region. The electrostatic solvation free energy is positive at very small charge region for the sodium chloride ion system, which PB or GB models fail to predict. Again, the linear term coefficients are approximately of the same magnitude but opposite sign.

Part IV

Conclusions

Chapter 11

Conclusions and future research directions

In this thesis, we have presented a few methods towards more robust and efficient calculation of Protein-Ligand binding affinities. To be specific, the WaterMap method, which focuses on the role of each individual water molecule in the binding pocket of protein to the binding affinity, and the more rigorous free energy perturbation (FEP) method, have been introduced, developed, and discussed. The accuracies and precisions associated with these methods are different so as the computational expenses. They are used in different stages of structural based drug design projects. In what follows, we will briefly summarize the main features of each method, and suggest future research directions.

We have shown that the proteins may adopt active site geometries that will destabilize the water molecules through hydrophobic enclosure and/or correlated hydrogen bonds. In the extreme cases, if the interactions for water molecules are very unfavorable, a void might be formed in the binding pocket. Displacement of these energetically and/or entropically unfavorable water molecules by ligand, and/or occupation of ligand atoms in the dry region of the binding pocket will provide the driven force for Protein-Ligand binding. The WaterMap method and the cavity contribution term consider the explicit driving force from the solvent. Through inhomogeneous solvation theory (ISM), the enthalpy and entropy of each individual water molecule in the binding pocket are calculated, and the free energy

difference between water in the binding pocket and that in bulk will give the contribution to the binding affinity when the water is displaced by ligand in the binding process. The semi-localized water molecules identified by WaterMap will provide rich physical insights in drug design. In addition to the numerical agreement with experiment for the relative binding affinity prediction, the WaterMap calculation provides a vivid picture about the thermodynamics of water in the binding pocket of protein, which can actively guide drug design projects. The cavity contribution term calculates the free energy difference to solvate a ligand heavy atom in the dry region of the protein binding pocket and that in bulk water, which gives the contribution to the binding affinity from the dry regions of binding pocket. For proteins with dry regions in the binding pocket, it is necessary to combine the WaterMap calculation with the cavity contribution term, providing a whole complete picture for contributions from both wet and dry regions of the binding pocket. Those calculations also allow us to suggest a general molecular recognition motif between the dry regions in the binding pocket and hydrophobic groups in the ligand. Calculations based on the WaterMap and the cavity term methods on many different proteins, some of which are of medicinal interest, have shown great success, and we expect these models will prove successful in more systems and will actively guide drug design projects.

The WaterMap method and the cavity contribution term only characterize the contribution from the displacement of solvent to the binding affinity. There are other terms which will contribute to the binding affinity, like the protein-ligand direct interaction energy, the protein and/or ligand strain energy, and the entropy change associated with the protein-ligand association. But for congeneric ligands, those other terms will approximately contribute equally to the binding affinity, so the WaterMap and the cavity contribution term are used to rank order the relative binding affinities between congeneric ligands. Clearly, there are many other augmentations that need to be made before the method can be used to robustly handle a wide variety of cases, particularly when the protein and ligand are not complementary, and the ligands are not congeneric. Here, we suggest the following research directions to further develop the WaterMap method:

- (1) The MM/GBSA or the linear interaction energy (LIE) model assess the Protein-Ligand binding affinity by treating the solvent as a dielectric media or by approximating

the free energy as a linear combination of the average interaction energies through a short molecular dynamics simulation of the Protein-Ligand complex. These methods are computationally much cheaper than the WaterMap calculation, and they include all the important terms to the binding affinity. So one possible research direction is to combine the WaterMap calculation with these methods, replacing the corresponding terms which are calculated with large errors in these methods by the WaterMap contribution and keeping the other terms. Recently, there are some effects to combine the WaterMap method with MM/GBSA, and encouraging results are obtained.[184] It might be possible and even easier to combine WaterMap with LIE model, since the direct interaction between the protein and the ligand can be easily estimated by the LIE model, which is missing in WaterMap calculation.

(2) It is known that the protein or the ligand strain free energy is another source of contribution to the binding affinity, and sometimes they affect the binding affinity in a nontrivial way. For example, in Chapter 7, we have shown that the protein strain free energies for the two binding complexes, Thrombin/CDA and Thrombin/CDB, differ by about 0.8 kcal/mol although the structures of the proteins are essentially the same. However, the protein and/or the ligand strain free energy is missed in both the WaterMap calculation and the MM/GBSA or LIE methods. For the WaterMap method to be able to robustly rank order the relative binding affinities among a wider sets of ligands and proteins, it is necessary to develop a method to estimate the protein and the ligand strain free energy. The well known method to estimate the configurational entropy of the macromolecule is the quasi-harmonic approximation, using multidimensional Gaussian distribution to approximate the configurational distribution of the macromolecule.[185] It is possible to use similar kind of technique to approximate the strain free energy of the protein or the ligand.

(3) The WaterMap method estimates the binding affinity based on the free energy difference between water in the binding pocket of protein and that in bulk, so to get a robust estimate of the binding affinity, the potential energy models of the water are critical. Recently, there is some evidence in the literature that whether a water molecule is present or absent in the binding pocket depends on the potential energy models of water used in the simulation.[186] So it is quite possible that the free energy calculated using WaterMap will be different using different models for water. Thus it is necessary to identify a potential

energy model for water that works best for the WaterMap calculation.

Free energy perturbation (FEP) simulations provide one of the most accurate simulation techniques to calculate the Protein-Ligand binding affinities, and it has a potentially large impact on drug design projects, especially in late stage lead optimization cycle. We have shown that the current implementations of FEP simulation methods can be substantially improved in sampling efficiency when the enhanced sampling method is incorporated. We have developed a new enhanced sampling technique called REST (replica exchange with solute tempering), and successfully combined it with an efficient schedule for lambda-hopping FEP (which we call FEP/REST) to solve the sampling problem in brute force FEP simulation. To be specific, by scaling the Hamiltonian of a specific region of interest in the system by a factor smaller than one and run all the replicas on the same temperature using Hamiltonian replica exchange method (HREM), a small number of replicas are sufficient to maintain a large exchange acceptance ratio, and enhanced sampling is achieved through the increased effective temperature of the “hot” region. We have shown that the improved version of REST (which we call REST2) also bypasses the poor scaling with system size of normal TREM (temperature replica exchange method), and it is more efficient than the original REST for sampling systems with large conformational changes. The FEP/REST method doesn’t require the prior known slow degrees of the freedom of the system, and superior convergence of the free energy are demonstrated both by consistency of the results (independence from the starting conformation) and agreement with experimental binding affinity data. We have shown in two cases that the FEP/REST facilitates the sampling of different conformations separated by large energy barriers, one in the protein and the other in the ligand. We expect that this method will demonstrate its ability in a wider set of proteins where the energy barrier separating the relative conformations is large enough to cause sampling problem using brute force FEP.

The FEP/REST protocol provides an efficient way to calculate the relative binding affinity between two ligands and treat local structural reorganization effect. If the energy barriers separating the different conformational states are very high, a very long equilibration time is required to equilibrate the generalized ensemble and to converge the free energy calculation. In addition, to rank order the relative binding affinity between a set of N lig-

and molecules, we need as least $N-1$ simulations, which is computationally very expensive. Here, we suggest possible research directions to further develop the FEP/REST method to overcome these limitations:

(1) In the current implementation of FEP/REST, the “hot” region include the ligands and protein residues surrounding the binding pocket, assuming that the slow degrees of freedom are within this region. It is sufficient to treat localized conformational reorganization during alchemical transformation from one ligand to another, but might not be sufficient for treating delocalized conformational changes (allosteric regulation). A possible procedure to treat this problem is the following: (a) include a larger hot region in a first round FEP/REST simulation, and find those key residues responsible for the allosteric regulation; (b) run a second round FEP/REST just including those key residues in the hot region.

(2) In the Thrombin case shown in Chapter 7, the time required to equilibrate between the two conformations of the ligands is relatively long (about 1.5 ns). For more complicated systems, longer equilibration time might be required. Recently, there is some efforts in the literature to reweigh the configurations sampled in a simulation (not necessarily Boltzmann distribution) to a Boltzmann distribution.[187] With this kind of technique, we can do two short FEP/REST simulations starting from different conformations and then combine the trajectories and reweigh each configuration sampled to estimate the free energy difference. In this way, we don't need to fully equilibrate the generalized ensemble to estimate the free energy.

(3) In the current form of FEP/REST simulation, from each simulation, we can only get the relative binding affinity between two ligands. So a total number of $N-1$ simulations are required to rank order a set of N ligands. In future implementations of FEP/REST, we can set up the mutation path in such a way that all FEP/REST simulations for a set of ligands binding to the same receptor share a common immediate state, and the relative binding affinity is compared to a common immediate state.[188] (A-O-B, C-O-D,... relative binding affinity is compared with a common immediate state O) In this way, from one FEP/REST simulation, we can get relative binding affinities of the two ligands (initial and final states) compared to a common immediate state, and a total number of $N/2$ simulations are sufficient to rank order a set of N ligands.

In summary, while the WaterMap method and the FEP/REST method have demonstrated many important advantages compared with previous existing methods to calculate Protein-Ligand binding affinities, and represent significant breakthrough in this field, a lot of problems remain to be solved before these techniques be robustly and routinely applied in everyday structural based drug design projects.

Part V

Bibliography

Bibliography

- [1] Bruce J. Berne, John D. Weeks, and Ruhong Zhou. Dewetting and hydrophobic interaction in physical and biological systems. *Annu. Rev. Phys. Chem.*, 60:85–103, 2009.
- [2] Robert Abel, Tom Young, Ramy Farid, Bruce J. Berne, and Richard A. Friesner. Role of the active-site solvent in the thermodynamics of factor xa ligand binding. *J. Am. Chem. Soc.*, 130(9):2817–2831, 2 2008.
- [3] Tom Young, Robert Abel, Byungchan Kim, Bruce J. Berne, and Richard A. Friesner. Motifs for molecular recognition exploiting hydrophobic enclosure in protein-ligand binding. *Proc. Nat. Acad. Sci. USA*, 104(3):808–813, MAR 2007.
- [4] T. Lazaridis. Inhomogeneous fluid approach to solvation thermodynamics. 1. theory. *J. Phys. Chem. B*, 102(18):3531–3541, APR 1998.
- [5] T. Lazaridis and M. E. Paulattis. Entropy of hydrophobic hydration - a new statistical mechanical formulation. *J. Phys. Chem.*, 96(9):3847–3855, APR 1992.
- [6] T. Lazaridis and M. Karplus. Orientational correlations and entropy in liquid water. *J. Chem. Phys.*, 105(10):4294–4316, SEP 1996.
- [7] Z. Li and T. Lazaridis. Thermodynamics of buried water clusters at a protein-ligand binding interface. *J. Phys. Chem. B*, 110(3):1464–1475, JAN 2006.
- [8] Z. Li and T. Lazaridis. The effect of water displacement on binding thermodynamics: Concanavalin a. *J. Phys. Chem. B*, 109(1):662–670, JAN 2005.

