

SCALABLE MOLECULAR DESIGN USING REVERSIBLE JUMP MCMC AND STOCHASTIC APPROXIMATION

A Dissertation

Presented to the Faculty of the Weill Cornell Graduate School
of Medical Sciences
in Partial Fulfillment of the Requirements for the Degree of
Doctor of Philosophy

by

Patrick B. Grinaway

May 2019

© 2019 Patrick B. Grinaway

ALL RIGHTS RESERVED

SCALABLE MOLECULAR DESIGN USING REVERSIBLE JUMP MCMC AND STOCHASTIC APPROXIMATION

Patrick B. Grinaway, Ph.D.

Cornell University 2019

Despite the existence of useful models for atomic-scale interactions, designing novel molecules (such as drugs) using this information has been extremely difficult. The difficulty results from the nature of the model, which is both extremely high dimensional and multimodal, as well as the nature of the objective function, which is an expectation under this model. In prior work, these models would be leveraged by computing individual expectation values and using these to rank the various molecules. Here, we introduce a joint probability distribution of configurations and chemical states, allowing the simulation to visit different molecular identities as well as different configurations. This introduces the requirement for reversible jump MCMC, as a change in molecular identity results in a change in dimensionality of the configurations. We then combine this approach with the Self-Adjusted Mixture Sampling (SAMS) technique developed by Tan to achieve sampling of chemical identities according to an arbitrary prespecified distribution. We then sought to use the relative free energies of each state as the target distribution, effectively prioritizing favorable chemical states. However, this requires the estimation of free energies. To resolve this, we simultaneously run another MCMC chain that provides online estimates of the necessary free energies. We additionally generate a transdimensional version of nonequilibrium switching amenable to highly parallel hardware setups.

BIOGRAPHICAL SKETCH

Patrick Grinaway graduated from Lafayette College in 2011 with a Bachelor of Arts in Engineering Studies, focusing on coursework in Electrical and Computer Engineering. While he was attending Lafayette, he became fascinated with modeling complex systems and understanding biology. Under the mentorship of Professor Laurie Caslake, he studied the radiation tolerance of bacteria, and under the mentorship of Professor Michael Kelly, he studied the . After graduating, he briefly worked with Professor William Terzaghi of Wilkes University, studying the role of noncoding RNA in fine-tuning the expression of genes regulated by the gene phytochrome A in *O. sativa*. He joined the Physiology, Biophysics, and Systems Biology program in 2012 and became a member of the Chodera lab in 2013, where his work has focused on the development of Markov chain Monte Carlo algorithms for scalable free energy calculations.

This thesis is dedicated to Florencia Fama, without whose constant support this work would have been impossible.

ACKNOWLEDGEMENTS

I would like to acknowledge my advisor, John D. Chodera, for introducing me to the world of Markov chain Monte Carlo, as well as other Chodera lab members, in particular, Julie Behr, Ivy Zhang, and Hannah Bruce MacDonald, who contributed greatly to this project and Bas Rustenburg for excellent scientific discussions. I would also like to thank Silicon Therapeutics for providing funding for part of this work and in particular, Woody Sherman and Bryce Allen, who provided substantial intellectual contribution. Finally, it is critical that I acknowledge my thesis committee, which included Professors Alex Kentsis, Thomas Fuchs, and Zhiqiang Tan, as well as those attending my defense, including Professors Daniel Heller and Olivier Elemento.

TABLE OF CONTENTS

Biographical Sketch	iii
Dedication	iv
Acknowledgements	v
Table of Contents	vi
List of Tables	ix
List of Figures	x
1 Introduction	1
1.1 The role of molecular design in drug discovery	1
1.1.1 Physical modeling of drug-target interaction	2
1.1.2 Free energy calculations	2
1.1.3 Typical approaches to computation of binding free energy	4
1.1.4 Expanded Ensembles	4
1.1.5 The search space for designs is very large	6
1.1.6 Monte Carlo methods can aid such a large search space	7
1.2 Chemical Space Sampling	7
1.2.1 Chemical Monte Carlo	7
1.2.2 Lambda dynamics	8
1.2.3 Nonequilibrium Switching	9
1.3 Summary	10
2 Reversible Jump MCMC for Molecular Simulation	11
2.1 A Rigorous and Efficient Formulation of Chemical Monte Carlo	11
2.1.1 Key difficulties	11
2.1.2 Transdimensional Nonequilibrium Switching	14
3 Exploration of Chemical Space	15
3.1 Introduction	15
3.2 Quantitative Metrics	16
3.3 Practical Approaches	18
3.3.1 Chemical State Proposal Algorithm	18
3.3.2 MCSS Hyperparameters	20
3.4 Factors that affect acceptance probability	20
3.4.1 An exploration of the effect of various atom mapping options on acceptance probability	24
3.5 Future work	26
3.5.1 Atom map learning	26
3.5.2 Relieving the restriction to a finite set of chemicals	28

4	Geometry Proposals	29
4.1	Introduction to dimension matching	29
4.2	Problems faced by dimension matching in chemical space	30
4.2.1	The target is highly multimodal	30
4.2.2	The proposal must be drawn exactly	30
4.2.3	The proposal must be associated with a normalized probability	31
4.2.4	Atomic positions are correlated	32
4.2.5	Propose one atom at a time	32
4.2.6	Naive bond/angle/torsion doesnt work	33
4.3	CBMC-like Algorithm	34
4.3.1	Description of algorithm	34
4.3.2	Performance on simple cases	37
4.3.3	Shortcomings	40
4.4	Future directions: Particle Filtering	41
4.4.1	Basic Overview	41
4.4.2	Advantages	41
4.4.3	Algorithm Hyperparameters	42
4.5	Tuning of Dimension Matching Parameters in the context of NCMC	43
4.5.1	Tradeoffs	44
5	NCMC Switching	46
5.1	Introduction	46
5.2	Description of Algorithm	46
5.2.1	Construction of Hybrid System	47
5.3	A note on ring closure	53
5.4	Tuning NCMC Protocol length	57
5.4.1	Limitations of tuning only length	57
5.4.2	Tuning Annealing Schedule	58
5.4.3	Limitations to tuning annealing schedule	59
5.4.4	Relationship to Geometry Tuning	60
6	Stochastic Approximation	62
6.1	Stochastic Approximation for Molecular Simulation	62
6.1.1	SAMS	63
6.2	Doubly-recursive SAMS	64
6.2.1	Description of Algorithm	65
6.3	Toy Examples	65
6.3.1	Performance	67
6.4	Limitations	67
6.4.1	Convergence Rate	67
6.4.2	Number of States	69
6.5	Weight initialization	69
6.5.1	Implicit volume term	70

6.5.2	Hydration	70
6.5.3	Complex weight initialization	71
7	Alternative to SA: Highly Parallel Nonequilibrium Switching	74
7.1	Introduction	74
7.2	Advantage of RJ Nonequilibrium Switching	75
7.3	Use of cloud computing resources	76
7.4	Pricing and Economics of Cloud Computing	77
7.4.1	Factors affecting performance	79
8	Hydration free energy	80
8.1	Introduction	80
8.2	Thermodynamic Cycle	80
8.3	Vacuum Phase	82
8.3.1	Explicit Phase	82
8.3.2	Nonequilibrium switching protocol length	82
8.4	Analysis of the work of the entire move	83
9	Conclusion	85
9.1	Contributions of work	85
9.2	Exploration of Chemical Space	85
9.2.1	Chemical State Proposals	86
9.2.2	Dimension Matching	86
9.2.3	NCMC Switching	88
9.2.4	Stochastic Approximation	90
9.3	Transdimensional Nonequilibrium Switching	91
9.4	Conclusion	92
A	Appendix A	93
A.1	Derivation of the Acceptance Probability	93
A.1.1	Proposal of Chemical State Jump	93
A.1.2	Derivation of Acceptance Probability	94
	Bibliography	95

LIST OF TABLES

4.1	Comparison of features of various proposal schemes for dimension matching.	44
-----	--	----

LIST OF FIGURES

3.1	An example of an apparently reasonable atom map between benzene and cyclohexane.	20
3.2	3D structures of benzene and cyclohexane. Note that benzene is flat, while cyclohexane adopts a "chair" conformation.	21
3.3	A geometry that was constructed from an atom map that left no space to properly close the rings.	22
3.4	A geometry resulting from the inclusion of bond order as a criterion for mapping atoms. More of the system is rebuilt, so a better geometry can be created.	23
3.5	This map between kinase inhibitors may contain more atoms than another, but as it is partially breaking a ring, energies after the proposal are often highly unfavorable	24
3.6	The set of trial molecules used for empirical exploration of atom map criteria. This set was used to demonstrate because it contains very similar cores, but different substituents that might affect the performance of atom maps.	24
3.7	Forward (blue) and negative reverse (red) log acceptance probability distributions for benzene to chlorobenzene in vacuum under different mapping schemes. Here, it is clearly visible that the choice of mapping scheme has a profound effect on the ultimate quality of the proposal.	25
3.8	Forward (blue) and negative reverse (red) log acceptance probability distributions for toluene to phenol in vacuum under different mapping schemes. Here, the choice of the mapping scheme does not seem to have as significant of an impact on the quality of the proposal.	26
4.1	An example of a unimodal proposal (blue) with a multimodal target(red). Note that there are many regions of the target that are poorly covered by the proposal.	31
4.2	An example of an internal coordinate description of atomic positions. The atom highlighted in green is being proposed, based on the bond (r), angle (θ), and torsion (ϕ).	33
4.3	Examples of substituted benzenes	38
4.4	Examples of n-alkanes	38
4.5	The average $\ln P_{\text{accept}}$ vs. the number of degrees of freedom added. Note that the variance is quite low for small (chlorobenzene and dichlorobenzene) changes, but quickly grows with the addition of flexible chains. Note that anything beyond the 98th percentile was clipped for this figure.	40

5.1	Chlorobenzene, with atoms colored by charge. If the ring is mapped to benzene (where all partial charges are equal), the abrupt change is very unfavorable.	52
5.2	Nonbonded interpolation scheme, with the schedule of each individual force's λ parameter plotted against the global λ . Note that sterics are always eliminated after electrostatics, and are always added before electrostatics. This provides protection against unshielded charges.	52
5.3	The distributions of the logP_to_hybrid contribution under different initial bond softening parameters. The initial (darkest and broadest) distribution is with no softening; moving to lighter colors, we see the distribution narrow considerably as the initial force constants for bonds connecting unique atoms is scaled by 0.1, 0.01, 0.001, and 0.0001	54
5.4	The distributions of the logP_to_hybrid contribution under different initial angle softening parameters. The initial (darkest and broadest) distribution is with no softening; moving to lighter colors, we do not see the same profound effect as the initial force constants for angles connecting unique atoms is scaled by 0.1, 0.01, 0.001, and 0.0001	55
5.5	The distributions of the logP_to_hybrid contribution under different initial angle and bond softening parameters. The initial (darkest and broadest) distribution is with no softening; moving to lighter colors, we see that the distribution grows considerably narrower as both the angle and bond terms are softened to the same degree according to the schedule 0.1, 0.01, 0.001, 0.0001. . .	55
5.6	The standard deviation of the logP_to_hybrid component under different combinations of bond and angle softening constants. It is noteworthy that the angles need not be softened nearly as much as the bonds. All quantities are in effective units of $k_B T$. .	56
5.7	The standard deviations of the standard deviations of logP_to_hybrid under different combinations of bond and angle softening terms. All quantities are in effective units of $k_B T$	56
6.1	The target weights of the various harmonic oscillators (solid) and the true value (dotted line). Note that the higher free energy states take significantly longer to converge to their true values	66
6.2	Convergence of binary SAMS from a single realization of harmonic oscillators with only a small separation in means.	68

7.1	Illustration of a nonequilibrium switching free energy calculation. Note the requirement for simulation at the alchemical endpoints. The horizontal green lines represent the equilibrium simulations at each endpoint, while the magenta lines depict the switching trajectories from samples of each state to the other.	76
8.1	The thermodynamic cycle of the hydration free energy calculation. Blue arrows denote the data available in FreeSolv; magenta arrows denote the legs performed in this work.	81
8.2	The standard deviation of the work for an NCMC switching trajectory between benzene and naphthalene vs. number of steps, with a 1fs timestep. Top: work for transforming benzene into naphthalene. Bottom: work for transforming naphthalene into benzene.	83
8.3	The forward and reverse (blue and red, respectively) log acceptance probability distributions for transitions between hexachlorobenzene and ethylbenzene. Although the variances of the distributions appear high, the overlap is very good, allowing a low-error free energy estimate.	84
8.4	The forward and reverse (blue and red, respectively) log acceptance probability distributions for transformations between naphthalene-2-amine and hexachlorobenzene. Despite having to create or break a fused ring, the overlap is still quite good, allowing low variance free energy estimates	84

CHAPTER 1

INTRODUCTION

1.1 The role of molecular design in drug discovery

Drug discovery includes many facets, from identifying the relevant biology to performing studies ensuring that a drug of interest is effective for the target population [66]. Typically, the process starts with the identification of the so-called target molecule based on biological studies [66]. Once this target molecule (usually a protein [73]) is found, the task of molecular design begins. Often, the target of interest is an enzyme [73] or a receptor where a natural small molecule binds. The task of the designer, then, is to produce a molecule that can outcompete the natural ligand, partially or fully ablating the activity of the target. Of course, molecular design in drug discovery is not limited to this case. There are many other cases where a small molecule is sought that binds to an alternative, or allosteric binding site [41, 63], or even one that causes the association of the target with another protein, as in PROTAC [13, 85]. In many of these cases, rational design begins with the determination of the structure of the target, either by modeling or by experiment [48], after which the task of actually designing the molecule can begin. In the last several decades, computers have become instrumental in aiding this rational design phase [46, 89], allowing chemists to filter out unwanted compounds before testing or synthesis, and providing directions for novel ideas. There are many computer-based approaches, ranging from machine learning [45, 10] to physical modeling [52, 82, 6, 75]. In this work, I will focus on physical modeling approaches, though even in that realm there are many opportunities for machine learning approaches to make valuable inroads.

1.1.1 Physical modeling of drug-target interaction

An attractive avenue for computer-assisted drug discovery is to model, at an atomic scale, how a hypothetical compound and its target might interact.

Docking

Given that we often begin with a snapshot of the target’s structure, one approach, known as molecular docking [78, 11], holds the target (and often the small molecule) rigid, and uses a scoring function to assign an optimal position for the atoms of the small molecule. Scoring many such molecules like this allows us to rank hypothetical molecules, a process known as virtual high-throughput screening [77]. This practice has found wide adoption in drug discovery [77] due to its relative computational simplicity and its intuitive nature (one can examine scored poses manually). However, docking does suffer from an important shortcoming: the fact that the receptor is held rigid means that the flexibility of the target is not taken into account [91]. Fortunately, statistical mechanics offers us a path forward to overcome this limitation: free energy calculations [28].