- [9] J. Zielkiewicz. Two-particle entropy and structural ordering in liquid water. *J. Phys. Chem. B*, 112(26):7810–7815, JUL 2008.
- [10] J. Zielkiewicz. Structural properties of water: Comparison of the spc, spce, tip4p, and tip5p models of water. *J. Chem. Phys.*, 123(10):104501, SEP 2005.
- [11] R. Esposito, F. Saija, A. M. Saitta, and P. V. Giaquinta. Entropy-based measure of structural order in water. *Phys. Rev. E*, 73(4):040502, APR 2006.
- [12] K. A. T. Silverstein, K. A. Dill, and A. D. J. Haymet. Hydrophobicity in a simple model of water: Entropy penalty as a sum of competing terms via full, angular expansion. *J. Chem. Phys.*, 114(14):6303–6314, APR 2001.
- [13] F. Saija, A. M. Saitta, and P. V. Giaquinta. Statistical entropy and density maximum anomaly in liquid water. *J. Chem. Phys.*, 119(7):3587–3589, AUG 2003.
- [14] H. Singh, N. Misra, V. Hnizdo, A. Fedorowicz, and E. Demchuk. Nearest neighbor estimates of entropy. *Am. J. Math. Manag. Sci.*, 23(3/4):301–322, 2003.
- [15] H.S. Green. *Molecular theory of fluids*. Chapter III. North-Holland: Amsterdam, 1952.
- [16] Harold J. Raveché. Entropy and molecular correlation functions in open systems. i. derivation. *J. Chem. Phys.*, 55(DOI 10):2242–2250, SEP 1971.
- [17] Duane C. Wallace. On the role of density fluctuations in the entropy of a fluid. *J. Chem. Phys.*, 87(4):2282–2284, APR 1987.
- [18] C. E. Shannon. A mathematical theory of communication. *Bell. Syst. Tech. J.*, 27(3):379–423, 1948.
- [19] I. Z. Fisher and B. L. Kopeliovich. Improvement of superposition approximation in the theory of liquids. *Dokl. Akad. Nauk SSSR [Sov. Phys. Dokl. 5, 761 (1960)]*., 133(1):81–83, 1960.
- [20] H. Reiss. Superposition approximations from a variation principle. *J. Stat. Phys.*, 6(1):39–47, 1972.

- [21] A. Singer. Maximum entropy formulation of the kirkwood superposition approximation. *J. Chem. Phys.*, 121(8):3657–3666, AUG 2004.
- [22] B. J. Killian, J. Y. Kravitz, and M. K. Gilson. Extraction of configurational entropy from molecular simulations via an expansion approximation. *J. Chem. Phys.*, 127(2):024107, JUL 2007.
- [23] V. Hnizdo, E. Darian, A. Fedorowicz, E. Demchuk, S. Li, and H. Singh. Nearest-neighbor nonparametric method for estimating the configurational entropy of complex molecules. *J. Comput. Chem.*, 28(3):655–668, FEB 2007.
- [24] Hiroyuki Matsuda. Physical nature of higher-order mutual information: Intrinsic correlations and frustration. *Phys. Rev. E*, 62(3):3096–3102, SEP 2000.
- [25] D. O. Loftsgaarden and C. P. Quesenberry. A nonparametric estimate of a multivariate density function. *Ann. Math. Statist.*, 36(3):1049–1051, 1965.
- [26] Sunil Arya and David M. Mount. Approximate nearest neighbor queries in fixed dimensions. In *SODA '93: Proceedings of the fourth annual ACM-SIAM Symposium on Discrete algorithms*, pages 271–280, Philadelphia, PA, USA, 1993. Society for Industrial and Applied Mathematics.
- [27] Jerome H. Freidman, Jon Louis Bentley, and Raphael Ari Finkel. An algorithm for finding best matches in logarithmic expected time. *ACM Trans. Math. Softw.*, 3(3):209–226, SEP 1977.
- [28] D. E. Smith and A. D. Haymet. Free energy, entropy, and internal energy of hydrophobic interactions: Computer simulations. *J. Chem. Phys.*, 98(8):6445–6454, APR 1993.
- [29] S. Z. Wan, R. H. Stote, and M. Karplus. Calculation of the aqueous solvation energy and entropy, as well as free energy, of simple polar solutes. *J. Chem. Phys.*, 121(19):9539–9548, NOV 2004.
- [30] Charles H. Bennett. Efficient estimation of free energy differences from monte carlo data. *J. Comput. Phys.*, 22(2):245–268, 1976.

- [31] Kevin J. Bowers, Edmond Chow, Huafeng Xu, Ron O. Dror, Michael P. Eastwood, Brent A. Gregersen, John L. Klepeis, Istvan Kolossvary, Mark A. Moraes, Federico D. Sacerdoti, John K. Salmon, Yibing Shan, and David E. Shaw. Scalable algorithms for molecular dynamics simulations on commodity clusters. In *SC '06: Proceedings of the 2006 ACM/IEEE conference on Supercomputing*, page 84, New York, NY, USA, 2006. ACM.
- [32] William L. Jorgensen, Jayaraman Chandrasekhar, Jeffrey D. Madura, Roger W. Impey, and Michael L. Klein. Comparison of simple potential functions for simulating liquid water. *J. Chem. Phys.*, 79(10):926–935, 1983.
- [33] Shuichi Nose. A unified formulation of the constant temperature molecular dynamics methods. *J. Chem. Phys.*, 81(1):511–519, JUL 1984.
- [34] William G. Hoover. Canonical dynamics: Equilibrium phase-space distributions. *Phys. Rev. A*, 31(3):1695–1697, 1985.
- [35] Glenn J. Martyna, Douglas J. Tobias, and Michael L. Klein. Constant pressure molecular dynamics algorithms. *J. Chem. Phys.*, 101(10):4177–4189, 1994.
- [36] Tuckerman, B.J. Berne, and G.J. Martyna. Reversible multiple time scale molecular dynamics. *J. Chem. Phys.*, 97(3):1990–2001, AUG 1992.
- [37] Tom Darden, Darrin York, and Lee Pedersen. Particle mesh ewald: An $n \cdot \log(n)$ method for ewald sums in large systems. *J. Chem. Phys.*, 98(10):10089–10092, JUN 1993.
- [38] H.J.C. Berendsen, J.P.M. Postma, W.F. van Gunsteren, and J. Hermans. *Intermolecular Forces*, pages 331–342. Interaction Models for Water in Relation to Protein Hydration. Reidel: Dordrecht, 1981.
- [39] H. J. C. Berendsen, J. R. Grigera, and T. P. Straatsma. The missing term in effective pair potentials. *J. Phys. Chem.*, 91(24):6269–6271, NOV 1987.
- [40] 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(20):9665–9678, MAY 2004.
- [41] R. H. Henchman. Free energy of liquid water from a computer simulation via cell theory. *J. Chem. Phys.*, 126(6):064504, FEB 2007.
- [42] M. R. Shirts and V. S. Pande. Solvation free energies of amino acid side chain analogs for common molecular mechanics water models. *J. Chem. Phys.*, 122(14):134508, JUL 2005.
- [43] Duane C. Wallace. Statistical mechanical theory of liquid entropy. *Int. J. Quantum Chem*, 52(2):425–435, SEP 1994.
- [44] P. V. Giaquinta and G. Giunta. About entropy and correlations in a fluid of hard spheres. *Physica A*, 187(1-2):145–158, AUG 1992.
- [45] W. L. Jorgensen and C. Jenson. Temperature dependence of tip3p, spc, and tip4p water from npt monte carlo simulations: Seeking temperatures of maximum density. *J. Comput. Chem.*, 19(10):1179–1186, JUL 1998.
- [46] L. A. Baez and P. Clancy. Existence of a density maximum in extended simple point charge water. *J. Chem. Phys.*, 101(11):9837–9840, DEC 1994.
- [47] W. Wagner and A. Pruß. The iapws formulation 1995 for the thermodynamic properties of ordinary water substance for general and scientific use. *J. Phys. Chem. Ref. Data*, 31(2):387–478, JUN 2002.
- [48] Michael K. Gilson and Huan-Xiang Zhou. Calculation of protein-ligand binding affinities. *Annu. Rev. Biophys. Biomol. Struct.*, 36(1):21–42, 2007.
- [49] David L. Moble and Ken A. Dill. Binding of small-molecule ligands to proteins: “what you see” is not always “what you get”. *Structure*, 17:489–498, 2009.
- [50] Huan-Xiang Zhou and Michael K. Gilson. Theory of free energy and entropy in noncovalent binding. *Chem. Rev.*, 109(9):4092–4107, 2009. PMID: 19588959.

- [51] Olgun Guvench and Alexander D MacKerell Jr. Computational evaluation of protein-small molecule binding. *Curr. Opin. Struct. Biol.*, 19(1):56–61, 2009.
- [52] Rasmus P. Clausen, Peter Naur, Anders S. Kristensen, Jeremy R. Greenwood, Mette Strange, Hans Braluner-Osborne, Anders A. Jensen, Anne Sophie T. Nielsen, Ulla Geneser, Lone M. Ringgaard, Birgitte Nielsen, Darryl S. Pickering, Lotte Brehm, Michael Gajhede, Povl Krogsgaard-Larsen, and Jette S. Kastrup. The glutamate receptor glur5 agonist (s)-2-amino-3-(3-hydroxy-7,8-dihydro-6h-cyclohepta[d]isoxazol-4-yl)propionic acid and the 8-methyl analogue: Synthesis, molecular pharmacology, and biostructural characterization. *J. Med. Chem.*, 52(15):4911–4922, 2009. PMID: 19588945.
- [53] T. Beuming, R. Farid, and W. Sherman. High-energy water sites determine peptide binding affinity and specificity of pdz domains. *Protein Sci.*, 18:1609–1619, 2009.
- [54] D. Robinson, W. Sherman, and R. Farid. Understanding kinase selectivity through energetic analysis of binding site waters. *ChemMedChem*, 5:618–627, 2010.
- [55] Cristiano R. W. Guimaraes and Alan M. Mathiowetz. Addressing limitations with the mm-gb/sa scoring procedure using the watermap method and free energy perturbation calculations. *J. Chem. Inf. Mod.*, 50(4):547–559, 2010.
- [56] R. A. Pearlstein, Q.Y. Hu, J. Zhou, D. Yowe, J. Levell, B. Dale, V. K. Kaushik, D. Daniels, S. Hanrahan, W. Sherman, and R. Abel. New hypotheses about the structure-function of proprotein convertase subtilisin/kexin type 9: Analysis of the epidermal growth factor-like repeat a docking site using watermap. *Proteins*, 78:2571–2586, 2010.
- [57] Jill E. Chrencik, Akshay Patny, Iris K. Leung, Brian Korniski, Thomas L. Emmons, Troii Hall, Robin A. Weinberg, Jennifer A. Gormley, Jennifer M. Williams, Jacqueline E. Day, Jeffrey L. Hirsch, James R. Kiefer, Joseph W. Leone, H. David Fischer, Cynthia D. Sommers, Horng-Chih Huang, E.J. Jacobsen, Ruth E. Tenbrink, Alfredo G. Tomasselli, and Timothy E. Benson. Structural and thermodynamic char-

- acterization of the tyk2 and jak3 kinase domains in complex with cp-690550 and cmp-6. *J. Mol. Bio.*, 400(3):413 – 433, 2010.
- [58] David L. Mobley, John D. Chodera, and Ken A. Dill. Confine-and-release method: Obtaining correct binding free energies in the presence of protein conformational change. *J. Chem. Theory Comput.*, 3(4):1231–1235, 2007. PMID: 18843379.
- [59] Yuqing Deng and Benoit Roux. Computations of standard binding free energies with molecular dynamics simulations. *J. Phys. Chem. B*, 113(8):2234–2246, 2009.
- [60] Jessica M.J. Swanson, Richard H. Henchman, and J. Andrew McCammon. Revisiting free energy calculations: A theoretical connection to mm/pbsa and direct calculation of the association free energy. *Biophys. J.*, 86(1):67 – 74, 2004.
- [61] Niu Huang, Chakrapani Kalyanaraman, Katarzyna Bernacki, and Matthew P. Jacobson. Molecular mechanics methods for predicting protein-ligand binding. *Phys. Chem. Chem. Phys.*, 8:5166–5177, 2006.
- [62] Cristiano R. W. Guimares and Mario Cardozo. Mm-gb/sa rescoring of docking poses in structure-based lead optimization. *Journal of Chemical Information and Modeling*, 48(5):958–970, 2008. PMID: 18422307.
- [63] David L. Mobley, Christopher I. Bayly, Matthew D. Cooper, Michael R. Shirts, and Ken A. Dill. Small molecule hydration free energies in explicit solvent: An extensive test of fixed-charge atomistic simulations. *J. Chem. Theory Comput.*, 5(2):350–358, 2009.
- [64] E. Gallicchio, M. M. Kubo, and R. M. Levy. Enthalpy-entropy and cavity decomposition of alkane hydration free energies: Numerical results and implications for theories of hydrophobic solvation. *J. Phys. Chem. B*, 104(26):6271–6285, 2000.
- [65] Schrodinger LLC New York NY. Maestro, version 8.5, 2008.
- [66] Lingle Wang, Robert Abel, Richard A. Friesner, and B. J. Berne. Thermodynamic properties of liquid water: An application of a nonparametric approach to computing the entropy of a neat fluid. *J. Chem. Theory Comput.*, 5(6):1462–1473, 2009.

- [67] Thomas C. Beutler, Alan E. Mark, Rene C. van Schaik, Paul R. Gerber, and Wilfred F. van Gunsteren. Avoiding singularities and numerical instabilities in free energy calculations based on molecular simulations. *Chem. Phys. Lett.*, 222(6):529 – 539, 1994.
- [68] Michael L. Connolly. The molecular surface package, version 3.9.3, 2006.
- [69] G. Hummer, S. Garde, A. E. Garcia, M. E. Paulaitis, and L. R. Pratt. Hydrophobic effects on a molecular scale. *J. Phys. Chem. B*, 102(51):10469–10482, 1998.
- [70] Noel T. Southall, Ken A. Dill, and A. D. J. Haymet. A view of the hydrophobic effect. *J. Phys. Chem. B*, 106(3):521–533, 2002.
- [71] Steve W. Homans. Water, water everywhere – except where it matters? *Drug Discovery Today*, 12(13-14):534 – 539, 2007.
- [72] KA Sharp, A Nicholls, RF Fine, and B Honig. Reconciling the magnitude of the microscopic and macroscopic hydrophobic effects. *Science (New York, N. Y.)*, 252(5002), 1991.
- [73] Irina Massova and Peter Kollman. Combined molecular mechanical and continuum solvent approach (mm-pbsa/gbsa) to predict ligand binding. *Perspect. Drug Discovery Des.*, 18:113–135, 2000. 10.1023/A:1008763014207.
- [74] Richard A. Friesner, Robert B. Murphy, Matthew P. Repasky, Leah L. Frye, Jeremy R. Greenwood, Thomas A. Halgren, Paul C. Sanschagrín, and Daniel T. Mainz. Extra precision glide: Docking and scoring incorporating a model of hydrophobic enclosure for protein ligand complexes. *J. Med. Chem.*, 49(21):6177–6196, 2006.
- [75] Tom Young, Lan Hua, Xuhui Huang, Robert Abel, Richard A. Friesner, and B. J. Berne. Dewetting transitions in protein cavities. *Proteins*, 78:1856–1869, 2010.
- [76] Richard Malham, Sarah Johnstone, Richard J. Bingham, Elizabeth Barratt, Simon E. V. Phillips, Charles A. Laughton, and Steve W. Homans. Strong solute-solute dispersive interactions in a protein-ligand complex. *J. Am. Chem. Soc.*, 127(48):17061–17067, 2005.

- [77] Tomi T. Airene, Heidi Kidron, Yvonne Nymalm, Matts Nylund, Gun West, Peter Mattjus, and Tiina A. Salminen. Structural evidence for adaptive ligand binding of glycolipid transfer protein. *J. Mol. Bio.*, 355(2):224–236, 2006.
- [78] Paula I. Lario, Richard A. Pfuetzner, Elizabeth A. Frey, Louise Creagh, Charles Haynes, Anthony T. Maurelli, and Natalie C. J. Strynadka. Structure and biochemical analysis of a secretin pilot protein. *EMBO J. (European Molecular Biology Organization)*, 24(1):1111–1121, 2005.
- [79] Elizabeth Barratt, Richard J. Bingham, Daniel J. Warner, Charles A. Laughton, Simon E. V. Phillips, and Steve W. Homans. Van der waals interactions dominate ligand-protein association in a protein binding site occluded from solvent water. *J. Am. Chem. Soc.*, 127(33):11827–11834, 2005.
- [80] Julien Michel, Julian Tirado-Rives, and William L. Jorgensen. Prediction of the water content in protein binding sites. *J. Phys. Chem. B*, 113(40):13337–13346, 2009.
- [81] Richard J. Bingham, John B. C. Findlay, Shih-Yang Hsieh, Arnout P. Kalverda, Alexandra Kjellberg, Chiara Perazzolo, Simon E. V. Phillips, Kothandaraman Seshadri, Chi H. Trinh, W. Bruce Turnbull, Geoffrey Bodenhausen, and Steve W. Homans. Thermodynamics of binding of 2-methoxy-3-isopropylpyrazine and 2-methoxy-3-isobutylpyrazine to the major urinary protein. *J. Am. Chem. Soc.*, 126(6):1675–1681, 2004.
- [82] David E. Timm, L.J. Baker, Heather Mueller, Lukas Zidek, and Milos V. Novotny. Structural basis of pheromone binding to mouse major urinary protein (mup-i). *Protein Sci.*, 10:997–1004, 2001.
- [83] Scott D. Sharrow, Milos V. Novotny, and Martin J. Stone. Thermodynamic analysis of binding between mouse major urinary protein-I and the pheromone 2-sec-butyl-4,5-dihydrothiazole. *Biochemistry*, 42(20):6302–6309, 2003.
- [84] Lingle Wang, Richard A. Friesner, and B. J. Berne. Hydrophobic interactions in model enclosures from small to large length scales: non-additivity in explicit and implicit solvent models. *Faraday Discuss.*, 146:246–262, 2010.

- [85] Xavier Siebert and Gerhard Hummer. Hydrophobicity maps of the n-peptide coiled coil of hiv-1 gp41. *Biochemistry*, 41(9):2956–2961, 2002.
- [86] L. L. C. Schrodinger. Schrodinger suite 2010 protein preparation wizard, 2010. New York, NY.
- [87] G. Hummer, S. Garde, A. E. Garcia, A. Pohorille, and L. R. Pratt. An information theory model of hydrophobic interactions. *Proc. Nat. Acad. Sci. USA*, 93(17):8951–8955, 1996.
- [88] Robert H. Swendsen and Jian Sheng Wang. Replica monte carlo simulation of spin-glasses. *Phys. Rev. Lett.*, 57(21):2607–2609, 1986.
- [89] Koji Hukushima and Koji Nemoto. Exchange monte carlo method and application to spin glass simulations. *J. Phys. Soc. Jpn.*, 65(6):1604–1608, 1996.
- [90] Angel E. Garcia and Kevin Y. Sanbonmatsu. Exploring the energy landscape of a β hairpin in explicit solvent. *Proteins Struct. Funct. Bioinf.*, 42(3):345–354, 2001.
- [91] Ruhong Zhou, Bruce J. Berne, and Robert Germain. The free energy landscape for β hairpin folding in explicit water. *Proc. Nat. Acad. Sci. USA*, 98(26):14931–14936, 2001.
- [92] Young Min Rhee and Vijay S. Pande. Multiplexed-replica exchange molecular dynamics method for protein folding simulation. *Biophys. J.*, 84(2):775–786, 2003.
- [93] Weihua Zheng, Michael Andrec, Emilio Gallicchio, and Ronald M. Levy. Simulating replica exchange simulations of protein folding with a kinetic network model. *Proc. Nat. Acad. Sci. USA*, 104(39):15340–15345, 2007.
- [94] Pu Liu, Byungchan Kim, Richard A. Friesner, and B. J. Berne. Replica exchange with solute tempering: A method for sampling biological systems in explicit water. *Proc. Nat. Acad. Sci. USA*, 102(39):13749–13754, 2005.
- [95] Xuhui Huang, Morten Hagen, Byungchan Kim, Richard A. Friesner, Ruhong Zhou, and B. J. Berne. Replica exchange with solute tempering: Efficiency in large scale systems. *J. Phys. Chem. B*, 111(19):5405–5410, 2007.

- [96] Samuel L. C. Moors, Servaas Michielssens, and Arnout Ceulemans. Improved replica exchange method for native-state protein sampling. *J. Chem. Theory Comput.*, 7(1):231–237, 2011.
- [97] Tsuyoshi Terakawa, Tomoshi Kameda, and Shoji Takada. On easy implementation of a variant of the replica exchange with solute tempering in gromacs. *J. Comput. Chem.*, 32(7):1228–1234, 2011.
- [98] Hiroaki Fukunishi, Osamu Watanabe, and Shoji Takada. On the hamiltonian replica exchange method for efficient sampling of biomolecular systems: Application to protein structure prediction. *J. Chem. Phys.*, 116(20):9058–9067, 2002.
- [99] Roman Affentranger, Ivano Tavernelli, and Ernesto E. Di Iorio. A novel hamiltonian replica exchange md protocol to enhance protein conformational space sampling. *J. Chem. Theory Comput.*, 2(2):217–228, 2006.
- [100] George A. Kaminski, Richard A. Friesner, Julian Tirado-Rives, and William L. Jorgensen. Evaluation and reparametrization of the opl-aa force field for proteins via comparison with accurate quantum chemical calculations on peptides. *J. Phys. Chem. B*, 105(28):6474–6487, 2001.
- [101] Jonathan W. Neidigh, R. Matthew Fesinmeyer, and Niels H. Andersen. Designing a 20-residue protein. *Nat. Struct. Mol. Biol.*, 9(6):425, 2002.
- [102] William L. Jorgensen. Efficient drug lead discovery and optimization. *Acc. Chem. Res.*, 42(6):724–733, 2009. PMID: 19317443.
- [103] Emilio Gallicchio and Ronald M Levy. Advances in all atom sampling methods for modeling protein-ligand binding affinities. *Curr. Opin. Struct. Biol.*, 21(2):161–166, 2011.
- [104] John D. Chodera, David L. Mobley, Michael R. Shirts, Richard W. Dixon, Kim Branson, and Vijay S. Pande. Alchemical free energy methods for drug discovery: progress and challenges. *Curr. Opin. Struct. Biol.*, 21(2):150–160, 2011.