1.1.2 Free energy calculations

Given a model of atomic-scale interactions known as a forcefield, we can model the configurational ensembles of ligand-target association and rigorously (within the confines of the forcefield) compute the strength of the association. For the simple bimolecular association reaction,



the so-called absolute free energy of binding, ΔG , can be computed in terms of a ratio of partition coefficients [28],

$$\Delta G = -k_B T \ln \left[\frac{Z_{PL}}{Z_L Z_P} \right] + c \quad (1.2)$$

where c is an additive constant and the individual partition coefficients Z_* are given by,

$$Z_* = \int dx e^{-u_*(x)} \quad (1.3)$$

Here, $x \in \mathbb{R}^{3N}$ is the positions of the atoms, k_B is the Boltzmann constant, T is the absolute temperature, and $u_*(x) = \beta U_*(x)$ is the reduced potential [74] for species $*$, with $\beta = (k_B T)^{-1}$ denoting the inverse thermal energy and $U_*(x)$ the corresponding potential energy. Although this form is apparently very convenient, the integral in is extremely high dimensional, as each atom contributes three degrees of freedom. For realistic systems, however, the dimensionality can easily exceed 100,000 degrees of freedom, making deterministic approximations completely infeasible. However, it is nonetheless possible to estimate the ratio above using stochastic methods such as Markov chain Monte Carlo (MCMC) [43]. These methods have seen a massive upsurge in interest in recent years as commercial applications have burgeoned and the availability of cheap computing power has grown [89]. Despite this, it remains very costly to perform a free energy calculation. On modern hardware, it can easily take 8 to 12 hours to perform a rough free energy calculation [89]. Coupled with the short timetables on which industrial medicinal chemists operate (based on personal communication with industrial scientists), extending free energy calculations to practical use is a formidable challenge.

1.1.3 Typical approaches to computation of binding free energy

There are several Monte Carlo-based approaches commonly used to estimate the binding free energy of a small molecule to a target protein. Since we are interested purely in the ratio [27, 28, 74], the estimators used here will directly estimate the ratio. [74, 50, 27] The quality of this ratio estimate will depend strongly on the overlap between the probability distributions at each endpoint [50, 27]. Because of this strong dependence on the overlap between the endpoints (which is often quite low and a cause of high Monte Carlo errors) [27], it is prudent to construct intermediate nonphysical or "alchemical" states that contain a mixture of the parameters of either endpoint distribution. In this way, neighboring distributions can be constructed with high overlap, enabling efficient estimation. For the purpose of this work, we will denote the construction of neighboring distribution using a (possibly multidimensional) control parameter λ , which indexes over intermediate distributions as $p(x; \lambda) \propto \exp(-u(x; \lambda))$. The control parameter is defined so that $\lambda = 0$ refers to the distribution in the denominator of the partition function ratio, and $\lambda = 1$ refers to the distribution in the numerator. There are many methods by which one could exploit this approach; in this work, we will focus on the method of expanded ensembles [47], as well as the method of nonequilibrium switching [36, 1].

1.1.4 Expanded Ensembles

One approach that is becoming increasingly common is the method of expanded ensembles [47]. In this method, one constructs a joint distribution $p(x, k)$, where $x \in \mathbb{R}^{3N}$ again denotes atomic positions, and the discrete index $k \in \{1, 2, \dots, K\}$

indexes over some set of K distributions between which one would want relative free energies:

$$p(x, k) \propto e^{-u_k(x) + g_k} \quad (1.4)$$

Note that we have also introduced K free parameters, or log weights, g_k , $k \in \{1, 2, \dots, K\}$, which can be used to bias the sampling of the mixture components. Having constructed this joint distribution, one could sample a sequence of iterates (x_t, s_t) , $t \in \{1, \dots, T\}$, from the expanded ensemble defined by $p(x, k)$ by the following algorithm:

$$x_{t+1} \sim p(x | s_t) \propto e^{-u_{s_t}(x)} \quad (1.5)$$

$$s_{t+1} \sim p(s | x_{t+1}) = \frac{e^{-u_s(x_{t+1}) + g_s}}{\sum_{k=1}^K e^{-u_k(x_{t+1}) + g_k}} \quad (1.6)$$

where $u_k(x)$ is the reduced potential [74, 12] of configuration x in thermodynamic state k . Following such a simulation, it should be apparent that

$$\frac{p(k = m)}{p(k = n)} = \frac{\int dx p(x, m)}{\int dx p(x, n)} = \frac{Z_m}{Z_n} \quad (1.7)$$

where Z_n is defined as in Eq. 1.3. However, since the relative free energies of the different components of the mixtures are the logarithm of the relative populations,

$$\Delta G = G_m - G_n = -\ln \frac{Z_m}{Z_n} = -\ln \frac{p(k = m)}{p(k = n)} \quad (1.8)$$

even a small free energy difference can result in very large population differences. These large population differences make it very difficult to achieve enough samples of each state such that a reliable free energy estimate can be made. Recall that the expanded ensemble defined in Eq. 1.4 contains log weights g_k that allow us to bias sampling. There are several tools available to adaptively reweight the mixture components via the bias term, compensating for potentially

large differences in populations [42, 88, 84] One of the most recent of these tools is known as Self-Adjusted Mixture Sampling (SAMS) [84]. This algorithm for adapting the biasing terms g_k is appended to the sampling algorithm in Eq. 1.5 as:

$$g_k^{(t-1/2)} = g_k^{(t-1)} - t^{-1} \frac{\delta_{s_t,k}}{\pi_k} \quad (1.9)$$

$$g_k^{(t)} = g_k^{(t-1/2)} - g_1^{(t-1/2)} \quad (1.10)$$

where $\delta_{s_t,k}$ is unity if the sampler state s_t is currently visiting state k and zero otherwise, and π_k represents the desired target probability for state k . This stochastic approximation technique provably minimizes the asymptotic variance of the weights [84], and if all target probabilities π_k are set such that $\pi_k = \frac{1}{K}$, where K is the number of thermodynamic states, the log weights will asymptotically converge to the relative free energies of the different states [84]. By compensating for the differences in free energies between different mixture components, we can now efficiently sample the expanded ensemble and post-process the data to get high quality free energy estimates. However, the issue remains that each expensive calculation only yields a single relative free energy estimate, where many are needed to sort through the almost innumerable potential designs.

1.1.5 The search space for designs is very large

How many synthetically accessible chemical species are there? Relatively compact databases such as GDB-17 [72] provide over 100 billion feasible compounds of a certain size. This highlights the other significant challenge—in addition to exploring and integrating out the complex landscape of molecular conformations, we must also explore the even more poorly characterized space of possible

chemicals. This space is not only large, but also discrete, with no immediately obvious ordering. However, all hope is not lost—the Monte Carlo methods discussed above will be brought to bear on this problem as well.

1.1.6 Monte Carlo methods can aid such a large search space

The above presentation of the method of expanded ensembles almost immediately suggests an interesting idea—rather than using the discrete parameter to index over intermediate distributions, use it instead to index over different chemical species [67]. This brings to bear the power of Monte Carlo algorithms on the molecular design problem. However, each of the previous attempts contain limitations that will be overcome in this work.

1.2 Chemical Space Sampling

At its core, the concept of chemical space sampling in the context of molecular simulation first consists of defining the expanded ensemble in Eq. 1.4 to be indexed over different chemical identities, rather than intermediate distributions. The additional parameter, g_s , enables us to apply arbitrary weighting to different states, which will in turn enable the adaptive work contained in this thesis.

1.2.1 Chemical Monte Carlo

In [67], it was suggested that one could simply define the different distributions indexed by s as the ensembles of corresponding chemicals. Then, by applying a

clever weighting, one could even bias the simulation toward favorable chemical states. However, this approach suffered from several drawbacks. First of all, when one is attempting to jump from one chemical state to another, one must account for the fact that the number of atoms (and hence the dimensionality) of the two states is likely different. Second, it may be computationally infeasible in many cases to derive the clever weighting that causes sampling to favor the states of interest (for instance, those that bind tightly to the target, or those that bind tightly to one target but not another). Finally, even if the above problems were surmounted, it can suffer difficulty getting acceptances in the more accurate explicit solvent, as newly-introduced atoms are likely to clash with other atoms.

1.2.2 Lambda dynamics

Another interesting concept is that of Lambda Dynamics [38] and its related technique, Multisite Lambda Dynamics [39, 22]. This technique replaces the discrete index k with a continuous (and possibly multidimensional) parameter λ , and performs dynamics on that parameter as well as on the configurational degrees of freedom. The λ parameter now allows the simulation to visit not just the interesting endpoints, but intermediate states as well. Since it visits the intermediate states, it can potentially overcome the issue that chemical Monte Carlo has with explicit solvent. However, several problems still remain. One issue is that a biasing potential must be used to prevent the λ parameter from spending all of its time in uninteresting intermediate states. One may grow concerned as well that because the simulation spends so much time away from the endpoints, reweighting to recover an estimate of the true free energy difference may be too difficult. Another is that when the technique is extended to multiple sites, one

must discover an efficient protocol through the now-high dimensional λ space. Simply exploring this space (or indeed even adaptively adding biasing potentials as in metadynamics) will quickly become very difficult.

1.2.3 Nonequilibrium Switching

Another method that has gained in popularity for free energy calculations is known as nonequilibrium switching [1, 36, 60]. In this approach, equilibrium simulations at the endpoints of a set of intermediate distributions are first conducted. Then, from the equilibrium snapshots, a nonequilibrium switching move is carried out: that is, one takes a step of dynamics, then alters the control parameter λ slightly, then another step, and so on. Each time the control parameter is incremented, we compute the change in the potential energy. This change is then added to a work value. At the conclusion of the switch (when the control parameters have reached the other endpoint), the work values in each direction are used to estimate the free energy difference [60, 36, 1, 21]. In this scheme, the hyperparameters of the algorithm include the length of the protocol (how large each step of the control parameters is) and the schedule of changes to the control parameters. The emergence of massively parallel hardware and cloud computing has made this approach interesting, as the nonequilibrium switching trajectories are computationally independent of each other. This means that with the availability of many processing units, one could potentially spare a considerable amount of wall clock time by running nonequilibrium switching trajectories simultaneously.

1.3 Summary

Over the past several years, free energy calculations have reached a mature stage. However, there is still a considerable amount of work to be done. Calculations using the techniques described above are still very expensive, and can only compute relative free energies over pairs or, in some cases, small sets of congeneric series. In this work, I aim to accomplish three goals. First, I aim to develop a rigorous formulation of a chemical state sampling algorithm that can be used not only to compute relative free energies, but can be used to adaptively prioritize sampling based on a free energy based objective. Second, and as a component of the first, I aim to develop a scheme that allows jumps to molecules of quite different topology and geometry, adding complex structure such as fused rings. This, too, is developed in a theoretically rigorous formalism that allows the practitioner to be assured that the algorithm is asymptotically correct. Finally, I aim to develop a formulation of the methods developed here that is efficient in the highly parallel and heterogeneous computing environments that dominate modern high-performance computing. The development of theoretically rigorous foundations for each of these additionally enables future work to proceed straightforwardly from here.

CHAPTER 2

REVERSIBLE JUMP MCMC FOR MOLECULAR SIMULATION

2.1 A Rigorous and Efficient Formulation of Chemical Monte Carlo

As an alternative to the above-mentioned schemes, I propose here an algorithm for sampling from an expanded ensemble with each mixture component representing a different chemical species. I furthermore provide a proof that this algorithm preserves the expanded ensemble as an invariant distribution. Finally, I apply an extension of the Eq. 1.9 algorithm to adaptively achieve target weights that prioritize more favorable chemical species.

2.1.1 Key difficulties

In deriving and implementing a rigorous formulation of chemical space sampling that is also efficient, there are several key difficulties. One difficulty, as discussed above, is that the space itself is discrete without a good idea of neighborhoods. Unlike a continuous-valued parameter, it is not obvious which proposals should be given a high probability. Additionally, the jumps between chemical states will involve a change in the number of atoms, necessitating additional corrections to the acceptance criterion and the development of efficient dimension-matching algorithms. Compounding the issue is the need to insert atoms into a condensed phase system with a reasonable acceptance probability. For this, I resort to Nonequilibrium Candidate Monte Carlo [62], which carries with it a collection

of its own fascinating algorithmic challenges. Finally, I extend existing stochastic approximation algorithms to adaptively reweight the different chemical states in order to achieve target sampling probabilities that match the free energy objective of interest.

How are neighbors defined in chemical space?

The first of these key difficulties is how to define neighbors in chemical space. Ordinarily, the design of an efficient MCMC algorithm relies on a suitable proposal distribution[43]; that is, a probability distribution from which one can easily sample and which is capable of providing both a suitably high acceptance probability and a suitable degree of exploration of the state space. However, in the case of chemical state sampling, it is not clear, given that a chain is currently visiting chemical state k , which k' would result in a suitable acceptance probability. Furthermore, as discussed later, the "neighborhood" of chemicals will change depending on the target, making optimization of these parameters difficult to perform in the general case. In this work, I resort to a heuristic proposal distribution that is simple to compute and reasonably effective.

Transdimensional space

Even if one can hop through chemical space with enough speed, how does one actually perform the proposed jumps? For this, I resort to Reversible Jump Markov chain Monte Carlo (RJ-MCMC) [30], a formulation of Metropolis-Hastings which includes correction factors for when the current and proposed states are not defined on spaces of the same dimensionality. Although the theory presented

in [30] is extremely general, performing this task efficiently (while still meeting the necessary requirements for preserving the target distribution) is quite difficult. Additionally, and related to the chemical neighborhood problem, the dimension-jump challenge also includes the question of which degrees of freedom should be considered to be in common between the current and proposed system? On one hand, including more degrees of freedom as common results in a more straightforward task for the dimension-matching algorithm. On the other, if a set of atoms is said to be in common, but occupies a very different configurational ensemble, the acceptance probability may be very poor. In this work, I discuss several approaches to solve these issues in the context of molecular simulation.

How do we favor "good" molecules in a large set?

Lastly, even if we can construct the Markov kernel that allows us to sample the distribution of interest, how do we ensure we spend our simulation time wisely? Large free energy differences between different chemical states will all but ensure that we have great difficulty exploring chemical space. In the past, practitioners have used stochastic approximation successfully to adapt the various weights g_s of the different chemical states, typically to achieve even sampling. [88] More recently, an asymptotically optimal stochastic approximation algorithm known as SAMS [84] has become available. However, when sampling a large chemical space, one would prefer to gravitate toward favorable states, not even sampling. To that end, in this work I develop an extension of the SAMS algorithm that allows us to couple multiple MCMC chains and achieve a target sampling distribution based on a free-energy objective.

2.1.2 Transdimensional Nonequilibrium Switching

In addition to allowing a rigorous formulation of the chemical Monte Carlo expanded ensemble framework, as a byproduct, the algorithm developed here also permits a transdimensional version of nonequilibrium switching. This approach, discussed in chapter 7, allows a highly parallel application of the reversible jump algorithm presented here. Since the atoms present in one molecule but absent in another are added as part of the proposal attempt, this allows one to more efficiently use equilibrium simulation. Additionally, it allows the possibility of ring forming and breaking, which is otherwise quite difficult in free energy calculations, which has been a problem considered difficult for some time [44]. Finally, allowing for a considerable degree of parallelism, this allows one to use the reversibly jump algorithm along with a very large number of processing units to minimize the wait in wall clock time before predictions are available.

CHAPTER 3

EXPLORATION OF CHEMICAL SPACE

3.1 Introduction

In order to jump from one chemical state to another, we must be able to propose a chemical state to visit. In the language of the Metropolis [33] algorithm, I will denote this proposal density $q(\cdot | k)$. For the purpose of the correctness of this algorithm, the formal requirements that I will place on $q(\cdot | k)$ are:

- The proposal must be reversible: $q(k'|k) > 0 \Rightarrow q(k|k') > 0$
- For every pair of states k_0 and k_N , $\exists \{k_0, \dots, k_N\}$ s.t. $\prod_{i=0}^{N-1} q(k_{i+1}|k_i) > 0$

The first condition implies reversibility: that is, if the algorithm can propose to jump forward, it can (possibly with a different probability) propose to jump backwards. The second condition implies that there are no states that are ultimately unreachable from any other state. Although this seems straightforward, and many obvious choices are apparent, the quality of the proposal distribution has a profound impact on the performance of the algorithm [86] and many adaptive techniques have been developed in other MCMC applications for this reason [71]. Not only does the proposal distribution need to satisfy the above basic requirements, it also must propose states to which a transition is feasible (in simulation terms) as well as provide the opportunity for far enough jumps that the simulation effectively explores chemical space. Feasibility in simulation terms has multiple aspects. One such aspect is that the remainder of the algorithm (the dimension matching and the nonequilibrium switching, specifically)

must be capable of completing the proposal. Importantly, this involves how the atoms will be mapped between the current and proposed states. Recall that the systems corresponding to chemical states k and k' do not typically have the same dimensionality. However, it is a choice of the algorithm to choose which dimensions should be considered in common between the two endpoint systems, and which should be considered to be unique to the endpoints. This choice—the atom mapping problem—will have a profound impact on the efficiency of the overall algorithm, as discussed below. However, before concerning ourselves with the efficiency of the algorithm, the atom map must at least ensure that the proposal can be completed. The requirement imposed by the dimension matching algorithm (discussed in Chapter 4) is that there must always be at least three atoms with positions (that is, either mapped or proposed by the dimension matching) forming a dihedral angle with each unique atom. The nonequilibrium switching implementation requires that constraint lengths do not change. In principle, this is not a formal requirement, but to avoid computing the necessary Jacobian, we do not map hydrogens, as typically only bonds to hydrogens are constrained. To add to these challenges, in the terms of the simulation’s operator, one would like to explore regions of chemical space that are feasible to purchase or synthesize as determined by models [90], or otherwise avoid regions that are unfavorable for reasons besides binding free energy.

3.2 Quantitative Metrics

In principle, one might choose to base proposal probability from state k to state k' on the thermodynamic length of the path from $p(x, k)$ to $p(x', k')$. This quantity,

described at length in [17, 80, 96], is given by:

$$\mathcal{L} \equiv \int_0^1 dt \sqrt{\left. \frac{d\lambda_i}{dt} \right|_{\lambda} g_{ij}(\lambda) \left. \frac{d\lambda_j}{dt} \right|_{\lambda}} \quad (3.1)$$

$$g_{ij}(\lambda) \equiv \mathbb{E} \left[\frac{\partial \ln p(x)}{\partial \lambda_i} \frac{\partial \ln p(x)}{\partial \lambda_j} \right] = \int dx p(x) \frac{\partial \ln p(x)}{\partial \lambda_i} \frac{\partial \ln p(x)}{\partial \lambda_j} \quad (3.2)$$

where $\lambda(t)$ is a multidimensional control parameter that in this case interpolates between $p(x, k)$ for $t = 0$ and $p(x', k')$ for $t = 1$ along a predefined path. The quantity in Eq. 3.2 is clearly the Fisher information metric; the quantity in Eq. 3.1 is known as the thermodynamic length, and it captures a sense of how much work will be done by the transition attempt. Importantly, the Cramer-Rao inequality reminds us that the variance of any unbiased estimator is lower-bounded by the inverse of the Fisher information [69]; to put this into terms relating to this work, the larger the differences between distributions in each step of the transition from k to k' , the higher the lower bound on the variance of the acceptance probability. It seems intuitive then that we would want to choose transition protocols that minimize this quantity; however, a quick inspection reveals that it is at least as difficult to compute as the relative free energy. It thus seems impractical to propose states based on the thermodynamic length, since this is in practice very difficult to compute. However, as an aside, it may be possible to learn this quantity from data. Additionally, some other works have devised sophisticated approximate sampling algorithms to find near-optimal paths [29]. Having completed many simulations, the thermodynamic length for each may be estimated, and fed to a machine learning model to predict thermodynamic lengths in future calculations. Although this sounds appealing, it is important to note that the thermodynamic length depends not just on the identity of the chemical ligand, but also on its environment (for instance, a protein), which

will vary depending on the calculation. More work is required to explore this approach.

3.3 Practical Approaches

Despite the inherent difficulty in choosing chemical states as above, there exist several useful heuristic approaches to which we may resort. Note that in the process of transitioning from one chemical species to another, we must identify which degrees of freedom belong to both systems, and which must be added or deleted. We can use the number of degrees of freedom in common as a basis for the proposal distribution, under the assumption that the more degrees of freedom in common between the endpoint systems, the more likely the transition is to be accepted.

3.3.1 Chemical State Proposal Algorithm

In order to propose a chemical state k' from a current state k , we first must derive the proposal probabilities. Since atom mapping will be essential to the performance of the algorithm, we choose to use the number of atoms in an atom map as proportional to the proposal probability. The atom maps themselves are derived from a maximum common substructure search [25] as implemented in the OpenEye toolkit. This algorithm attempts to find the largest contiguous substructure that is present in both molecule graphs, and has several hyperparameters that can be tuned or optimized. More specifically, upon initialization, the proposal

algorithm performs the following steps, given a set of molecules $\{\mathcal{M}_0, \dots, \mathcal{M}_N\}$ corresponding to chemical states $\{k_0, \dots, k_N\}$:

- For each pair of molecules $(\mathcal{M}_i, \mathcal{M}_j)$, perform a maximum common sub-structure search (MCSS) [25] to obtain a set of maps. In order to be valid, the map must contain at least three contiguous atoms, and must not map constrained bonds of different lengths.
- Choose the mapping of atoms (that is, the list of which atoms should be considered degrees of freedom in common) with the maximum number of atoms.
- The probability of proposing j given current state i , termed $q(j | i)$, is proportional to the number of atoms in common in the MCSS map.

When running the simulation, the algorithm performs a proposal as follows:

- Perceive the current chemical state of the simulation k
- Propose a state k' from $q(\cdot | k)$
- Compute atom map and $\ln P_{\text{forward}} = q(k' | k)$
- Compute $\ln P_{\text{reverse}} = q(k | k')$

As described above, there are several hyperparameters here that can profoundly impact the performance of the algorithm. Foremost among them are the hyperparameters of the MCSS search: which atoms should we count as being in common? Second, should the algorithm allow the creation and deletion of partial rings? Although this would result in a larger atom map, it may be less favorable. Third, the algorithm can optionally also sample from the list of atom maps, if

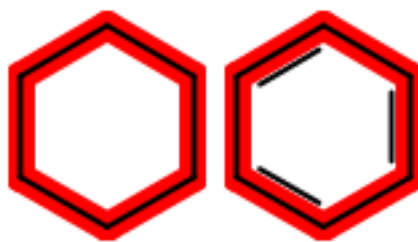


Figure 3.1: An example of an apparently reasonable atom map between benzene and cyclohexane.

there is more than one result. Finally, the map produced must be usable by the following components of the algorithm.

3.3.2 MCSS Hyperparameters

When performing an MCSS search on chemical graphs, one must recognize that it is not simply the graph structure that is important, but also the labels on the edges and vertices, which correspond to bonds and atoms, respectively. To illustrate the consequence of improper atom maps (for instance, that disregard bond order), note Figure 3.3. The simple inclusion of the bond order as a criterion for atom mapping greatly enhances the geometry that is produced, as in Figure 3.4

3.4 Factors that affect acceptance probability

For example, in Figure 3.1, there are two molecules with six-membered carbon rings. Mapping these would seem intuitive based on graph structure, but the nature of the edges makes this prohibitively unfavorable. The 3D structures of these molecules, visualized in Figure 3.2, demonstrate that this is highly



Figure 3.2: 3D structures of benzene and cyclohexane. Note that benzene is flat, while cyclohexane adopts a "chair" conformation.

infeasible—the initial structure of either will be very unfavorable using the potential of the other. In this case, chemical intuition would make the problem obvious: benzene is aromatic, while cyclohexane is not. In many other cases, the distinction may not be clearly obvious. In Figure 3.5, we see an example of the MCSS algorithm maximizing the overlap between two molecules in terms of the atom map. However, it may be naive to attempt this, as the dimension-

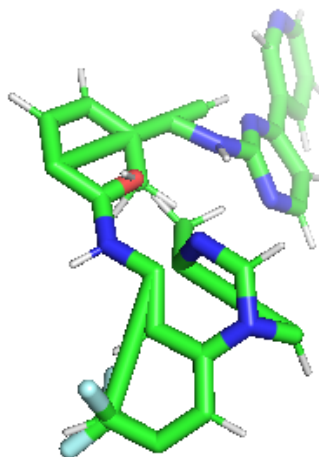


Figure 3.3: A geometry that was constructed from an atom map that left no space to properly close the rings.

matching component of the algorithm will now have to complete the ring, and the annealing will have to gradually introduce part of the ring. In this instance, the random choice of atom maps would be helpful, as the algorithm would occasionally not choose the ring-breaking map. An attractive alternative to randomly choosing an atom map is to simply filter out maps that partially map rings. While the toolkit that I used for this project unfortunately did not provide this as an option, it is possible to simply eliminate the maps that break rings. The algorithm for determining whether rings are broken is as follows, with a return value of True signifying that rings are broken:

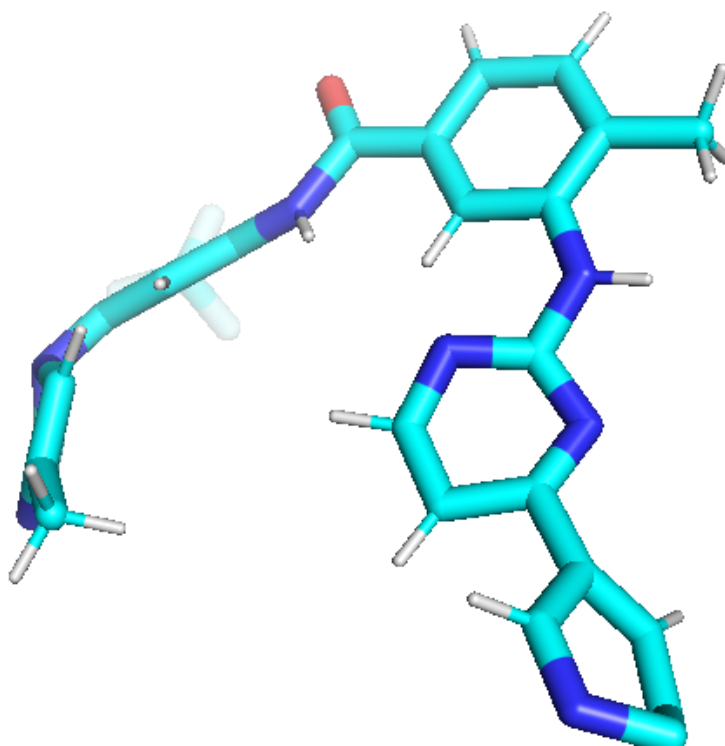


Figure 3.4: A geometry resulting from the inclusion of bond order as a criterion for mapping atoms. More of the system is rebuilt, so a better geometry can be created.

- For a pair of molecules (k, k') , enumerate the cycle basis of both molecular graphs [65]
- For each bond in a cycle of molecular graph k , check that it is in graph k' , if any are not, return True
- If all bonds in cycles of k are also in k' , return False

This algorithm checks that rings are not created or destroyed, and is done using the NetworkX toolkit [32]. However, building an entire ring is not infeasible, and so the algorithm could be modified to check that if a ring is created or broken, an entire ring is created or broken.

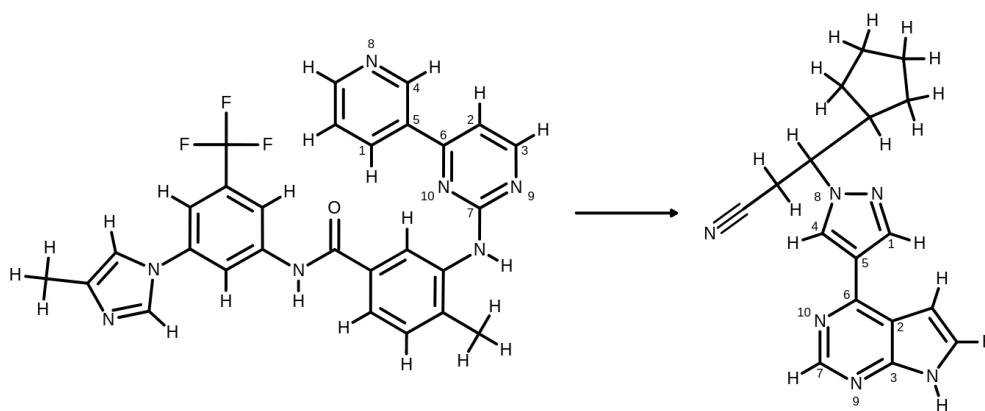


Figure 3.5: This map between kinase inhibitors may contain more atoms than another, but as it is partially breaking a ring, energies after the proposal are often highly unfavorable

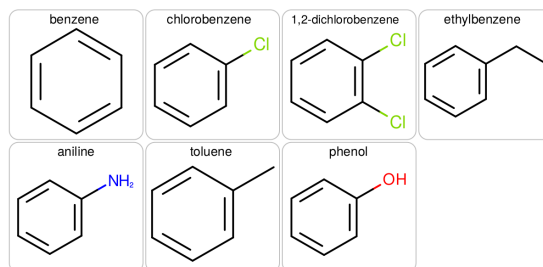


Figure 3.6: The set of trial molecules used for empirical exploration of atom map criteria. This set was used to demonstrate because it contains very similar cores, but different substituents that might affect the performance of atom maps.