- [105] W. L. Jorgensen and Tirado J. Rives. The opls potential functions for proteins - energy minimizations for crystals of cyclic-peptides and crambin. *J. Am. Chem. Soc.*, 110(6):1657–1666, 1988.
- [106] A. D. MacKerell, D. Bashford, Bellott, R. L. Dunbrack, J. D. Evanseck, M. J. Field, S. Fischer, J. Gao, H. Guo, S. Ha, D. Joseph-McCarthy, L. Kuchnir, K. Kuczera, F. T. K. Lau, C. Mattos, S. Michnick, T. Ngo, D. T. Nguyen, B. Prodhom, W. E. Reiher, B. Roux, M. Schlenkrich, J. C. Smith, R. Stote, J. Straub, M. Watanabe, J. Wiorkiewicz-Kuczera, D. Yin, and M. Karplus. All-atom empirical potential for molecular modeling and dynamics studies of proteins. *J. Phys. Chem. B*, 102(18):3586–3616, 1998.
- [107] Junmei Wang, Romain M. Wolf, James W. Caldwell, Peter A. Kollman, and David A. Case. Development and testing of a general amber force field. *J. Comput. Chem.*, 25(9):1157–1174, 2004.
- [108] Jay L. Banks, Hege S. Beard, Yixiang Cao, Art E. Cho, Wolfgang Damm, Ramy Farid, Anthony K. Felts, Thomas A. Halgren, Daniel T. Mainz, Jon R. Maple, Robert Murphy, Dean M. Philipp, Matthew P. Repasky, Linda Y. Zhang, Bruce J. Berne, Richard A. Friesner, Emilio Gallicchio, and Ronald M. Levy. Integrated modeling program, applied chemical theory (impact). *J. Comput. Chem.*, 26(16):1752–1780, 2005.
- [109] Wei Jiang and Benoit Roux. Free energy perturbation hamiltonian replica-exchange molecular dynamics (fep/h-remd) for absolute ligand binding free energy calculations. *J. Chem. Theory Comput.*, 6(9):2559–2565, 2010.
- [110] Julien Michel, Julian Tirado-Rives, and William L. Jorgensen. Energetics of displacing water molecules from protein binding sites: Consequences for ligand optimization. *J. Am. Chem. Soc.*, 131(42):15403–15411, 2009. PMID: 19778066.
- [111] Lingle Wang, B. J. Berne, and R. A. Friesner. Ligand binding to protein-binding pockets with wet and dry regions. *Proc. Nat. Acad. Sci. USA*, 108(4):1326–1330, 2011.

- [112] Lingle Wang, Richard A. Friesner, and B. J. Berne. Replica exchange with solute scaling: A more efficient version of replica exchange with solute tempering (rest2). *J. Phys. Chem. B*, 115(30):9431–9438, 2011.
- [113] Andrew Morton and Brian W. Matthews. Specificity of ligand binding in a buried nonpolar cavity of t4 lysozyme: Linkage of dynamics and structural plasticity. *Biochemistry*, 34(27):8576–8588, 1995. PMID: 7612599.
- [114] Andrew Morton, Walter A. Baase, and Brian W. Matthews. Energetic origins of specificity of ligand binding in an interior nonpolar cavity of t4 lysozyme. *Biochemistry*, 34(27):8564–8575, 1995. PMID: 7612598.
- [115] William C. Lumma, Keith M. Witherup, Thomas J. Tucker, Steven F. Brady, John T. Sisko, Adel M. Naylor-Olsen, S. Dale Lewis, Bobby J. Lucas, and Joseph P. Vacca. Design of novel, potent, noncovalent inhibitors of thrombin with nonbasic p-1 substructures: Rapid structure-activity studies by solid-phase synthesis. *J. Med. Chem.*, 41(7):1011–1013, 1998.
- [116] Shaun R. Coughlin. Thrombin signalling and protease-activated receptors. *Nature*, 407(14):258–264, 2000.
- [117] Christopher S. Burgey, Kyle A. Robinson, Terry A. Lyle, Philip E. J. Sanderson, S. Dale Lewis, Bobby J. Lucas, Julie A. Krueger, Rominder Singh, Cynthia Miller-Stein, Rebecca B. White, Bradley Wong, Elizabeth A. Lyle, Peter D. Williams, Craig A. Coburn, Bruce D. Dorsey, James C. Barrow, Maria T. Stranieri, Marie A. Holahan, Gary R. Sitko, Jacquelynn J. Cook, Daniel R. McMasters, Colleen M. McDonough, William M. Sanders, Audrey A. Wallace, Franklin C. Clayton, Dennis Bohn, Yvonne M. Leonard, Theodore J. Detwiler, Joseph J. Lynch, Youwei Yan, Zhongguo Chen, Lawrence Kuo, Stephen J. Gardell, Jules A. Shafer, and Joseph P. Vacca. Metabolism-directed optimization of 3-aminopyrazinone acetamide thrombin inhibitors. development of an orally bioavailable series containing p1 and p3 pyridines. *J. Med. Chem.*, 46(4):461–473, 2003.

- [118] Yuqing Deng and Benoit Roux. Calculation of standard binding free energies: Aromatic molecules in the t4 lysozyme 199a mutant. *J. Chem. Theory Comput.*, 2(5):1255–1273, 2006.
- [119] Emilio Gallicchio, Mauro Lapelosa, and Ronald M. Levy. Binding energy distribution analysis method (bedam) for estimation of protein-ligand binding affinities. *J. Chem. Theory Comput.*, 6(9):2961–2977, 2010.
- [120] Donghong Min, Hongzhi Li, Guohui Li, Ryan Bitetti-Putzer, and Wei Yang. Synergistic approach to improve alchemical free energy calculation in rugged energy surface. *J. Chem. Phys.*, 126(14):144109, 2007.
- [121] Ed Harder, Jon Maple, Wolfgang Damm, Mark Reboul, Jeremy Greenwood, and R. A. Friesner. Developmet and testing of the OPLS 2.0 force field. Manuscript in Prep., 2011.
- [122] Woody Sherman, Tyler Day, Matthew P. Jacobson, Richard A. Friesner, and Ramy Farid. Novel procedure for modeling ligand/receptor induced fit effects. *J. Med. Chem.*, 49(2):534–553, 2006.
- [123] Andrew Pohorille, Christopher Jarzynski, and Christophe Chipot. Good practices in free-energy calculations. *J. Phys. Chem. B*, 114(32):10235–10253, 2010.
- [124] W. Kauzmann. Some factors in the interpretation of protein denaturation. *Adv. Prot. Chem.*, 14:1–63, 1959.
- [125] Ken A. Dill. Dominant forces in protein folding. *Biochemistry*, 29:7133–7155, 1990.
- [126] S. Shimizu and H. S. Chan. Anti-cooperativity and cooperativity in hydrophobic interactions: Three-body free energy landscapes and comparison with implicit-solvent potential functions for proteins. *Proteins: Struct. Funct. Bioinf.*, 48(1):15–30, 2002.
- [127] X. Li and J. Liang. Geometric cooperativity and anticooperativity of three-body interactions in native proteins. *Proteins: Struct. Funct. Bioinf.*, 60(1):46–65, 2005.

- [128] George Némethy and Harold A. Scheraga. The structure of water and hydrophobic bonding in proteins. iii. the thermodynamic properties of hydrophobic bonds in proteins^{1,2}. *J. Phys. Chem.*, 66(10):1773–1789, 2002.
- [129] F. Brugè, S. L. Fornili, G. G. Malenkov, M. B. Palma-Vittorelli, and M. U. Palma. Solvent-induced forces on a molecular scale: non-additivity, modulation and causal relation to hydration. *Chem. Phys. Lett.*, 254(5-6):283 – 291, 1996.
- [130] Jeffrey A. Rank and David Baker. A desolvation barrier to hydrophobic cluster formation may contribute to the rate-limiting step in protein folding. *Protein Sci.*, 6(2):347–354, 1997.
- [131] Seishi Shimizu and Hue Sun Chan. Anti-cooperativity in hydrophobic interactions: A simulation study of spatial dependence of three-body effects and beyond. *J. Chem. Phys.*, 115(3):1414–1421, 2001.
- [132] Cezary Czaplewski, Sylwia Rodziewicz-Motowidło, Adam Liwo, Daniel R. Ripoll, Ryszard J. Wawak, and Harold A. Scheraga. Molecular simulation study of cooperativity in hydrophobic association. *Protein Sci.*, 9(6):1235–1245, 2000.
- [133] Cezary Czaplewski, Daniel R. Ripoll, Adam Liwo, Sylwia Rodziewicz-Motowidło, Ryszard J. Wawak, and Harold A. Scheraga. Can cooperativity in hydrophobic association be reproduced correctly by implicit solvation models? *Int. J. Quantum Chem.*, 88(1):41–55, 2002.
- [134] Cezary Czaplewski, Adam Liwo, Daniel R. Ripoll, and Harold A. Scheraga. Molecular origin of anticooperativity in hydrophobic association. *J. Phys. Chem. B*, 109(16):8108–8119, 2005.
- [135] Cezary Czaplewski, Sylwia Rodziewicz-Motowidło, Adam Liwo, Daniel R. Ripoll, Ryszard J. Wawak, and Harold A. Scheraga. Comment on “anti-cooperativity in hydrophobic interactions: A simulation study of spatial dependence of three-body effects and beyond” [j. chem. phys. [bold 115], 1414 (2001)]. *J. Chem. Phys.*, 116(6):2665–2667, 2002.

- [136] Seishi Shimizu and Hue Sun Chan. Reply to “comment on ‘anti-cooperativity in hydrophobic interactions: A simulation study of spatial dependence of three-body effects and beyond’ ” [j. chem. phys. [bold 116], 2665 (2002)]. *J. Chem. Phys.*, 116(6):2668–2669, 2002.
- [137] Seishi Shimizu, Maria Sabaye Moghaddam, and Hue Sun Chan. Comment on “molecular origin of anticooperativity in hydrophobic association”. *J. Phys. Chem. B*, 109(44):21220–21221, 2005.
- [138] Cezary Czaplewski, Sebastian Kalinowski, Adam Liwo, Daniel R. Ripoll, and Harold A. Scheraga. Reply to “comment on ‘molecular origin of anticooperativity in hydrophobic association’ ”. *J. Phys. Chem. B*, 109(44):21222–21224, 2005.
- [139] Cezary Czaplewski, Sebastian Kalinowski, Adam Liwo, and Harold A. Scheraga. Comparison of two approaches to potential of mean force calculations of hydrophobic association: particle insertion and weighted histogram analysis methods. *Mol. Phys.*, 103(21-23):3153 – 3167, 2005.
- [140] Seishi Shimizu and Hue Sun Chan. Temperature dependence of hydrophobic interactions: A mean force perspective, effects of water density, and nonadditivity of thermodynamic signatures. *J. Chem. Phys.*, 113(11):4683–4700, 2000.
- [141] Maria Sabaye Moghaddam, Seishi Shimizu, and Hue Sun Chan. Temperature dependence of three-body hydrophobic interactions: Potential of mean force, enthalpy, entropy, heat capacity, and nonadditivity. *J. Am. Chem. Soc.*, 127(1):303–316, 2005.
- [142] T. Ghosh, A. E.Garcia, and S. Garde. Water-mediated three-particle interactions between hydrophobic solutes: Size, pressure, and salt effects. *J. Phys. Chem. B*, 107(2):612–617, 2002.
- [143] Tuhin Ghosh, Amrit Kalra, and Shekhar Garde. On the salt induced stabilization of pair and many-body hydrophobic interactions. *J. Phys. Chem. B*, 109(1):642–651, 2004.