3.4.1 An exploration of the effect of various atom mapping options on acceptance probability

As an interesting exploration, I measured the empirical effect of different atom mappings on the variance of the instantaneous acceptance probabilities between various substituted benzene molecules as shown in Figure 3.6. In this experiment,

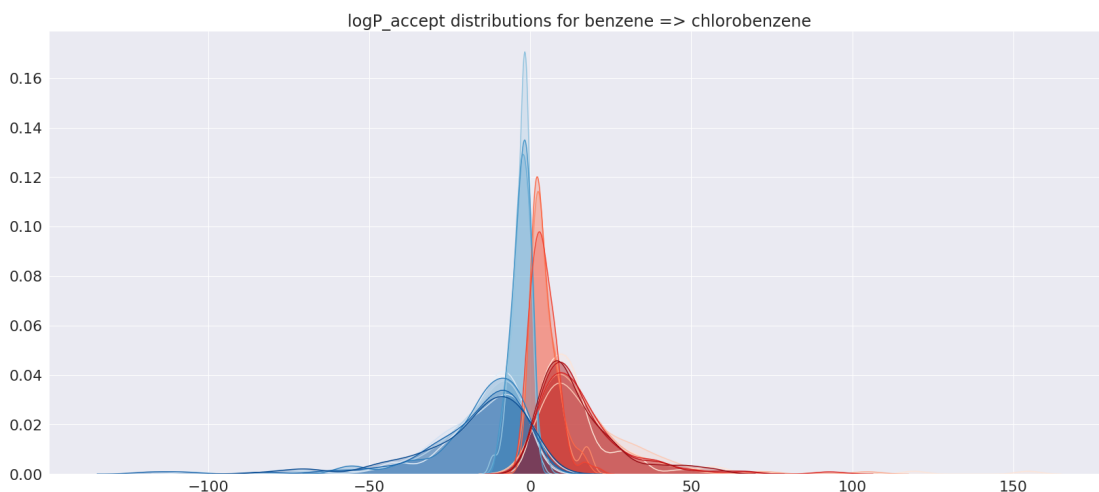


Figure 3.7: Forward (blue) and negative reverse (red) log acceptance probability distributions for benzene to chlorobenzene in vacuum under different mapping schemes. Here, it is clearly visible that the choice of mapping scheme has a profound effect on the ultimate quality of the proposal.

I ran the molecules at a timestep of 1 femtosecond for 1 nanosecond at 300 Kelvin in vacuum, taking a snapshot of the trajectory every picosecond. For each other molecule in the set, I attempted 100 reversible jump proposals under each atom matching scheme, starting from a configuration randomly drawn from the equilibrium cache of the starting molecule. As is readily apparent from Figure 3.7, the choice of map can have a significant effect on the quality of the proposal. However, for some other pairs, the choice of map does not appear as important. For example, in Figure 3.8, the various maps produce distributions of acceptance probabilities that cluster around the the same region. As such, it is an area of active investigation to determine how to map atoms in a way that maximizes the efficiency of the computation. It is clearly of interest in practical applications to understand what the ideal atom mapping criterion for explicit solvent is, and whether this differs substantially from the criteria chosen in vacuum.

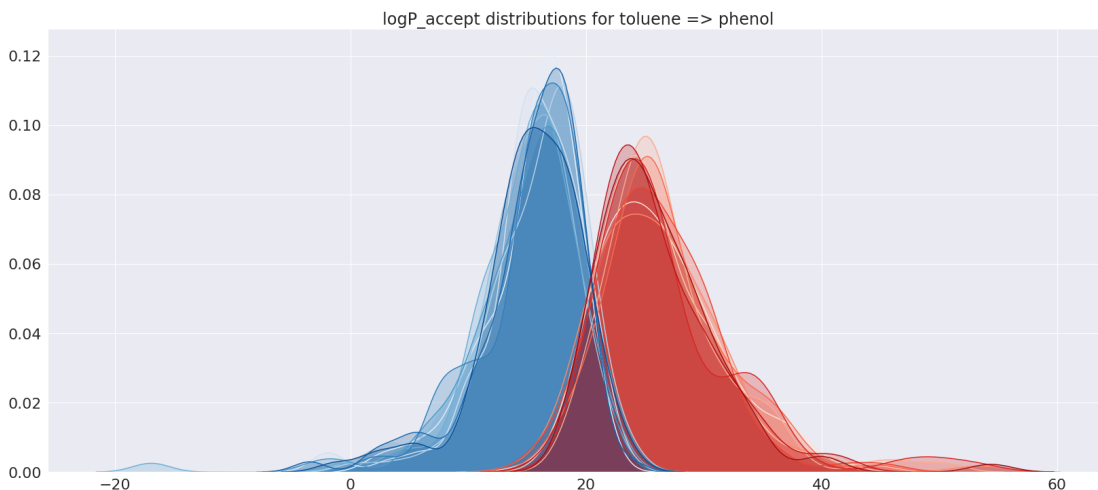


Figure 3.8: Forward (blue) and negative reverse (red) log acceptance probability distributions for toluene to phenol in vacuum under different mapping schemes. Here, the choice of the mapping scheme does not seem to have as significant of an impact on the quality of the proposal.

3.5 Future work

3.5.1 Atom map learning

For OpenEye MCSS, there are multiple atom map options. For instance, the user can choose to have the MCSS algorithm match atoms based on aromaticity, element, number of heavy bond partners, and more. In principle, these discrete choices can be enumerated and searched to find near-optimal combinations. This does leave open the question of what to use as the objective. In principle, the variance of the acceptance probability is ideal, though it may be difficult to estimate. Alternatively, one could try to find the set of discrete hyperparameters that maximize the average acceptance probability, which is straightforward

to compute. This does, however, leave the issue that a set of atom mapping parameters ideal for one set of molecules may not be performant for another.

Treat atom mapping as a Reinforcement Learning problem

As a result, one could imagine treating the atom mapping as a reinforcement learning problem, wherein an agent must select a set of parameters given the transformation at hand [81, 92]. In this way, one might hope that the algorithm could learn general rules that would be valid across many different sets of molecules. One advantage of this approach is that, under certain conditions [4, 2, 5], one could potentially design an adaptive MCMC algorithm that performs this learning as the calculation proceeds. This would potentially allow the user to avoid a costly pre-optimization (or the use of suboptimal parameters derived from a reference set).

Atom mapping with neural network

Finally, one might also simply leave the atom mapping choice to a neural network, leaving open the possibility to efficiently learn a function mapping from transformations (or even environments) to atom maps. Although the map itself is discrete, other sophisticated differentiable neural architectures such as the differentiable Forth interpreter [8] were able to overcome this.

3.5.2 Relieving the restriction to a finite set of chemicals

Ideally, the simulation would not be limited to a prespecified set of chemical states. Rather, one could imagine applying synthetic transformation rules to generate new, feasible chemicals. This has the advantage of resulting in chemical matter that is more likely to be feasible to synthesize (or at least more likely to be able to convince a chemist to try). Alternatively, one could imagine using an RNN such as in [31], where a neural network might learn to generate new chemicals. In this way, we might not only adaptively solve the issue of which chemical states are "neighboring," but also push toward synthetic novelty in a way that humans have not been able to. Although there are challenges in this part of the algorithm, many of the challenges for expanding beyond a fixed set of chemical species lie in the weight adaptation algorithm, discussed later. The development of this algorithm opens up new avenues for exploration, as well as new avenues for multidisciplinary collaboration with other fields.

CHAPTER 4

GEOMETRY PROPOSALS

4.1 Introduction to dimension matching

Although Reversible Jump MCMC [30] provides a very general and convenient framework for sampling spaces with varying dimensionality, there are still many practical and problem-specific considerations. Neglecting other components of the acceptance probability, the contribution to the acceptance probability involving dimension matching in the proposal from $(x_{old}, k) \rightarrow (x_{new}, k')$ is given by

$$\ln P_{\text{geometry}} = [-u_{k'}(x_{new}) + u_k(x_{old})] + \ln \frac{\phi(x_{old} | x_{new})}{\phi(x_{new} | x_{old})} \quad (4.1)$$

where $u_{k'}$ is the reduced potential for chemical state k' , u_k is the reduced potential for chemical state k , x_{new} is the new proposed configuration, x_{old} is the old configuration, and $\phi(x, y)$ is the dimension matching distribution. Notably, x_{new} and x_{old} can have different dimensionalities. Intuitively, this involves ensuring that the spaces on which the proposal and the current distribution lie are the same. In other words, we must augment each probability distribution with another distribution that contributes the missing degrees of freedom for each. However, although the recipe is in principle straightforward, it is difficult to construct an efficient dimension matching distribution for the case of chemical space sampling. Among the challenges for constructing an efficient distribution are the multimodality of the target distribution, the need to draw an exact sample, and the need to be able to calculate a normalized probability for the proposal.

4.2 Problems faced by dimension matching in chemical space

4.2.1 The target is highly multimodal

Consider a flexible alkane: a terminal methyl group could be rotated to produce configurations of roughly even probability. Although one might not consider this to be a significant problem, consider the terms in Eq. 4.1. The positive contribution of the reverse proposal probability (that is, the probability of placing the atoms being deleted where they were found) means that the proposal distribution must not simply place atoms in a favorable position; it must also recognize favorable positions as such. To use a more concrete example, consider the pair of distributions in Figure 4.1. Although the proposal distribution would clearly place a sample in a favorable location, a sample drawn from the target (such as from a Langevin dynamics simulation) would likely have low probability under the proposal. This prevents algorithms that, for instance, use a single reference torsion for proposals. As a result of these complex interactions and multimodal distributions, many simple sampling schemes are ruled out.

4.2.2 The proposal must be drawn exactly

Another choice that may seem attractive is to simply run an MCMC chain for the missing degrees of freedom and use that sample for the new coordinates. However, in order for this to be drawn *exactly* from the target distribution, the Markov chain would have to be run infinitely long (except in the case of perfect MCMC [9], which is not applicable here). This is clearly infeasible, and leaves

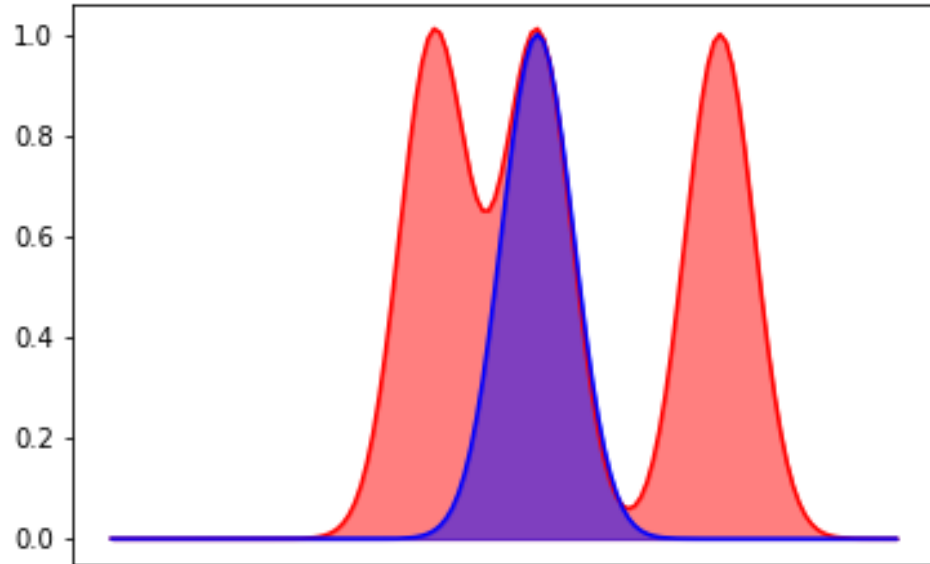


Figure 4.1: An example of a unimodal proposal (blue) with a multimodal target (red). Note that there are many regions of the target that are poorly covered by the proposal.

us without familiar tools that we use to sample high-dimensional multimodal probability distributions.

4.2.3 The proposal must be associated with a normalized probability

Yet another temptation at this point may be to simply draw from a tractable distribution as in the Hastings algorithm [33] (for example, a Gaussian), and then minimize the configuration's energy. However, performing such a nonlinear transformation on the proposed configuration would require that we be able to

invert and differentiate the transformation—the determinant of the Jacobian matrix of the transformation is required to calculate the probability. Since we cannot straightforwardly compute such a Jacobian, techniques involving complex nonlinear transformations are largely impractical in this case.

4.2.4 Atomic positions are correlated

Further complicating the matter is the fact that various degrees of freedom in a molecule are correlated. For instance, steric hindrance prevents atoms from being placed close to one another. This means that although *a priori* an atom might have a multimodal probability distribution, conditioned on a previous atom, some modes are now precluded.

4.2.5 Propose one atom at a time

Since proposing a set of atoms at once is fraught with difficulty as described above, one alternative is to propose one atom at a time, significantly reducing the dimensionality of the proposal distribution. However, the challenges described above (especially multimodality and correlated atomic positions) still remain. Following other literature in Grand Canonical Monte Carlo (GCMC) [79] as well as in other molecular simulation literature and statistical literature [15], we propose positions in so-called internal coordinates. In this coordinate system, atomic positions are defined by a bond length r , a bond angle θ , and a dihedral angle ϕ as described in Figure 4.2.

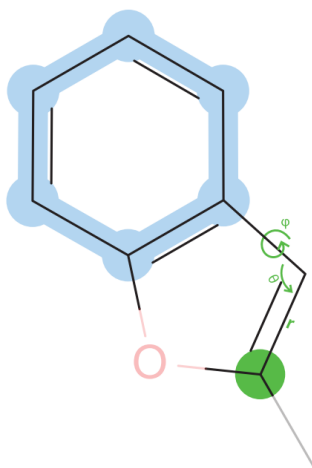


Figure 4.2: An example of an internal coordinate description of atomic positions. The atom highlighted in green is being proposed, based on the bond (r), angle (θ), and torsion (ϕ).

4.2.6 Naive bond/angle/torsion doesn't work

Having moved to internal coordinates, the most obvious proposal scheme would be to choose a dihedral angle associated with each atom, and directly propose a bond length, bond angle, and dihedral angle based on the forcefield parameters. However, this has a major pitfall: Although it works with simple geometries and cases where very few atoms need to be added, it fails quickly with complex geometries such as rings, as well as cases where many atoms must be added. In the latter case, one can imagine that neglecting other dihedral and bond angle terms can easily lead to very unfavorable configurations for complex molecules. As a refinement of this algorithm, we include other angle, dihedral, and even nonbonded terms in the proposal scheme.

4.3 CBMC-like Algorithm

For the common dimension-matching tasks, I implemented and refined a CBMC-like [79, 3, 15] algorithm that can take into account forcefield terms besides one bond, angle, and torsion term.

4.3.1 Description of algorithm

The CBMC-like algorithm with guide torsions has several components as follows.

At a high level, for a given transformation $k \rightarrow k'$:

- Determine which degrees of freedom, based on the atom map, are present in k' not in k
- Separate atomic degrees of freedom into heavy atoms and hydrogens
- Perform a breadth-first search of the molecular graph, identifying the order in which heavy atoms can be proposed
- Repeat for hydrogen atoms
- For each new atom, propose a bond length, bond angle, and dihedral angle, and convert to cartesian coordinates
- Accumulate the log probability of the acceptance probability (described in detail below)

Determining the atomic proposal order

After determining which atoms are not mapped, the dimension-matching algorithm must first determine the order in which to propose the new atoms, as well

as which internal frame of reference will be used. In greater detail, beginning with a $\ln P_{\text{choice}} = 0$, a set of atoms with positions and without positions, and until all new atoms are added:

- Add atoms which do not have positions, but are connected to a bond, angle, and dihedral partner with positions to `atoms_eligible_for_proposal`
- Uniformly choose without replacement atoms from the `atoms_eligible_for_proposal` list, and add the probability of choosing each atom to $\ln P_{\text{choice}}$
- For each atom that is chosen, uniformly choose a dihedral angle, and add the probability of this to $\ln P_{\text{choice}}$
- Add each atom that is chosen to `atoms_with_positions`

The above algorithm is first conducted for the heavy atoms, and then for the hydrogens. When the reverse proposal probability is being evaluated, the same procedure is repeated with stochasticity. It is noteworthy to add here that the use of this algorithm for determining proposal order imposes an extra requirement on the atom map: it must contain at least one atom with a dihedral that can be used for proposal.

Coordinate proposal algorithm

Once the proposal order and the reference frame for each atom has been chosen, the algorithm can now stochastically propose coordinates for the new atoms. For each new atom, along with its corresponding dihedral, the proposal algorithm works as follows:

- If the bond is not constrained, propose bond length $r \sim p(r) = \mathcal{N}(r_0, \sigma_r^2)$, where $\mathcal{N}(\mu, \sigma^2)$ denotes the normal density with mean μ and variance σ^2 , and $\sigma_r = (\beta K_r)^{-1/2}$ where K_r is the bond force constant; if the bond is constrained, set r to its constraint length r_0
- Propose bond angle $\theta \sim \mathcal{N}(\theta_0, \sigma_\theta^2)$, where $\sigma_\theta = (\beta K_\theta)^{-1/2}$, where K_θ is the angle force constant
- Propose dihedral angle $\phi \sim p(\phi; r, \theta)$, where $p(\phi; r, \theta)$ is defined below
- Calculate $(x, y, z), \ln \det J(r, \theta, \phi) \leftarrow \text{internal_to_cartesian}$, where $J(r, \theta, \phi)$ is the Jacobian describing the hypervolume $dx dy dz$ given $dr d\theta d\phi$, where `internal_to_cartesian` converts from internal coordinates to cartesian coordinates.
- Compute the log contribution to the acceptance probability, $\ln P_{\text{atom}} = \ln [p(r) p(\theta) p(\phi; r, \theta) J(r, \theta, \phi)]$

Note that in Algorithm 4.3.1, we are able to use the forcefield's distribution (harmonic) for the bond and angle terms. However, the dihedral conditional density $p(\phi; r, \theta)$ —in the absence of any other nonbonded interactions involving the placed atom—has the form

$$p(\phi; r, \theta) \propto \exp \left[-\beta \sum_{l=1}^L \frac{K_{\phi,l}}{2} \cos(n_l \phi + \gamma_l) \right] \quad (4.2)$$

which does not have a closed-form normalizing constant; here, $K_{\phi,l}$ is a barrier height, n_l is an integral periodicity, and γ_l is a phase for the l th Fourier term. It is straightforward to numerically normalize this, however, and draw samples using rejection sampling, since the distribution is only one dimensional. Simply drawing from the conditional dihedral term is attractive, but generally affords poor performance, as it neglects most other terms. For linear alkanes and small

additions, this approach is feasible. However, for more complex molecules, this approach suffers from several drawbacks. First, it neglects other valence terms such as the other dihedrals and angles in which the atom in question is involved. Second, it will very often fail to properly create rings, since without the full set of bonds, intermediate dihedral angles would be just as likely to leave the ring open, causing extremely unfavorable valence interactions. Finally, it does not afford the dimension matching scheme an opportunity to incorporate local nonbonded terms.

Improvements to dihedral proposal density

Rather than simply normalizing a single dihedral term and proposing from that, one can instead perform a drive of the angle about its dihedral, and compute the potential energy at each point. This potential may be the full potential, or, for efficiency reasons, a subset of the full potential. Then, having performed the dihedral scan, one can normalize the resulting potential to form a proposal density. An additional hyperparameter of this algorithm is how finely discretized the dihedral angles must be.

4.3.2 Performance on simple cases

There are a number of different transformations that could be examined.

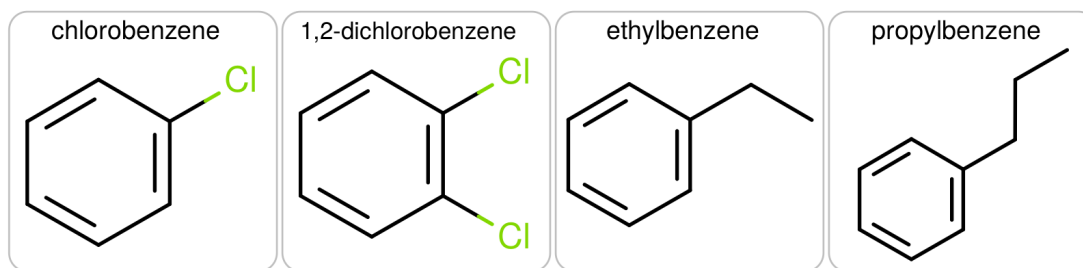


Figure 4.3: Examples of substituted benzenes

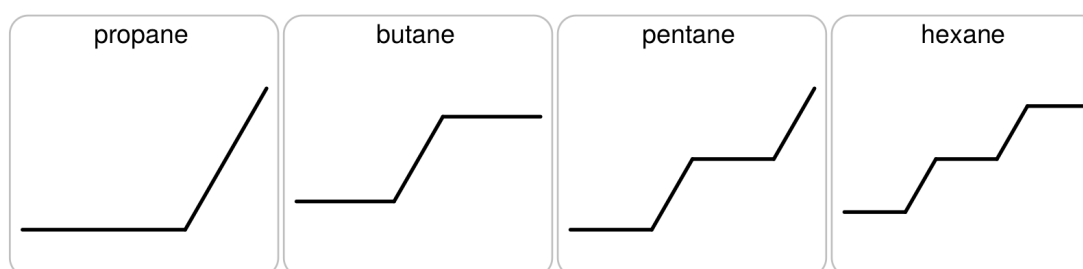


Figure 4.4: Examples of n-alkanes

Substituted Benzenes

Among the simplest are transformations between different substituted benzenes as shown in Figure 4.3. These transformations are straightforward since the molecules are fairly rigid, and there are only several atoms to insert.

Alkanes

Another type of transformation that is relatively straightforward is that between alkanes, as shown in Figure 4.4. However, these molecules are more flexible, and thus the dimension matching distribution must accurately capture the multimodality of the configurational probability distribution.

More complex transformations

Of course, while alkanes and substituted benzenes are useful model systems and the building blocks of many interesting compounds, we are far more interested in complex transformations, such as the addition of rings or jump between different drug-like molecules. In order to achieve these, however, we must further improve our strategy. One immediate difficulty, even with the improved dihedral scan, is that building a ring will often result in failure. The failure to close a ring results from the fact that when the intermediate ring atoms are being inserted, there are no terms that cause the ring to close, resulting in a grotesque geometry. To remedy this, we can add "guide torsions" that encourage the ring to close by removing the dihedral modes that correspond to an open ring as in [93]. The algorithm for the addition of guide torsions works as follows:

- Convert the molecule into the OpenEye OEMol format
- Perceive rigid bonds
- Generate a reference configuration using OpenEye Omega
- Measure the dihedral angles about the rigid bonds
- Add a dihedral term with a mode at the angle in the reference conformation

By applying the steps in Algorithm 4.3.2, we can significantly increase the chance of proposing a closed ring, resulting in a more realistic geometry.

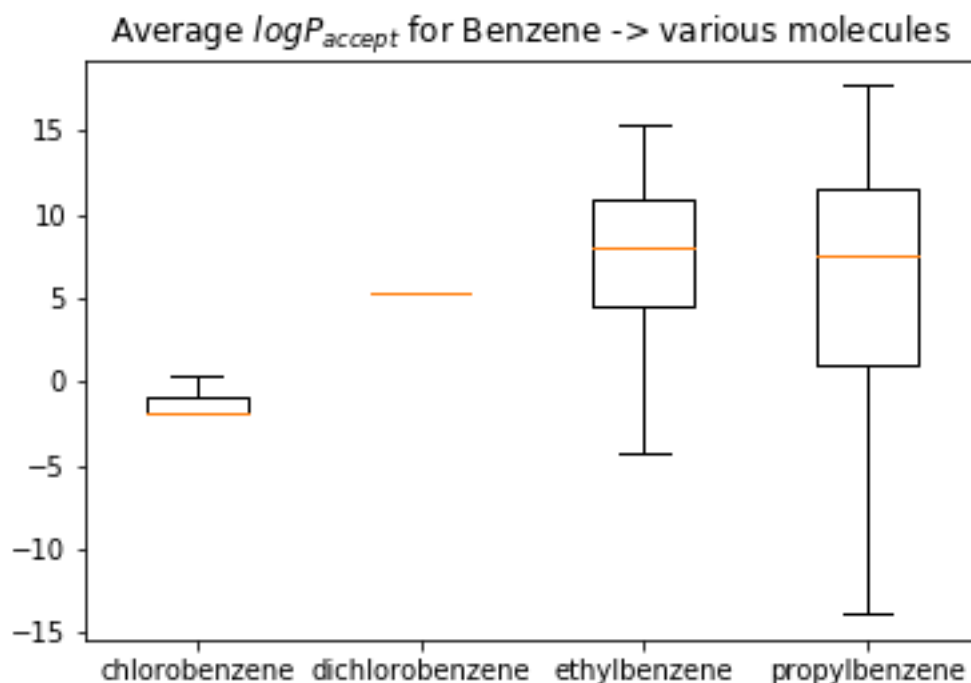


Figure 4.5: The average $\ln P_{\text{accept}}$ vs. the number of degrees of freedom added. Note that the variance is quite low for small (chlorobenzene and dichlorobenzene) changes, but quickly grows with the addition of flexible chains. Note that anything beyond the 98th percentile was clipped for this figure.

4.3.3 Shortcomings

However, with very large changes (involving the addition of many atoms), this approach is still not terribly performant. As an example of Sequential Importance Sampling (SIS), it suffers from the compounding of poor choices made earlier in the sequence. As the number of degrees of freedom that must be added grows, this problem becomes more serious, demonstrated in Figure 4.5. When the approach is tried on even more complex transformations, such as between kinase inhibitors, we observe even more severe performance loss, making acceptance a

highly unlikely proposition. However, there are yet remedies to these serious problems.

4.4 Future directions: Particle Filtering

4.4.1 Basic Overview

Instead of simply choosing one dihedral angle at each atom in 4.3.1, we can choose N dihedral angles and calculate an importance weight based on the difference between the probability of that atom's placement under the target distribution. Then, after a specified number of atoms are inserted, the algorithm resamples N of the configurations based on their weights. This technique is called particle filtering [58, 34], and is commonly used in state space models.

4.4.2 Advantages

Particle filtering can enable the algorithm to eliminate proposals that were reasonable at stage m but become very unfavorable at stage $m + n$. For example, if a particle (as the replicates of configurations are known in statistics) does not close a ring, when the ring closure should happen, its weight will be very small and it will be lost in the resampling. On the other hand, naive multinomial resampling carries with it its own risks: since at every resampling step there is a nonzero probability for favorable configurations to be lost, after enough resampling steps, the probability that favorable configurations are lost can grow quite high. Mitigating this are other resampling algorithms, which aim to reduce

the deleterious effect of multinomial sampling In addition to these advantages, particle filtering can be parallelized on tightly-coupled parallel hardware such as graphics processors [40].

4.4.3 Algorithm Hyperparameters

As briefly mentioned above, the particle filter has a number of hyperparameters that can be tuned and tweaked. Among these are the number of particles or replicates, the method used for resampling, and, as above, the method used for proposing the individual bonds, angles, and dihedrals.

Number of Particles

The most obvious algorithm hyperparameter is the number of particles. The larger this number is, the better the performance of the dimension matching algorithm. However, each particle requires the computation of the target density, which can be quite expensive. Because cost scales relatively quickly in this case, it would be ideal to minimize the number of particles needed. Furthermore, a very large number of particles can cause serious numerical instabilities [59], primarily from the need to normalize weights before resampling.

Resampling Algorithm

One can choose various approaches for the resampling step of the particle filter. The simplest is known as multinomial resampling, and consists of simply drawing N particles with replacement from the previous N according to their

weights. This is known to have a relatively high variance [23]. Another approach is so-called stratified resampling, which This also requires weight normalization, however. In situations with large numbers of particles, it may be feasible to resort to a scheme that is “almost” exact, such as Metropolis resampling. [59] In the Metropolis resampling scheme, one proposes to sample a particular ancestor and accepts/rejects. One then repeats this procedure many times until approximate convergence. These steps must be performed at each resampling attempt. This is not technically exact, as Metropolis MCMC only converges in the asymptotic limit, not for finite samples. However, in some studies [59], it actually has less bias than the exact methods, because weight normalization is not required and certain numerical issues are resolved Parallelization is also facilitated by Metropolis resampling, since the weights do not need to be normalized (a reduction operation) before resampling.