- [144] Cezary Czaplewski, Sylwia Rodziewicz-Motowidlo, Magdalena Dabal, Adam Liwo, Daniel R. Ripoll, and Harold A. Scheraga. Molecular simulation study of cooperativity in hydrophobic association: clusters of four hydrophobic particles. *Biophys. Chem.*, 105(2-3):339 – 359, 2003. Walter Kauzmann s 85th Birthday.
- [145] D. A. Mcquarrie. *Statistical Mechanics*. Chapter XIII. Sausalito, CA: University Science Books, 2000.
- [146] John G. Kirkwood. Statistical mechanics of fluid mixtures. *J. Chem. Phys.*, 3:300–313, 1935.
- [147] John G. Kirkwood. Molecular distribution in liquids. *J. Chem. Phys.*, 7:919–925, 1939.
- [148] John G. Kirkwood and Elizabeth Monroe Boggs. The radial distribution function in liquids. *J. Chem. Phys.*, 10:394–402, 1942.
- [149] William L. Jorgensen, Jeffry D. Madura, and Carol J. Swenson. Optimized intermolecular potential functions for liquid hydrocarbons. *J. Am. Chem. Soc.*, 106(22):6638–6646, 2002.
- [150] Benoit Roux and Thomas Simonson. Implicit solvent models. *Biophys. Chem.*, 78(1-2):1 – 20, 1999.
- [151] B. Lee and F. M. Richards. The interpretation of protein structures: Estimation of static accessibility. *J. Mol. Biol.*, 55(3):379 – 380, 1971.
- [152] M. L. Connolly. Solvent-accessible surfaces of proteins and nucleic acids. *Science*, 221(4612):709–713, 1983.
- [153] David J. Tannor, Bryan Marten, Robert Murphy, Richard A. Friesner, Doree Sitkoff, Anthony Nicholls, Barry Honig, Murco Ringnalda, and William A. Goddard. Accurate first principles calculation of molecular charge distributions and solvation energies from ab initio quantum mechanics and continuum dielectric theory. *J. Am. Chem. Soc.*, 116:11875–11882, 1994.

- [154] Ka Lum, David Chandler, and John D. Weeks. Hydrophobicity at small and large length scales. *J. Phys. Chem. B*, 103:4570–4577, 1998.
- [155] Robert Abel, Lingle Wang, Richard A. Friesner, and B. J. Berne. A displaced-solvent functional analysis of model hydrophobic enclosures. *J. Chem. Theory Comput.*, 6:2924–2934, 2010.
- [156] C. Clementi, H. Nymeyer, and J. N. Onuchic. Topological and energetic factors: What determines the structural details of the transition state ensemble and “en-route” intermediates for protein folding? an investigation for small globular proteins. *J. Mol. Biol.*, 298(5):937–953, 2000.
- [157] Frauke Gräter, Pascal Heider, Ronen Zangi, and B. J. Berne. Dissecting entropic coiling and poor solvent effects in protein collapse. *J. Am. Chem. Soc.*, 116:11578–11579, 2008.
- [158] Rahul Godawat, Sumanth N Jamadagni, and Shekhar Garde. Characterizing hydrophobicity of interfaces by using cavity formation, solute binding, and water correlations. *Proc. Nat. Acad. Sci. U.S.A.*, 106(36):15119–15124, 2009.
- [159] Nicolas Giovambattista, Pablo G. Debenedetti, and Peter J. Rossky. Hydration behavior under confinement by nanoscale surfaces with patterned hydrophobicity and hydrophilicity. *J. Phys. Chem. C*, 111:1323–1332, 2007.
- [160] A. Wallqvist and B. J. Berne. Computer simulation of hydrophobic hydration forces on stacked plates at short range. *J. Phys. Chem.*, 99(10):2893–2899, 1995.
- [161] Ruhong Zhou, Xuhui Huang, Claudio J. Margulis, and Bruce J. Berne. Hydrophobic collapse in multidomain protein folding. *Science*, 305(5690):1605 – 1609, 2004.
- [162] Pu Liu, Xuhui Huang, Ruhong Zhou, and Bruce J. Berne. Observation of a dewetting transition in the collapse of the melittin tetramer. *Nature*, 437:159 – 162, 2005.
- [163] J. Dzubiella and J.-P. Hansen. Reduction of the hydrophobic attraction between charged solutes in water. *J. Chem. Phys.*, 119(23):12049 – 12052, 2003.

- [164] J. Dzubiella and J.-P. Hansen. Competition of hydrophobic and coulombic interactions between nanosized solutes. *J. Chem. Phys.*, 121(11):5514 – 5530, 2004.
- [165] Alexander A. Rashin and Barry Honig. Reevaluation of the born model of ion hydration. *J. Phys. Chem.*, 89(26):5588–5593, 1985.
- [166] B. Jayaram, Richard Fine, Kim Sharp, and Barry Honig. Free energy calculations of ion hydration: an analysis of the born model in terms of microscopic simulations. *J. Phys. Chem.*, 93(10):4320–4327, 1989.
- [167] Ronald M. Levy, Mahfoud Belhadj, and Douglas B. Kitchen. Gaussian fluctuation formula for electrostatic free-energy changes in solution. *J. Phys. Chem.*, 95(5):3627–3633, 1991.
- [168] Francisco Figueirido, Gabriela S. Del Buono, and Ronald M. Levy. Molecular mechanics and electrostatic effects. *Biophys. Chem.*, 51(2-3):235–241, 1994.
- [169] Steven W. Rick and Bruce J. Berne. The aqueous solvation of water: A comparison of continuum methods with molecular dynamics. *J. Am. Chem. Soc.*, 116:3949–3954, 1994.
- [170] Gerhard Hummer, Lawrence R. Pratt, and Angel E. Garcia. Free energy of ionic hydration. *J. Phys. Chem.*, 100(10):1206–1215, 1996.
- [171] Gerhard Hummer, Lawrence R. Pratt, and Angel E. Garcia. Molecular theories and simulation of ions and polar molecules in water. *J. Phys. Chem. A*, 102(10):7885–7895, 1998.
- [172] Subramanian Vaitheeswaran, Hao Yin, and Jayendran C. Rasaiah. Water between plates in the presence of an electric field in an open system. *J. Phys. Chem. B*, 109(10):6629–6635, 2005.
- [173] Nicolas Giovambattista, Pablo G. Debenedetti, and Peter J. Rossky. Enhanced surface hydrophobicity by coupling of surface polarity and topography. *Proc. Nat. Acad. Sci. U.S.A.*, 106(36):15181–15185, 2009.

- [174] Ronen Zangi, Morten Hagen, and B. J. Berne. Effect of ions on the hydrophobic interaction between two plates. *J. Am. Chem. Soc.*, 129(10):4678–4686, 2007.
- [175] M. Born. Volumes and heats of hydration of ions. *Z. Physik*, 1(1):45–48, 1920.
- [176] W. Clark Still, Anna Tempezyk, Ronald C. Hawley, and Thomas Hendrickson. Semi-analytical treatment of solvation for molecular mechanics and dynamics. *J. Am. Chem. Soc.*, 112(10):6127–6129, 1990.
- [177] John David Jackson. *Classical Electrodynamics*. John Wiley and Sons: New York, 1962.
- [178] W. Rocchia, E. Alexov, and B. Honig. Extending the applicability of the nonlinear poisson-boltzmann equation: Multiple dielectric constants and multivalent ions. *J. Phys. Chem. B*, 105(28):6507–6514, 2001.
- [179] Jayaraman Chandrasekhar, David C. Spellmeyer, and William L. Jorgensen. Energy component analysis for dilute aqueous solutions of lithium(1+), sodium(1+), fluoride(1-), and chloride(1-) ions. *J. Am. Chem. Soc.*, 106(4):903–910, 1984.
- [180] Mafalda Nina, Dmitri Beglov, and Benoit Roux. Atomic radii for continuum electrostatics calculations based on molecular dynamics free energy simulations. *J. Phys. Chem. B*, 101(26):5239–5248, 1997.
- [181] David M. Huang and David Chandler. The hydrophobic effect and the influence of solute solvent attractions. *J. Phys. Chem. B*, 106(8):2047–2053, 2002.
- [182] J. Dzubiella, J. M. J. Swanson, and J. A. McCammon. Coupling hydrophobicity, dispersion, and electrostatics in continuum solvent models. *Phys. Rev. Lett.*, 96(8):087802, 2006.
- [183] J. Dzubiella, J. M. J. Swanson, and J. A. McCammon. Coupling nonpolar and polar solvation free energies in implicit solvent models. *J. Chem. Phys.*, 124(8):084905, 2006.

- [184] Robert Abel, Noeris K. Salam, John Shelley, Ramy Farid, Richard A. Friesner, and Woody Sherman. Contribution of explicit solvent effects to the binding affinity of small-molecule inhibitors in blood coagulation factor serine proteases. *ChemMedChem*, 6(6):1049–1066, 2011.
- [185] Martin Karplus and Joseph N. Kushick. Method for estimating the configurational entropy of macromolecules. *Macromolecules*, 14(2):325–332, 1981.
- [186] Elisa Fadda and Robert J. Woods. On the role of water models in quantifying the binding free energy of highly conserved water molecules in proteins: The case of concanavalin a. *J. Chem. Theory Comput.*, 7(10):3391–3398, 2011.
- [187] F. Marty Ytreberg and Daniel M. Zuckerman. A black-box re-weighting analysis can correct flawed simulation data. *Proceedings of the National Academy of Sciences*, 105(23):7982–7987, 2008.
- [188] Ilja V. Khavrutskii and Anders Wallqvist. Computing relative free energies of solvation using single reference thermodynamic integration augmented with hamiltonian replica exchange. *Journal of Chemical Theory and Computation*, 6(11):3427–3441, 2010.
- [189] H. W. Horn, W. C. Swope, and J. W. Pitera. Characterization of the tip4p-ew water model: Vapor pressure and boiling point. *J. Chem. Phys.*, 123(19):194504, NOV 2005.
- [190] Shobana, Benoit Roux, and Olaf S. Andersen. Free energy simulations: Thermodynamic reversibility and variability. *J. Phys. Chem. B*, 104(21):5179–5190, 2000.

Part VI

Appendices

Appendix A

Error analysis in NN method and Constant pressure correction

A.1 $Var[\ln f(x)]$ for some special cases

A.1.1 Gaussian distribution

Assume the probability distribution is a Gaussian distribution with average u and variance σ^2 , $N(u, \sigma^2)$. That is

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x-u)^2}{2\sigma^2}\right) \quad (\text{A.1})$$

then

$$\ln f(x) = -\frac{(x-u)^2}{2\sigma^2} - \ln(\sqrt{2\pi}\sigma) \quad (\text{A.2})$$

$$Var[\ln f(x)] = \frac{Var[(x-u)^2]}{4\sigma^4} \quad (\text{A.3})$$

$$= \frac{1}{4\sigma^4} [E(x-u)^4 - (E(x-u)^2)^2] \quad (\text{A.4})$$

$$= \frac{1}{4\sigma^4} (3\sigma^4 - (\sigma^2)^2) \quad (\text{A.5})$$

$$= \frac{1}{2} \quad (\text{A.6})$$

So for Gaussian distribution function, the term $Var[\ln f(x)]$ is a constant $\frac{1}{2}$.