Overall comparison

Finally, below is a table which compares the features of various proposal schemes that are in principle correct for the dimension matching algorithm.

4.5 Tuning of Dimension Matching Parameters in the context of NCMC

It should be noted that although the parameters for the dimension matching proposal are very important, there are several tradeoffs between tuning the parameters of the dimension matching and those of the NCMC or annealing, de-

Table 4.1: **Comparison of features of various proposal schemes for dimension matching.**

Method	Multimodal	Exact	computable P	Rings	Correlated	Speed
Minimization	Yes	Yes	No	Yes	Yes	Fast
Unimodal	No	Yes	Yes	No	No	Fast
Multi-reference	Yes	Yes	Yes	Yes	No	Fast
CBMC-like	Yes	Yes	Yes	No	Moderate	Moderate
CBMC+guide	Yes	Yes	Yes	Yes	Moderate	Moderate
Particle Filtering	Yes	Yes	Yes	Yes	Yes	Slow

scribed in the next chapter. The implementation details of each are quite different, leading to interesting tradeoffs. For instance, although the energy computations in the dimension matching algorithm are much faster (they exclude the majority of the system’s interactions), the annealing runs on the GPU, which often more than compensates for the added terms in the energy computation. An additional point to note is that the dimension matching (especially particle filtering) can be subject to numerical stability issues. Because the annealing protocol is typically very slow, and steps of Langevin Dynamics are taken between changes in the control parameter, the NCMC component of the overall algorithm may be better behaved. For this reason, it is useful to jointly optimize the parameters.

4.5.1 Tradeoffs

There is to some extent a tradeoff between the hyperparameters of the dimension matching algorithm, and those of the annealing. Intuitively, one might imagine that additional grid points and the use of nonbonded forces in the dimension

matching scheme might result in the need for fewer steps of NCMC. However, due to the current setup of the NCMC, this is not terribly helpful. Nonbonded forces are initially disabled for the new atoms (except intramolecular nonbonded forces), and so using the contributions of intermolecular nonbonded forces to proposed new geometries is likely to result in an inferior acceptance probability. As to the number of grid points, this could potentially be reduced with longer NCMC by including valence softening in the NCMC. However, the default configuration of NCMC does not include this, and so the number of grid points is largely decoupled from the NCMC protocol length. However, another potentially fruitful approach is to allow the NCMC code to soften the bonds and angles at the initial value of the control parameter. This allows otherwise unfavorable geometry proposals to be accorded a favorable energy, then slowly anneal back to the forcefield's parameters for those terms. It should be noted that doing this would then introduce a dependency between the geometry proposal's hyperparameters and the NCMC protocol. For example, it may result in the need for longer protocols, or may cause very unfavorable reverse geometry proposal probabilities. To understand the latter point, consider that if the NCMC protocol is symmetric, the atoms to be deleted will conclude the protocol with soft bonds and angles, thus distorting the geometry from what the forcefield would otherwise cause. However, the geometry proposal distribution is simply using the forcefield's parameters, and as such will likely assign low probability to the ending configuration. This can be straightforwardly ameliorated by modifying the geometry engine to also be able to soften its bond and angle proposal distributions, thereby once again matching the distribution seen in the nonequilibrium switching.

CHAPTER 5

NCMC SWITCHING

5.1 Introduction

If new atoms are introduced into a condensed-phase simulation, it is highly likely that they will be placed unfavorably close to other atoms. Additionally, if a ligand is transforming to another ligand, valence and nonbonded terms may also need to transform, as demonstrated in Figure 5.1. In order to mitigate these very serious issues, we turn to Nonequilibrium Candidate Monte Carlo (NCMC) [62]. NCMC is a technique that allows us to gradually interpolate the system from one set of control parameters to another. In this particular case, atoms are introduced after dimension matching as non-interacting dummies, and core atoms are left with their original parameters. Then, the algorithm alternates between incrementing the control parameters and taking one step of Langevin dynamics (or other MCMC algorithm that either maintains the invariant distribution, or accounts for deviations by including an importance weight). This allows new atoms to be introduced smoothly, and allows core atoms to be smoothly interpolated from one set of parameters to another.

5.2 Description of Algorithm

In broad terms, the algorithm performs the following steps:

- Given system A, system B, and an atom map between them, generate hybrid system H with control parameters $\lambda \in [0, 1]$

- Perform annealing or nonequilibrium switching, changing control parameters from 0 to 1 gradually and accumulating the change in energy
- Transfer the positions to system B

5.2.1 Construction of Hybrid System

The first step is to construct a system that is the hybrid of the two endpoints. The object of constructing this system is to be able to smoothly interpolate from the first endpoint to the second. To that end, I build the system with modified energy terms: all of the terms are controllable through global parameters collectively known as λ . When $\lambda = 0$, the hybrid system behaves like the initial endpoint. When $\lambda = 1$, the hybrid system behaves like the final endpoint. The construction of such a system, as well as the path that one takes through the space of control parameters, are important factors in the performance of the overall algorithm. Here, I will describe in greater detail the way that each type of energy term for a classical forcefield is handled.

Bond force

In the AMBER forcefields [68], bond forces are represented by a harmonic term:

$$U(r) = \frac{K_r}{2}(r - r_0)^2 \quad (5.1)$$

where the K is a force constant and the r_0 is the equilibrium bond length. For bonds present in both endpoint systems in the region of the system that is being

modified, I use a slightly modified potential:

$$r_0(\lambda_{bonds}) = (1 - \lambda_{bonds})r_{0_A} + \lambda_{bonds}r_{0_B} \quad (5.2)$$

$$K_r(\lambda_{bonds}) = (1 - \lambda_{bonds})K_{r,A} + \lambda_{bonds}K_{r,B} \quad (5.3)$$

$$U(r; \lambda_{bonds}) = \frac{K_r(\lambda_{bonds})}{2} [r - r_0(\lambda_{bonds})]^2 \quad (5.4)$$

In this way, the hybrid system's bond energy linearly depends on the control parameter λ_{bonds} between the endpoints. We are also free to change the bond parameters nonlinearly with the master λ control parameter's schedule. For bonds that are not in the modified region, a standard harmonic term is used. Likewise, for bonds between atoms that only exist in one endpoint or the other, a standard harmonic term is used to prevent them from departing the region where they belong. A tricky issue arises when two atoms are present in both systems, but have a bond only in one. This arises in the case of transforming a non-ring into a ring. In this case, for the system without the ring, the bond force constant is set to 0. However, it should be noted that these proposals are likely not desirable, and should be remedied by reducing the number of atoms in the map. Finally, it is important to mention that when a constrained bond is encountered, its constraint length is not changed. This is to avoid the need to compute a Jacobian for this deterministic transformation, although it would in principle be possible. I have found that this issue can generally be avoided entirely by simply not mapping hydrogen atoms.

Angle force

Also included in AMBER force fields (as well as others) are angle terms between 3 atoms. These force terms are also harmonic, with the form:

$$U(r) = \frac{K_\theta}{2}(\theta - \theta_0)^2 \quad (5.5)$$

I treat these similarly to the bonds in 5.2:

$$\theta_0(\lambda_{angles}) = (1 - \lambda_{angles})\theta_{0_A} + \lambda_{angles}\theta_{0_B} \quad (5.6)$$

$$K_\theta(\lambda_{angles}) = (1 - \lambda_{angles})K_{\theta,A} + \lambda_{angles}K_{\theta,B} \quad (5.7)$$

$$U(\theta; \lambda_{angles}) = \frac{K_\theta(\lambda_{angles})}{2} [\theta - \theta_0(\lambda_{angles})]^2 \quad (5.8)$$

As with the bonds, I do not use the modified potential for angles that are not in the region being modified, nor do I use them for triplets of atoms that are only present in one endpoint. The ring closure issue is also present with the angles, and is handled similarly.

Dihedral force

Dihedral angles also have potential terms, of the form:

$$U(\phi) = \sum_{l=1}^L \frac{K_{\phi,l}}{2} \cos(n_l \phi - \gamma_l) \quad (5.9)$$

where ϕ is the dihedral angle, and for Fourier component $l \in \{1, \dots, L\}$, n_l is the periodicity, γ_l is the phase, and $K_{\phi,l}$ is the force constant. For the dihedrals present in both systems, I make the total energy a linear combination of the endpoint terms, dependent on a control parameter $\lambda_{torsions}$:

$$U(\phi; \lambda_{torsions}) = (1 - \lambda_{torsions})U_A(\phi) + \lambda_{torsions}U_B(\phi) \quad (5.10)$$

Similarly to other valence terms, dihedral terms that are between atoms solely in one endpoint are always active to prevent very unfavorable configurations from being reached.

Nonbonded terms

Perhaps most challenging of all the energy terms are the nonbonded terms. Not only is there a great diversity of schemes for implementing long-range nonbonded interactions, but these parameters can have a very profound impact on the acceptance probability. The nonbonded interactions of the AMBER [68] forcefield consist of two components: electrostatics and Lennard-Jones or sterics. The sterics component is a relatively short-range component with a form:

$$U(r) = 4\epsilon \left[\left(\frac{\sigma}{r} \right)^{12} - \left(\frac{\sigma}{r} \right)^6 \right] \quad (5.11)$$

where ϵ is the depth of the potential well (how favorable it is to be in the minimum energy configuration) and σ is the distance between the atoms where that minimum is located. The electrostatics component is long-ranged and takes the form of the Coulomb force:

$$U(r) = \frac{1}{4\pi\epsilon_0} \frac{q_i q_j}{r} \quad (5.12)$$

where $q_i q_j$ is the product of individual atomic partial charges and ϵ_0 is the vacuum permittivity. Although initially one might be tempted to simply linearly interpolate the nonbonded terms, there is an immediate issue: although the sterics provide a repulsive term as two atoms become very close, the electrostatics term will (if the charges are of opposite sign) provide greater and greater attraction, leading to a singularity. As such, I decouple the sterics and electrostatics interpolation so that there are six control parameters for the transformation of nonbonded terms: $\lambda_{\text{stericsdelete}}$, $\lambda_{\text{stericsinsert}}$, $\lambda_{\text{electrostaticsdelete}}$, $\lambda_{\text{electrostaticsinsert}}$, as well as two

additional parameters, λ_{sterics} , $\lambda_{\text{electrostatics}}$ for core atoms whose identity is changing but are not being inserted or deleted. Having divided up the control parameters as such, I then interpolate as in Figure 5.2.

Softcore sterics

In order to provide greater simulation stability, I also utilize softcore sterics [35].

In this form, the sterics potential becomes:

$$U(r; \lambda_{\text{sterics}}) = 4\epsilon \left[\left(\frac{\sigma}{r} \right)^{12} - \left(\frac{\sigma}{r} \right)^6 \right] \quad (5.13)$$

$$\lambda_{\alpha} = \text{dummy}_A(1 - \lambda_{\text{sterics}}) + \text{dummy}_B(\lambda_{\text{sterics}}) \quad (5.14)$$

$$\epsilon = (1 - \lambda_{\text{sterics}})\epsilon_A + \lambda_{\text{sterics}}\epsilon_B \quad (5.15)$$

$$\sigma = (1 - \lambda_{\text{sterics}})\sigma_A + \lambda_{\text{sterics}}\sigma_B \quad (5.16)$$

$$r_{\text{eff}} = \sigma * \left[\alpha \lambda_{\text{sterics}}(1 - \lambda_{\text{sterics}}) + \left(\frac{r}{\sigma} \right)^6 \right]^{1/6} \quad (5.17)$$

where dummy_A and dummy_B are parameters that determine whether the atom in question is a dummy atom at the initial endpoint or the final endpoint. Though I did not explore this, the formulation includes an additional adjustable softcore parameter, α .

Challenges with existing code bases

The treatment of electrostatics controlled by a global parameter can be difficult, as particle mesh ewald [19] has a nontrivial efficient implementation. In this work, I leverage the ability of OpenMM [24] to linearly interpolate charges while still calculating the full PME energy. Without this, I was forced to implement only the direct space of the electrostatics term. However, this resulted in the overlap

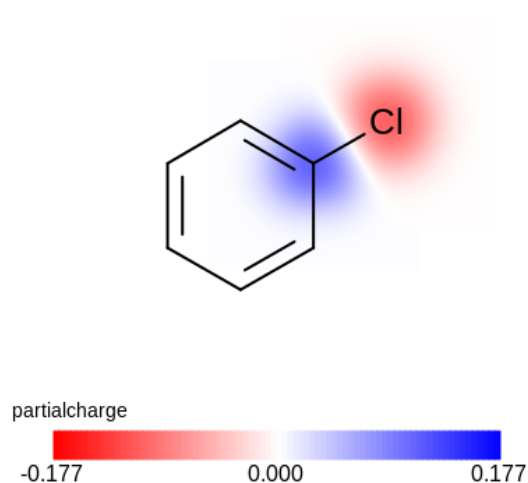


Figure 5.1: Chlorobenzene, with atoms colored by charge. If the ring is mapped to benzene (where all partial charges are equal), the abrupt change is very unfavorable.

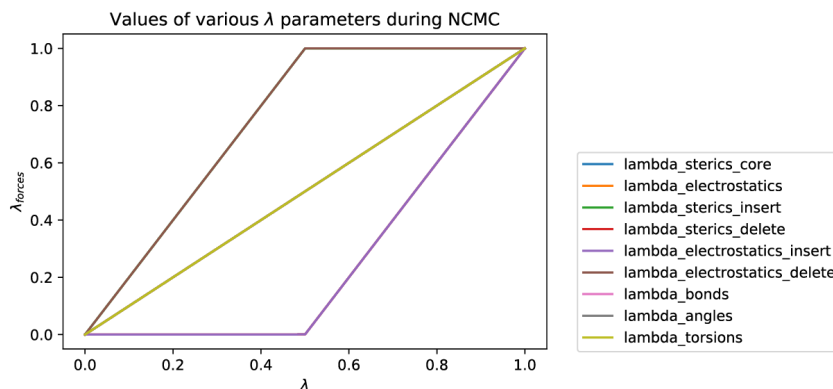


Figure 5.2: Nonbonded interpolation scheme, with the schedule of each individual force's λ parameter plotted against the global λ . Note that sterics are always eliminated after electrostatics, and are always added before electrostatics. This provides protection against unshielded charges.

between the hybrid system at the endpoints and the endpoint systems themselves being very poor, especially for charged ligands. Therefore, it is advantageous to ensure that the full PME can be represented in the intermediate states.