A.1.2 exponential distribution

Assume the probability distribution is an exponential distribution with parameter λ

$$f(x) = \begin{cases} \lambda \exp(-\lambda x) & x \geq 0 \\ 0 & x \leq 0 \end{cases} \quad (\text{A.7})$$

then

$$\ln f(x) = -\lambda x + \ln \lambda \quad (\text{A.8})$$

$$\text{Var}[\ln f(x)] = (\lambda^2) \text{Var}[x] \quad (\text{A.9})$$

$$= 1 \quad (\text{A.10})$$

So for exponential distribution function, the term $\text{Var}[\ln f(x)]$ is a constant 1.

A.1.3 exponential distribution in a finite range

Assume the probability is non-zero only in the range $[0, \alpha]$, (which is the case for real systems) and the distribution is exponential in this range. That is

$$f(x) = \begin{cases} b \exp(-ax) & x \in [0, \alpha] \\ 0 & x \notin [0, \alpha] \end{cases} \quad (\text{A.11})$$

where b is a normalization factor, which is equal to $(\frac{a}{1-\exp(-\alpha a)})$. Using the same procedure as above, we got the variance of $\ln f(x)$:

$$\begin{aligned} \text{Var}[\ln f(x)] &= \frac{1}{1 - \exp(-\alpha a)} [2 - \exp(-\alpha a)((\alpha a)^2 + 2\alpha a + 2)] \\ &\quad - \frac{1}{(1 - \exp(-\alpha a))^2} [1 - \exp(-\alpha a)(\alpha a + 1)]^2 \end{aligned} \quad (\text{A.12})$$

So the variance is in the range $[0, 1]$, and increases as αa increases. When αa goes to infinity, $\text{Var}[\ln f(x)]$ goes to 1, which is the case discussed in the previous section.

A.1.4 linear distribution in a finite range

Assume the distribution is a linear function in the range $[0, \alpha]$, and zero otherwise. That is

$$f(x) = \begin{cases} ax & x \in [0, \alpha] \\ 0 & x \notin [0, \alpha] \end{cases} \quad (\text{A.13})$$

where $a = \frac{2}{\alpha^2}$ is a normalization factor. Using the same trick, we get

$$E[\ln f(x)] = \ln \frac{2}{\alpha} - \frac{1}{2} \quad (\text{A.14})$$

$$E[(\ln f(x))^2] = (\ln \frac{2}{\alpha})^2 - \ln \frac{2}{\alpha} + \frac{1}{2} \quad (\text{A.15})$$

$$\text{Var}[\ln f(x)] = \frac{1}{4} \quad (\text{A.16})$$

So the variance is a constant $\frac{1}{4}$ for linear distributions with $f(x) = 0$ at one of the boundary. It is easy to show that for linear distributions in the range $[0, \alpha]$ with nonzero probability at the boundary point, ($f(x) \neq 0$ at the boundary points) the variance of $\ln f(x)$ is in the range $[0, \frac{1}{4}]$.

A.2 Determination of most proper weights

Given that x_1, x_2, \dots, x_n are independent variables with the same average u but different variance v_1, v_2, \dots, v_n , we may define $\bar{x} = \sum_{i=1}^n w_i x_i$, with constraint $\sum_{i=1}^n w_i = 1$. We may find the weights w_i such that the variance of \bar{x} is minimized:

$$\text{Var}[\bar{x}] = \sum_{i=1}^n (w_i)^2 v_i \quad (\text{A.17})$$

Using Lagrange multipliers we find:

$$w_i = \frac{\frac{1}{v_i}}{\sum_{i=1}^n \frac{1}{v_i}} \quad (\text{A.18})$$

and

$$\text{Var}[\bar{x}] = \frac{1}{\sum_{i=1}^n \frac{1}{v_i}} \quad (\text{A.19})$$

$$E\left[\sum_{i=1}^n w_i (x_i - \bar{x})^2\right] = E\left[\sum_{i=1}^n w_i ((x_i - u) - (\bar{x} - u))^2\right] \quad (\text{A.20})$$

$$= E\left[\sum_{i=1}^n w_i ((x_i - u)^2 - 2(x_i - u)(\bar{x} - u) + (\bar{x} - u)^2)\right] \quad (\text{A.21})$$

$$\begin{aligned} &= E\left[\sum_{i=1}^n w_i (x_i - u)^2\right] - 2E\left[\sum_{i=1}^n w_i (x_i - u)(\bar{x} - u)\right] \\ &\quad + E\left[\sum_{i=1}^n w_i (\bar{x} - u)^2\right] \end{aligned} \quad (\text{A.22})$$

By application of equation A.18 and $\sum_{i=1}^n w_i = 1$, we find:

$$E\left[\sum_{i=1}^n w_i (x_i - \bar{x})^2\right] = \frac{n-1}{\sum_{i=1}^n \frac{1}{v_i}} \quad (\text{A.23})$$

Thus we can approximate the variance of the weighted average by the estimator:

$$V = \frac{1}{n-1} \sum_{i=1}^n w_i (x_i - \bar{x})^2 \quad (\text{A.24})$$

A.3 Constant pressure correction to ΔG_{sim} for the FD entropy

In the FEP simulations, we turned on/off the interaction between one distinguished water molecule with the rest of the system at constant temperature T and constant pressure P_0 , over the series of several λ windows. The solvation free energy of the distinguished water molecule corresponds to the difference in the chemical potential μ between two phases: (1) the liquid phase, and (2) the ideal gas phase with the same temperature and number density as the liquid.[189] Ergo,

$$\Delta G_{sim}(T) = -kT \ln \frac{\tilde{\Delta}(\lambda=1)}{\tilde{\Delta}(\lambda=0)} = \mu_l(N, P_0, T) - \mu_g(N, P^*, T) \quad (\text{A.25})$$

where P^* is the pressure of the ideal gas with the same temperature T and number density as the simulated liquid at pressure P_0 , and $\tilde{\Delta}$ is the isobaric-isothermal partition function of the system specified by lambda. (For details, please see reference [189].)

The heat capacity of the ideal gas at constant pressure P^* is trivially constant with respect to temperature, and we may well approximate the heat capacity of liquid water to also be constant under constant pressure P_0 over the temperature range studied here. Then it follows

$$\Delta G(T) = \Delta H(T) - T\Delta S(T) \quad (\text{A.26})$$

$$\Delta H(T \pm \Delta T) = \Delta H(T) \pm \Delta C_P \Delta T \quad (\text{A.27})$$

$$\Delta S(T \pm \Delta T) = \Delta S(T) + \Delta C_P \ln \frac{T \pm \Delta T}{T} \quad (\text{A.28})$$

$$\Delta S(T) \approx -\frac{\Delta G(T + \Delta T) - \Delta G(T - \Delta T)}{2\Delta T} \quad (\text{A.29})$$

which are the typical equations of the finite difference method of computing the thermodynamic entropy. In these equations, all the Δ quantities correspond to the difference of the thermodynamic quantities between the liquid phase at P_0 and the ideal gas phase at P^* .

In similar simulations run at pressure P_0 but temperatures $T \pm \Delta T$ we analogously find

$$\Delta G_{sim}(T - \Delta T) = \mu_l(N, P_0, T - \Delta T) - \mu_g(N, P_1, T - \Delta T) \quad (\text{A.30})$$

$$\Delta G_{sim}(T + \Delta T) = \mu_l(N, P_0, T + \Delta T) - \mu_g(N, P_2, T + \Delta T) \quad (\text{A.31})$$

where P_1 and P_2 correspond to the ideal gas pressure with the same temperature and number density as the simulated liquids. Note that the ΔG values obtained from simulation differ from those occurring in equation A.29 because the reference gas phase free energies differ, and thus we must explicitly correct for this difference in reference state. By adding a correction term $\Delta G_{corr}(T \pm \Delta T)$ to the simulated free energy, we were able to use equation A.29 to calculate the entropy at temperature T , where:

$$\begin{aligned} \Delta G_{corr}(T - \Delta T) &= \mu_g(N, P_1, T - \Delta T) - \mu_g(N, P^*, T - \Delta T) \\ &= k(T - \Delta T) \ln \frac{P_1}{P^*} \end{aligned} \quad (\text{A.32})$$

$$\begin{aligned} \Delta G_{corr}(T + \Delta T) &= \mu_g(N, P_2, T + \Delta T) - \mu_g(N, P^*, T + \Delta T) \\ &= k(T + \Delta T) \ln \frac{P_2}{P^*} \end{aligned} \quad (\text{A.33})$$

and

$$\Delta S(T) = - \frac{\Delta G_{sim}(T + \Delta T) + \Delta G_{corr}(T + \Delta T) - \Delta G_{sim}(T - \Delta T) - \Delta G_{corr}(T - \Delta T)}{2\Delta T} \quad (\text{A.34})$$

These corrections, although small in magnitude, were systematically of opposite sign at temperatures $T \pm \Delta T$ because the thermal expansion coefficient of liquid water differs from the thermal expansion coefficient of the ideal gas. As a result, failure to apply these corrections will lead to a non-negligible systematical bias in the FD-FEP entropy.

The thermodynamic cycle indicating the whole process, including correction terms, is depicted in A.1. Note that in the cycle depicted in A.1, we must compute the correction terms at temperatures $T \pm \Delta T$ in order to compute the slope of ΔG with respect to T , ie

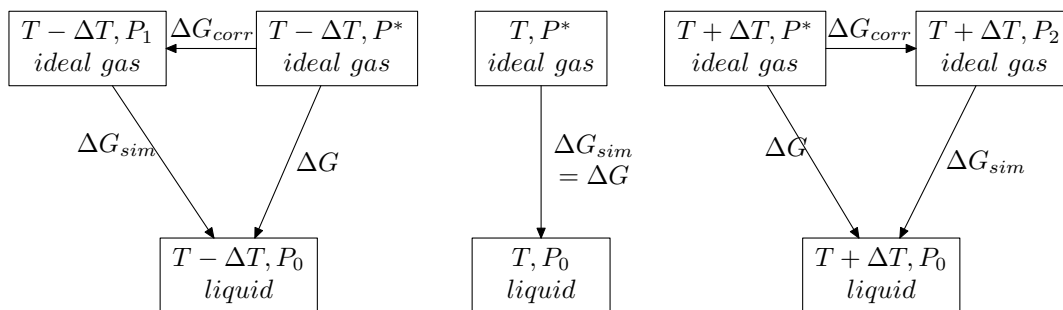


Figure A.1: Thermodynamic cycle depicting the constant pressure corrections to ΔG_{sim} at temperatures $T \pm \Delta T$ when computing the slope of ΔG_{sim} with respect to T .

the entropy associated with the solvation free energy of transferring the water molecule from the gas phase to the liquid phase at temperature T .

Appendix B

Comparison between REST and TREM

B.1 Can TREM be more efficient than REST1?