5.3 A note on ring closure

One of the original motivations for pursuing the reversible jump algorithm was the ability to add complex ring systems, a task generally considered quite difficult for relative free energy calculations [44]. However, as discussed in chapter 4, ring closure can be rather difficult, as the sequential importance sampling algorithm gradually produces weights that are increasingly poor. While resampling in the geometry proposal appears attractive, it is very expensive; the need to compute many replicas of energy terms would require a highly optimized code for that purpose, and furthermore may suffer from numerical issues [59]. However, an appealing alternative would consist of simply softening the bonds and angles in the initial state of the hybrid system. This would allow poorer proposals from the dimension matching system, instead cleaning up the poor geometries with annealing (which has the advantage of having tunable hyperparameters and running on the GPU). Before running the entire solvent system, it is useful to examine how the softening of bonds and angles might affect the contribution of the initial jump to the hybrid system, known in the code as `logP_to_hybrid`. In Figure 5.3, one can see that the softer bond parameters significantly reduced the variance of the contribution of the jump to the hybrid system. In addition to bonds, angles are also represented as harmonic terms, and so it behooves one to examine the effect of softening angles as well, under the same schedule. This time, we will hold the bond softening constant at 1.0 (no softening), and

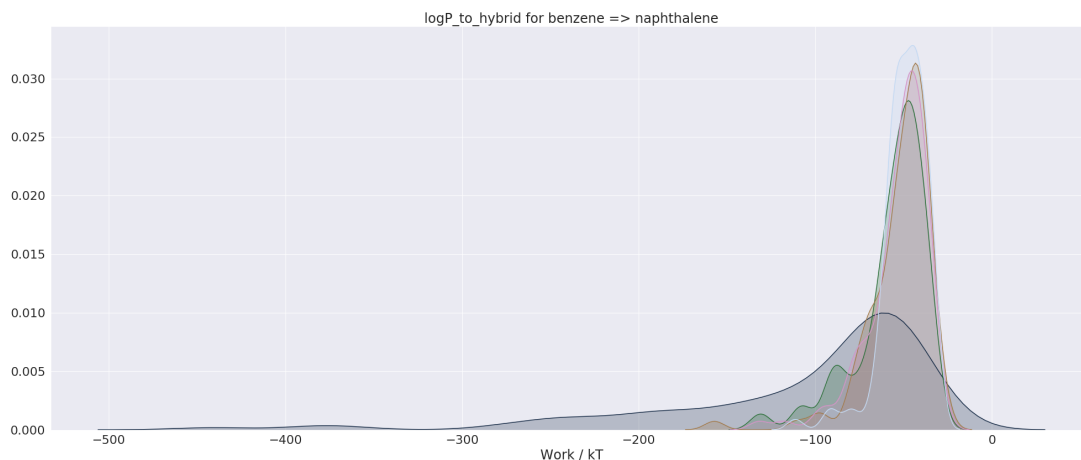


Figure 5.3: The distributions of the `logP_to_hybrid` contribution under different initial bond softening parameters. The initial (darkest and broadest) distribution is with no softening; moving to lighter colors, we see the distribution narrow considerably as the initial force constants for bonds connecting unique atoms is scaled by 0.1, 0.01, 0.001, and 0.0001

examine the effect of softening the angles in 5.4 Interestingly, the angle terms do not seem to have the same profound effect as the bond terms. An obvious explanation for this is that the force constants for the angles are considerably smaller, so it may require less softening. In Figure 5.5, I show the same softening schedule for both parameters. Taking a closer look at the effect on the standard deviation of `logP_to_hybrid`, Figure 5.6 shows the effect of the different softening parameters on the standard deviation of the `logP_to_hybrid` component. Notably, the angles need not be softened as much as the bonds to reduce the variance of this component. The standard deviation of the standard deviations, obtained by bootstrapping, can be found in Figure 5.7 Although we can greatly diminish the variance of this quantity by simply softening the bonds, it is important to keep in mind that the ultimate objective is not to minimize the variance of the jump to the hybrid system, but rather to minimize the variance of the ultimate acceptance

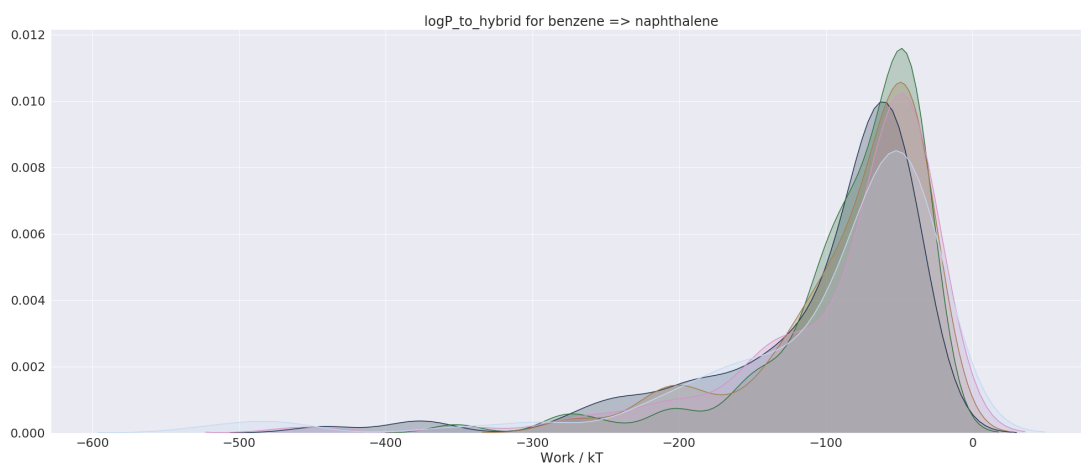


Figure 5.4: The distributions of the `logP_to_hybrid` contribution under different initial angle softening parameters. The initial (darkest and broadest) distribution is with no softening; moving to lighter colors, we do not see the same profound effect as the initial force constants for angles connecting unique atoms is scaled by 0.1, 0.01, 0.001, and 0.0001

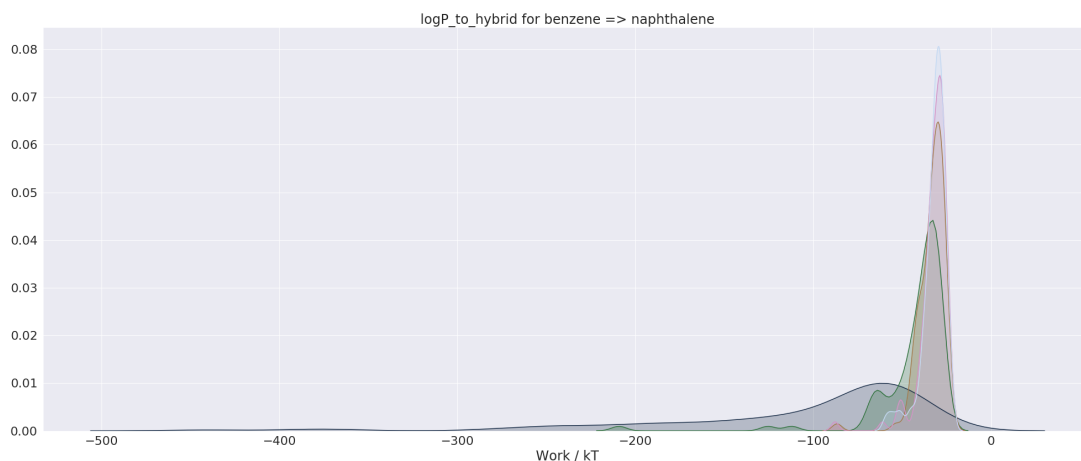


Figure 5.5: The distributions of the `logP_to_hybrid` contribution under different initial angle and bond softening parameters. The initial (darkest and broadest) distribution is with no softening; moving to lighter colors, we see that the distribution grows considerably narrower as both the angle and bond terms are softened to the same degree according to the schedule 0.1, 0.01, 0.001, 0.0001.

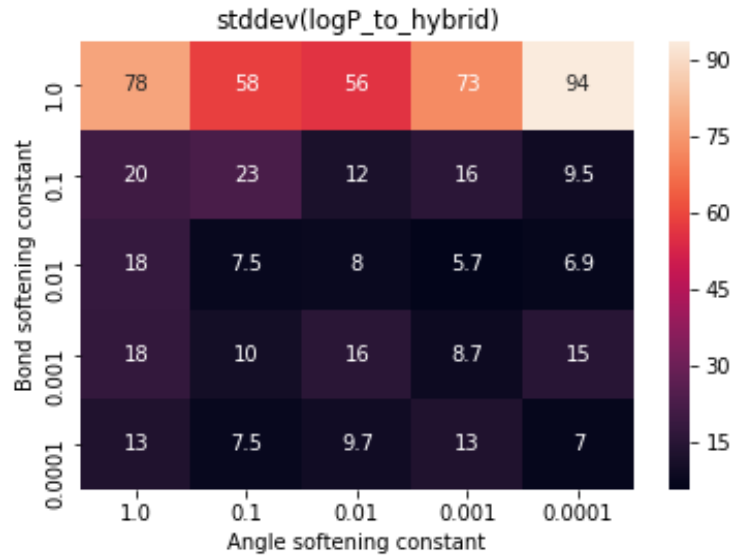


Figure 5.6: The standard deviation of the logP_to_hybrid component under different combinations of bond and angle softening constants. It is noteworthy that the angles need not be softened nearly as much as the bonds. All quantities are in effective units of $k_B T$.

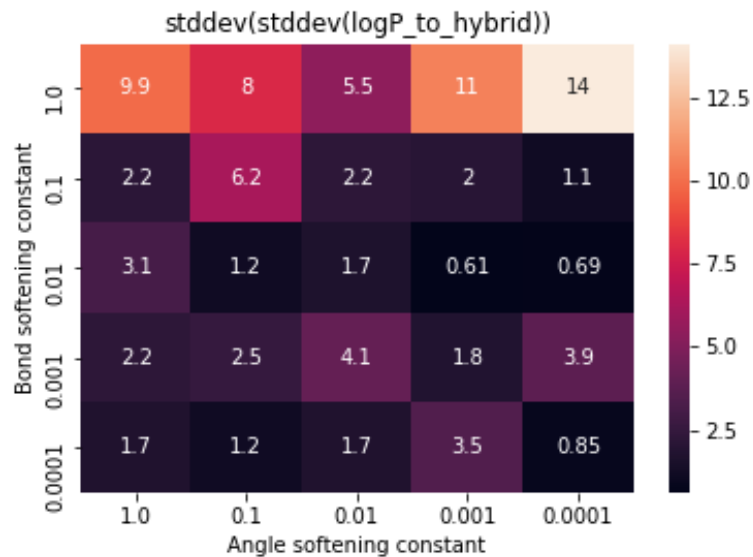


Figure 5.7: The standard deviations of the standard deviations of logP_to_hybrid under different combinations of bond and angle softening terms. All quantities are in effective units of $k_B T$.

probability. It is not difficult to imagine how excessive softening could negatively affect the overall acceptance probability, especially if the dimension matching distribution does not match the softened distribution.

5.4 Tuning NCMC Protocol length

The simplest hyperparameter to tune in NCMC is the length of the switching protocol. The longer this protocol is, the more favorable the work values. At the same time, the longer protocols require a greater amount of wall clock time. At a certain point, the additional wall clock time consumed by the NCMC simulation is not worth the diminishing benefit in terms of work (which contributes directly to the acceptance probability). Certain transformations may require a longer protocol than others; in this work, I use the same protocol length throughout a calculation.

5.4.1 Limitations of tuning only length

However, tuning the length alone, while straightforward, has limitations. For instance, with a naive protocol linearly switching λ , it may be necessary to extend the protocol to many steps to ensure reasonable acceptance probabilities. This approach does not give us the opportunity to more freely alter the schedule of the control parameters and potentially gain efficiency.

5.4.2 Tuning Annealing Schedule

However, with a protocol that may be nonlinear with the control parameter λ , we may afford ourselves an efficient protocol with many fewer steps. The practical upshot of this is that a shorter protocol can achieve reasonable performance, provided the protocol is reasonable. Unfortunately, as described previously, the thermodynamic length can be relatively difficult to calculate, and will vary depending

Theory

In order to perform the experiment above, I first computed the thermodynamic metric tensor along the one dimensional path. Then, I altered the schedule of control parameters to change λ more rapidly when the metric tensor was small (indicating that neighboring distributions are similar) and more slowly when it is large (indicating that neighboring distributions may be quite far apart). Although this sounds like a reasonable task, it is simplified in the case of the harmonic oscillator, where the metric tensor is available analytically. Unfortunately, as is typically the case, the metric tensor is not available analytically for any molecular simulation of interest. Thus, it must be estimated from sampling data, as described in [55]. This can be a rather expensive process, and may not be worth the added computation time.

5.4.3 Limitations to tuning annealing schedule

Despite the pleasant straightforwardness of using the metric tensor to define nonequilibrium switching protocols, there is another pitfall besides the computation cost: different pairs of ligands will have different metric tensors for their transformations, as well the same pair of ligands in different environments. This means that it is not straightforward to transfer the values learned from one calculation into another. Coupled with the computational cost, this may make the approach of using the thermodynamic metric tensor infeasible. However, there may be alternate schemes that allow some degree of transfer learning.

Simple solutions

The simplest solution would be to perform a calculation with a number of different ligands in a common environment (such as explicit water) and trying to estimate a metric tensor for each of those transformations. Then, one could potentially utilize a summary statistic of this collection of metric tensors in the hopes that on average, it will perform better than a naive protocol. These approaches need to be explored in greater detail with a large number of different conditions

General Solutions using Machine Learning

In the future, one promising path forward is to use advances in machine learning to predict favorable alchemical paths. In this paradigm, one might take many calculations in different environments, and learn a function $f_{\theta}(\mathcal{M}_1, \mathcal{M}_2)$ with optimizable parameters θ that maps a pair of molecules \mathcal{M}_1 and \mathcal{M}_2 to a

reasonable schedule of control parameters. Although appealing, it remains to be discovered how such a scheme would account for the myriad environments in which molecules find themselves in the course of a drug discovery program.

5.4.4 Relationship to Geometry Tuning

As mentioned earlier, there exists to some extent a tradeoff between the effort engaged in dimension matching and the effort spent in nonequilibrium switching. More precisely, the more similar the hybrid system at the endpoint is to the non-modified system at the corresponding endpoints, the more straightforward the task of the dimension matching algorithm. However, completely annihilating the valence interactions of newly-added atoms is also potentially dangerous. Atoms that are not anchored to the molecule can easily drift far away, creating additional work as the bonds and angles are scaled on. On the other hand, the dimension matching algorithm can essentially ignore the nonbonded interactions with surrounding atoms, as these are always decoupled at the beginning of the nonequilibrium switching protocol. One question that remains is whether it is worthwhile to include the nonbonded interactions for the sake of maintaining favorable intramolecular nonbonded interactions.