In a previous paper we noted that REST1 can sometimes be less efficient than TREM in systems in which there are large conformational energy changes between folded and unfolded structures. We attributed this to the absence of E_{ww} in the REST1 acceptance ratio formula, (Eq. (6.2)), a term that might be able to compensate for the large differences of $(E_{pp} + (1/2)E_{pw})$ between the folded and unfolded conformations.[95] On further analysis it appears that this may not be the reason for this behavior of REST1 in these systems. On the one hand, the acceptance ratio for an exchange from a folded structure to an unfolded structure is much lower in REST1 than in normal TREM if $\Delta(E_{pp} + (1/2)E_{pw})$ is much larger than $\Delta(E_{pp} + E_{pw} + E_{ww})$ between the folded and unfolded structures. (The acceptance ratio in normal TREM depends on $\Delta(E_{pp} + E_{pw} + E_{ww})$) On the other hand, the unfolded structure is sampled with much larger probability in REST1 than in normal TREM at higher temperatures. This is expected because the potential energy for water is scaled in REST1 making the unfolded structure more favorable. If high temperature replicas sample the whole conformation space efficiently, both unfolded and folded structures should be observed. Detailed balance then shows that the unfolding and folding rates at T_0 for TREM and REST1 will be identical and the correct distribution at T_0 should be maintained for

TREM and REST1. So the absence of E_{ww} is not responsible for the inefficient sampling of REST versus TREM as was thought previously. The problem is that in REST1 the replicas at higher temperatures sample the unfolded structure with dominating probability and the overlap of conformation space for replicas at lower and higher temperatures is very small. This is not a problem in REST2 where there is a smaller scaling factor of the E_{pw} term.

B.2 Why REST2 is more efficient than TREM.

There are two important factors that make REST2 more efficient than TREM. Firstly, TREM uses far more replicas than REST2. Secondly, the higher temperature trajectories in TREM are much slower with respect to cpu time than in REST2. This occurs because in the higher levels the atoms in TREM move much faster than in the higher levels of REST2, thus requiring that, for the same skin thickness, its nearest neighbor list must be updated much more frequently (or alternatively, for fixed update frequency, the skin thickness must be increased). In either case the TREM trajectory requires longer cpu times. Because replica exchange is attempted at a constant time interval, the speed in TREM is limited by the longer cpu times for the high temperature replicas. For example, for benchmark MD trajectories of the same duration run for the trpcage at 600K and at 300K, using DESMOND, the 300K trajectory requires half the cpu time that the 600K trajectory requires. Thus we save a factor of 48/10 from the smaller number of replicas and a factor of 2 for each trajectory (from the above difference in cpu times for trajectories of the same length) for a total speed up of at least $4.8 \times 2 = 9.6$ of REST2 over TREM for the trpcage. Although TREM and REST2 are rigorous sampling methods, TREM will converge much more slowly in cpu time than REST2.

Appendix C

Comparison between OPLS 2005 and OPLS 2.0 force fields for ligands CDA and CDB

C.1 Charge distributions on the Pyridine ring

The charge distributions on atoms of the P3 pyridine ring for ligand CDA from the two force fields with the atom numbers labeled as in Fig. C.1 are given in Table. C.1. It is clear that in the OPLS 2005 force field, the magnitude of the charges on the atoms of the P3 pyridine ring are very large, making it very polar and favoring its pointing into a polar solvent like water, while the OPLS 2.0 force field correctly assigns the charges on these atoms and the correct binding pose was sampled. Similar differences of charge distributions on the P1 pyridine ring and on ligand CDB were found for these two force fields.

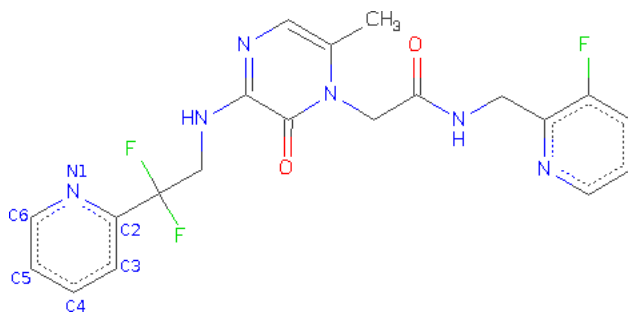


Figure C.1: The atom numbering on the P3 pyridine ring for ligand CDA

Table C.1: Charge distributions on atoms of the P3 pyridine ring for ligand CDA for two force fields.

Atom number	OPLS 2.0	OPLS 2005
N1	-0.431	-0.678
C2	0.045	0.370
C3	-0.138	-0.447
C4	-0.091	0.227
C5	-0.164	-0.447
C6	0.089	0.473

C.2 The distribution of dihedral angle involved in the flipping of the P1 pyridine ring

The distribution of the dihedral angle involved in the flipping of the P1 pyridine ring (N-C-C-C labeled in Fig 1 in text) determined from a REST simulation of the Thrombin/CDA complex using the OPLS 2005 force field for the CDA where the “hot region” is taken to include the ligand and the 10 residues surrounding the binding pocket is given in Fig. C.2. Two conformations corresponding to the crystal structure are found and an erroneous additional state with even larger probability was observed in the simulation. This erroneous state might be due to deficiencies in the force field for the ligand. This is validated from the distribution of the dihedral angle for ligand CDA in gas phase simulations using OPLS 2005 force field (See Fig. C.3), where the intrinsic potential energy of the ligand favors the erroneous state. This erroneous state does not appear in the simulation using the OPLS 2.0 force field for the ligand, and then the correct binding pose was sampled.

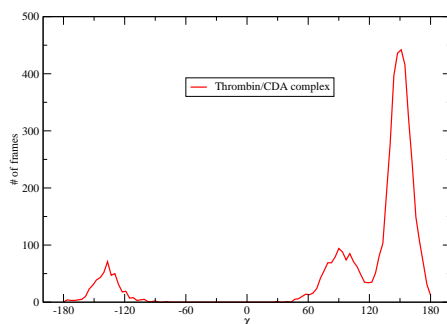


Figure C.2: The distribution of dihedral angle involved in the flipping of P1 pyridine ring (N-C-C-C labeled in Fig1 in text) for Thrombin/CDA complex using OPLS 2005 force field for the ligand.

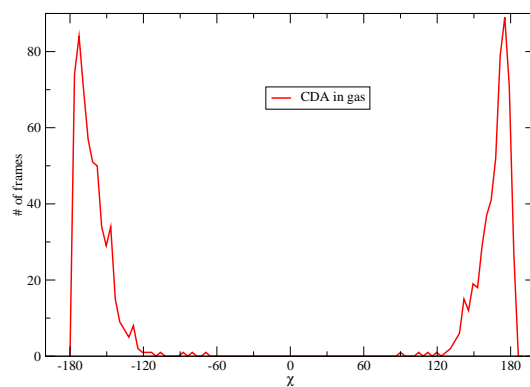


Figure C.3: The distribution of dihedral angle involved in the flipping of P1 pyridine ring (N-C-C-C labeled in Fig1 in text) for ligand CDA in gas phase using OPLS 2005 force field.

Appendix D

How to Treat bonded interactions in FEP simulation

The bonded interactions connecting the dummy atoms were treated differently from the default method in Desmond.[190] In this section, we give the details about how the bonded interactions involving the dummy atoms are treated in FEP/REST to avoid singularities and instabilities. As mentioned in the paper, this is a problem often not appreciated in the literature on FEP.

In the dual topology FEP method, depending on whether the bonded interactions between the dummy atoms and the rest of the mixed molecule are scaled or not, there are two different methods: “the ideal gas atom end state” method (scaled) and “the ideal gas molecule end state” method (non-scaled). In the ideal gas atom end state method, the dummy atom does not have any bonded interactions with the rest of the molecule and would move freely in the whole simulation volume, making the sampling very difficult. Thus most programs, including Desmond, use the ideal gas molecule end state method, in which the dummy atoms are bonded with the rest of the molecule. However, if there are more than one bonded stretch or bonded angle or bonded dihedral angle interactions between a dummy atom and the rest of the molecule, the distributions sampled for the mixed molecule (molecule with the dummy atoms) will be different from the molecule without the dummy atoms.[190] So in Desmond, only one bonded stretch, bonded angle and bonded dihedral

angle interactions involving a dummy atom are kept while all the other bonded interactions are scaled to zero at the end state.[190] In this way, the contributions of the dummy atoms to the free energies in the binding complex and in pure solvent will be identical and consequently the relative binding affinity will be independent of the dummy atoms. This follows because the relative binding affinity is equal to the difference of these two free energies. For most systems, this method works well, but for some systems, like the sets of ligands studied in this article, it will cause serious problems in the FEP simulation, which is explained in what follows.

Fig. D.1a displays how the structure of the mixed molecule is mutated from a benzene molecule to a p-xylene molecule. The two dummy hydrogen atoms which are mutated to the methyl groups all have two bonded angle and bonded dihedral angle interactions with the rest of the molecule. Take the dummy hydrogen atom numbered 6 in Fig. D.1a as an example. It has two bonded angle and bonded dihedral angle interactions with the rest of the molecule ($\theta_1, \theta_2, \phi_1, \phi_2$ labeled in Fig. D.1a). If all these bonded interactions are kept in the end state, the distribution of the angle formed by atoms numbered 2, 3, and 4 for the mixed molecule with the dummy atom will be different from the distribution of the real molecule without the dummy atom. So in the default setup of Desmond, only one bonded angle and one bonded dihedral angle interaction involving the dummy hydrogen atom (θ_1, ϕ_1) is kept in the end state, and the others (θ_2, ϕ_2) are scaled to zero. The energy difference ($U_1 - U_0$) (where U_1 is the potential energy in the previous lambda window and U_0 is the potential energy in the end state lambda window) sampled in the end state lambda window using Desmond's default setting is given in Fig. D.2. Clearly, the energy difference fluctuate about three different values (approximately -0.5, 21, and 40 kal/mol respectively).

From the simulated trajectories, we observed that the two dummy hydrogen atoms were located at the correct positions in the initial stage (Fig. D.1a), then one of them moved to the position which almost overlapped with a carbon atom on the ring (Fig. D.1b), and at the end the other hydrogen atom moved to the position which almost overlapped with another carbon atom on the ring (Fig. D.1c). These three configurations for the mixed molecule are located in the potential energy minima for the end state where only one bonded angle and one bonded dihedral angle interaction are kept (θ_1, ϕ_1) for the dummy hydrogen atoms (the

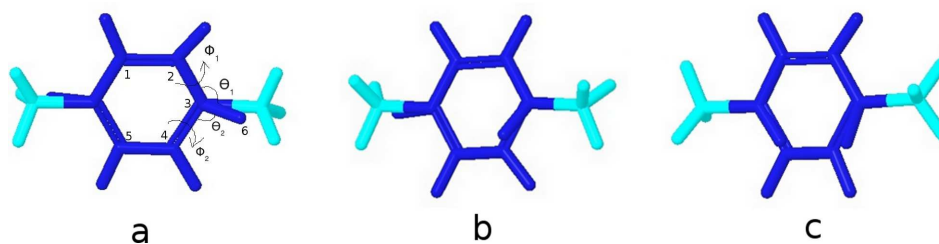


Figure D.1: The structure of the mixed molecule to mutate from benzene to p-xylene. Three different configurations are sampled in the end state when only one bonded angle and bonded dihedral angle interactions are kept for the dummy hydrogen atoms.

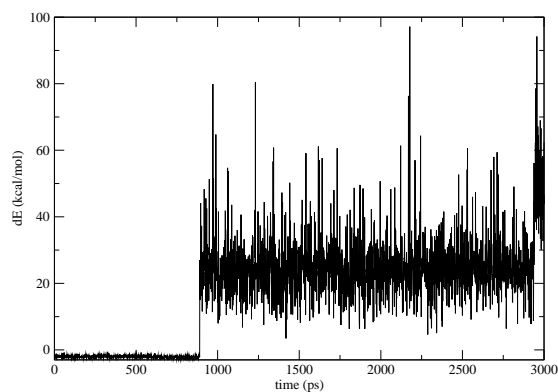


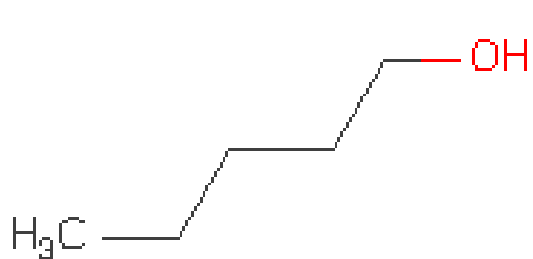
Figure D.2: The energy difference between the previous lambda window and the end state lambda window sampled by the end state lambda window for mutation from benzene to p-xylene when only one bonded angle and bonded dihedral angle interactions involving the dummy hydrogen atoms are kept in the end state.

default setup of Desmond). In the previous lambda window, however, when all the bonded interactions involving the dummy atoms (including θ_2, ϕ_2) are slowly turned on, the last two configurations (b and c) have large bonded angle interactions (also 1-4 pair interactions) leading to instabilities for the free energy calculation. This is the reason for large jumps in the energy difference profile displayed in Fig. D.2 and the three distinct values of energy difference correspond to these three configurations. (In extreme cases, when either dummy atom is located at the same position as the carbon atom on the ring, the 1-4 pair interaction will cause a singularity in the free energy calculation.)

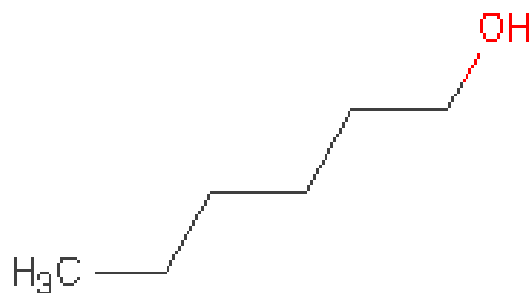
To avoid instabilities and singularities in the free energy calculations caused by the bonded interactions involving dummy atoms, we treated these bonded interactions in the end state differently in this article. We choose to keep all of the bonded stretch and bonded angle interactions involving the dummy atoms in the end state. For the dihedral angle interactions, only one bonded dihedral angle interaction involving a dummy atom is kept while the other dihedral angle interactions are scaled to zero in the end state. Take the dummy hydrogen atom numbered 6 in Fig. D.1 for example. We include two bonded angle terms (θ_1, θ_2) and one bonded dihedral angle term (ϕ_1) in the end state, while the other bonded dihedral angle term (ϕ_2) is scaled to 0 in the end state. Thus in the end state configurations b and c in Fig. D.1 are located in the high energy region in phase space, and thus will not be sampled. In this way, the instability and singularity problems encountered for the bonded interactions in FEP are eliminated. Although the distributions sampled for the “mixed molecule” (molecule including the dummy atoms) might be a bit different from the distributions for the molecule without the dummy atoms, the error introduced in this treatment is negligible because the fluctuations of bond angle are very small for the two sets of ligands studied in this article. In addition, the error in the relative binding affinity, which is the difference between the free energies in the binding complex and in free solvent, will be very small, because there is an excellent cancellation of errors in these two free energies.

Appendix E

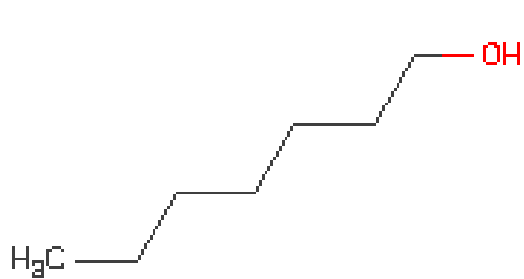
Structures of ligands studied in Chapter 4



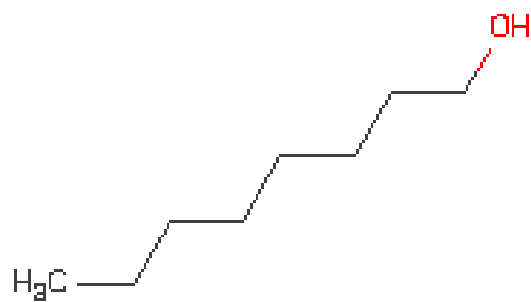
PE9:1ZND



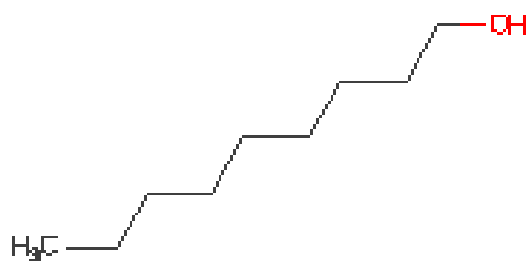
HE2:1ZNE



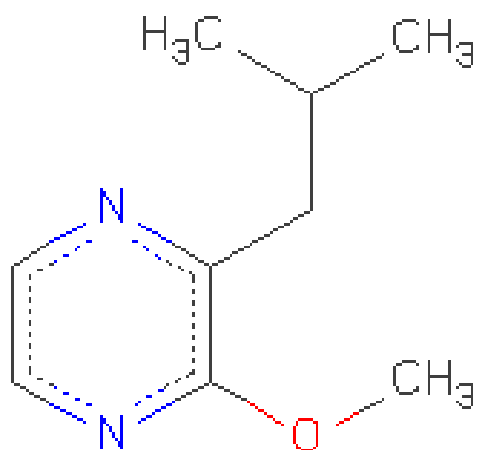
HE4:1ZNG



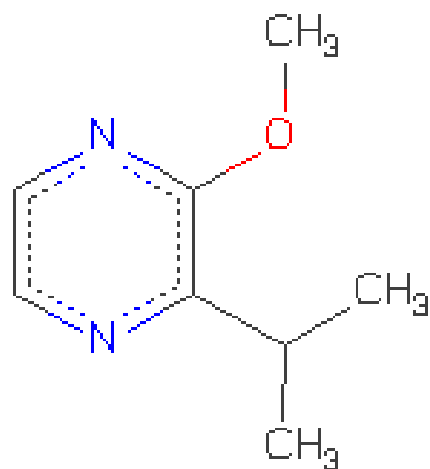
OC9:1ZNH



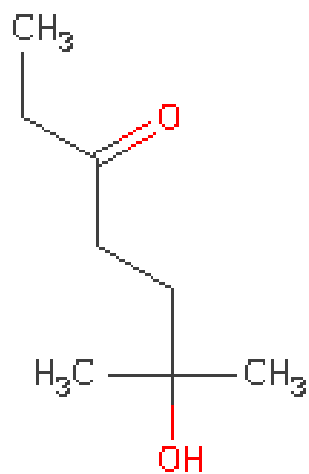
F09:1ZNG



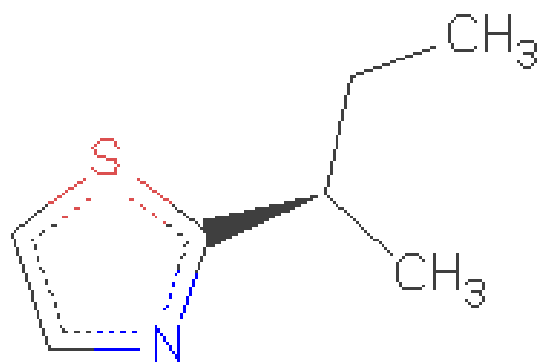
IBMP:1QY1(PMZ)



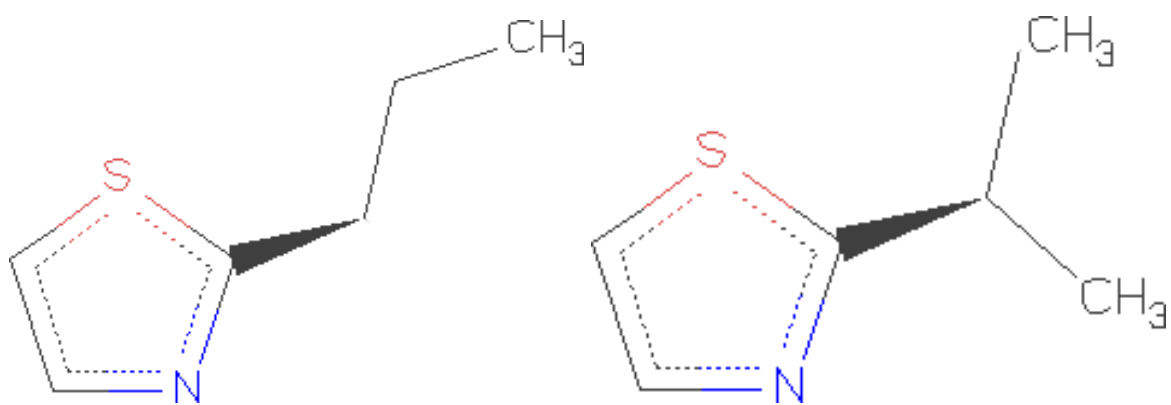
IPMP:1QY2(IPZ)



LTL:1I05(HMN)

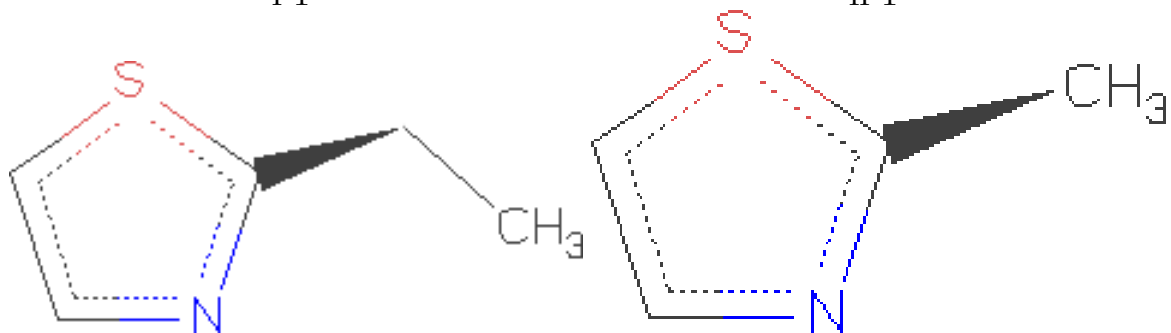


TZL:1I06(SBT)



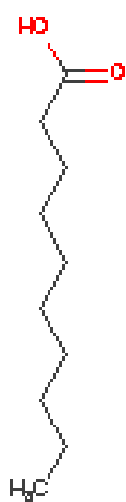
PT

IPT

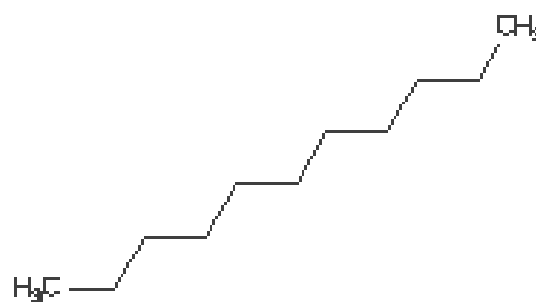


ET

MT



DKA:1WBE



UND:1Y9L