Possibility of resolving poor valence terms

Although initially it seems that softening bonds is deleterious to the acceptance probability, the use of the dimension matching scheme without resampling means that small deviations from the target distribution can easily compound for large proposals. As a result, it seems prudent to offer the choice of softening bonds, angles, and torsions between newly introduced atoms.

CHAPTER 6

STOCHASTIC APPROXIMATION

Having completed the components for preserving the correct equilibrium distribution, I now turn to the task of developing the on-line adaptation scheme for biasing toward states with desirable free energy properties. In order to do this, I turn to stochastic approximation, which has been used successfully in many contexts, especially Monte Carlo [70, 84, 42, 64].

6.1 Stochastic Approximation for Molecular Simulation

Before deriving the update rule for the present case, I will briefly review applications of stochastic approximation in molecular simulation. Often, when one is simulating from an expanded ensemble, the mixture components (individual discrete states) will have very different free energies. Since the relative population is given by

$$\frac{p(k = n)}{p(k = m)} = e^{-\beta(\Delta G_n - \Delta G_m)} \quad (6.1)$$

even small differences in free energies will result in very large differences in relative populations. Practically speaking, this results in only rarely calculating acceptance probabilities between certain states, and therefore poor statistics in the final estimate of relative free energies. To resolve this, one approach is to attempt to flatten the histogram of states. In order to achieve this flattening, one could use the expanded ensemble weights g_k add additional weight to those states with a lower free energy. This requires knowing the relative free energies, however, defeating the purpose of running the algorithm in the first place. Therefore,

many adaptive schemes have been developed to estimate these weights online, such as in [88] [84] [42].

6.1.1 SAMS

Recently, a formulation of stochastic approximation Monte Carlo was developed by Tan [84] that is asymptotically optimal in the sense of minimizing the variance of the weights g_k . This algorithm is performed as follows:

$$g_k^{(t-1/2)} = g_k^{(t-1)} - t^{-1} \frac{\delta_{s_t,k}}{\pi_k} \quad (6.2)$$

$$g_k^{(t)} = g_k^{(t-1/2)} - g_1^{(t-1/2)} \quad (6.3)$$

where each $g_k^{(t)}$ is the adapted weight of chemical state k at iteration t , s_t is the current state at iteration t , and π_k is the target probability for state k . Setting each $\pi_k = \frac{1}{n}$ for a system with n states will result in even sampling, and asymptotically will result in weights equal to the relative free energies [84]. Its simplicity, combined with its asymptotic optimality [84], makes it an appealing choice for adaptively reweighting the components of the expanded ensemble. Although this can be very useful for determining the entire set of pairwise relative free energies, it can still be very burdensome for large sets of chemical states. In general, we are not interested in merely even sampling, as we are not concerned with the relative free energies of two unfavorable chemical states. Therefore, we seek an adaptive algorithm that provides for greater sampling of chemical states that are favorable according to some free energy criterion.

6.2 Doubly-recursive SAMS

In general, in physical simulation tasks, we seek states that optimize a free energy difference between multiple states, not simply a single state. For instance, if we are optimizing to maximize binding affinity, the quantity we seek to maximize is (omitting constants) the association constant K_a :

$$K_a \propto \frac{Z_{PL}}{Z_P Z_L} \quad (6.4)$$

where Z_* represents the normalizing constant of the appropriate species $*$, P and L represent protein in solvent and ligand in solvent systems respectively, and PL represents the interacting protein-ligand complex system. Since we are trying to find the chemical state which maximizes this quantity, we are generally interested in so-called relative free energies. As a result, when comparing two ligands, L_1 and L_2 , we have:

$$K_{a,1} = \frac{Z_{PL_1}}{Z_P Z_{L_1}} \quad (6.5)$$

$$K_{a,2} = \frac{Z_{PL_2}}{Z_P Z_{L_2}} \quad (6.6)$$

$$K_{relative} = \frac{K_{a,1}}{K_{a,2}} = \frac{\frac{Z_{PL_1}}{Z_P Z_{L_1}}}{\frac{Z_{PL_2}}{Z_P Z_{L_2}}} = \frac{\frac{Z_{PL_1}}{Z_{L_1}}}{\frac{Z_{PL_2}}{Z_{L_2}}} = \frac{Z_{PL_1}}{Z_{PL_2}} \frac{Z_{L_2}}{Z_{L_1}} \quad (6.7)$$

Therefore, if we would like a chemical state to be sampled according to its relative binding affinity, we need to adapt the weights of the chemical states such that

$$p(k) \propto \frac{Z_{PL_k}}{Z_{L_k}} \quad (6.8)$$

That is, the target weights π must be adapted as well, since we do not know *a priori* what the relative free energies are.

6.2.1 Description of Algorithm

More generally, consider we have s different probability densities:

$$p_{ij}(x) = e^{g_{ij}^*} q_{ij}(x) , \quad i = 1, \dots, s , \quad j = 1, \dots, m \quad (6.9)$$

and we desire to design a chain where the marginal distributions of all s chains are

$$p_{ij} \propto \prod_{i'=1}^s e^{-\theta_s g_{i'j}^*} = \exp \left[- \sum_{i'=1}^s \theta_s g_{i'j}^* \right] \forall i = 1, \dots, s \quad (6.10)$$

where the *design vector* $\Theta \equiv \{\theta_1, \dots, \theta_s\}$ specifies how different targets and antitargets are used in weighting the design constraints. We postulate that we can do this by defining $\pi_i(Z)$ for $Z \equiv \{g_1, \dots, g_s\}$ as

$$p_{ij}(Z|\Theta) \propto \exp \left[- \sum_{i'=1}^s \theta_s g_{i'j} \right] \quad (6.11)$$

6.3 Toy Examples

Having described the algorithm in principle, I now present a simple toy example to empirically see that the quantity of interest is generated. In this example, we simulate a set of harmonic oscillators:

$$p_{i0} \propto e^{-\beta \frac{K_{i0}}{2} (x - x_{i0})^2} \quad (6.12)$$

$$p_{i1} \propto e^{-\beta \frac{K_{i1}}{2} (x - x_{i1})^2} \quad (6.13)$$

As a model of the binding free energy calculation, we would like:

$$p_1(k) \propto \frac{Z_{1k}}{Z_{2k}} \quad (6.14)$$

where the subscripts 1 and 2 indicate the different chains (analogous to complex and solvent for the binding free energy). As can be seen in Figure 6.1, the target

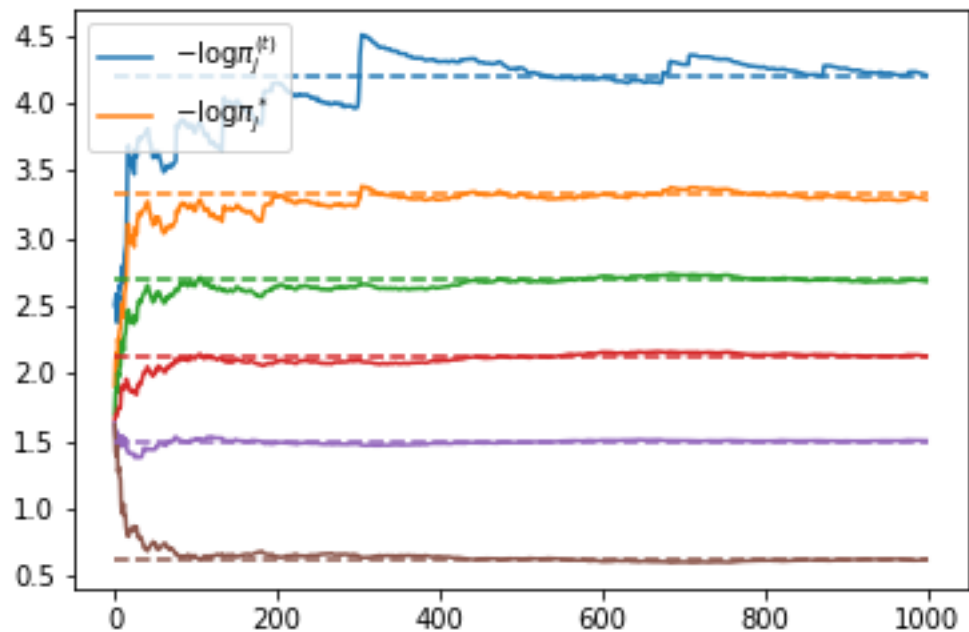


Figure 6.1: The target weights of the various harmonic oscillators (solid) and the true value (dotted line). Note that the higher free energy states take significantly longer to converge to their true values

weights of the favorable states (lower free energy) converge faster. This is what we desire, since we are not interested in how much more unfavorable one high free energy state is from another. We can extend this principle to multiple chains as well; as generalized above, this applies to cases such as multitargeting (finding chemical states with a favorable free energy in multiple chains) or selectivity (finding states that are favorable for one but not another).

6.3.1 Performance

Ultimately, the goal of the algorithm is to more quickly explore relevant regions of chemical space. However, there are some limitations to the stochastic approximation algorithm presented here that are worth discussion.

6.4 Limitations

Two issues arise when discussing the limitations of the stochastic approximation algorithm used here. The first is the convergence rate. Since the algorithm works by counting the number of times different states are visited, it can take a lengthy simulation before the chemical state space begins to be explored. The second limitation is that the SA algorithm presented here is designed for a fixed number of states. However, it is clear that to truly carry out molecular design, an algorithm must be capable of handling the situation where the number of states is not known *a priori*. This can arise, for instance, when the proposal distribution is a neural network.

6.4.1 Convergence Rate

The first topic of concern is the convergence rate. To examine this in greater detail, we observe behavior in a toy model. Even in the traditional binary SAMS update (described in [84]), it can take many iterations to begin to achieve even sampling, especially if the mixture components are kinetically separated. In Figure 6.2, one can see that the desired target weights (dotted lines) are quickly

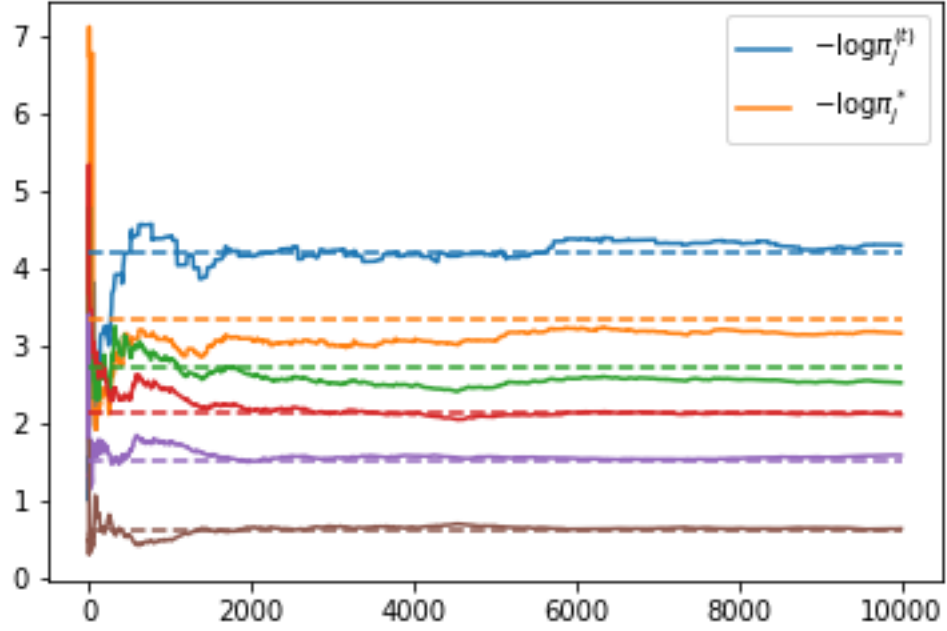


Figure 6.2: Convergence of binary SAMS from a single realization of harmonic oscillators with only a small separation in means.

approached by the adapted target weights (solid lines), and are more quickly approached for low free energy states than others. However, what happens if we add a small amount of kinetic separation to these harmonic oscillators? In principle, this issue is resolved by nonequilibrium switching; however, in practice, efficiently resolving these types of issues is quite difficult, especially when the states between which one must transition are so diverse. More work needs to be performed to not only initialize the weights well, but investigate the possibility of initial algorithms that may approach the desired target more quickly than the asymptotic algorithm.

6.4.2 Number of States

Ideally, we'd be able to jump to new (unseen) chemical states, which would enable molecular design in a truer sense. However, the formulation of the SA algorithm above clearly is designed for a fixed number of states. One way to resolve this is to imagine that, although more states may be added on the fly, we set an upper limit on how many states can be explored. In this case, we will eventually reach the maximum number of states, at which point the algorithm obviously becomes the double recursion SA above. However, there are some questions regarding this. We already saw that the convergence of the algorithm is highly sensitive to the mixing of the underlying sampler; what happens when the number of states changes? This is a subject that calls for further investigation and future work.

6.5 Weight initialization

One major challenge with SAMS algorithms is how the weights are initialized. The initial weights will have a profound effect on the performance, especially as weight adaptation diminishes. For this reason, many use a multiple-stage scheme as in [84]; this allows the weights to initially adapt much more quickly, and then revert to the asymptotically optimal gain decay after a certain condition is reached. In addition to this scheme, it is advantageous to determine whether there exists simple tricks that can initialize the SAMS weights nearer the correct target than zero.

6.5.1 Implicit volume term

One issue that the reversible jump-based simulation incurs that most others do not is that the dimensionality of phase space changes with the chemical state. As a result, depending on the units used, there is a volume term that is not accounted for in the acceptance probability. Since in free energy calculations we are always looking at differences, this is not a concern for accuracy, but is for efficiency, since this term will cause certain states to become very unfavorable based on the number of degrees of freedom in that state. Observing this discrepancy, I decided to initialize the stochastic approximation weights for each state using a very rough guess of the implicit volume term:

$$g_k^0 = n_{heavy} * 4.5 + n_{hydrogen} * 3.8 \quad (6.15)$$

The hydrogen term is smaller because it only has two degrees of freedom (bonds are constrained).

6.5.2 Hydration

A common environment in relative free energy calculations is the solvent phase, which consists of just the small molecules in solvent. Since the environment is essentially the same in solvent regardless of the target, it is advantageous to develop a quick initialization scheme that can be easily reused. Here, I implemented one scheme, however, there exists the potential for other schemes to be developed as well.

Initialize with minimized point energies

One interesting approach for initializing weights in the solvent phase is to use the minimized point energies of the corresponding molecule in implicit solvent. In [56] this approach was explored as an approximation to the hydration free energy of the molecule. Since we do not rely on the initialization of weights for correctness, only for efficiency, a simple scheme like this one is useful. It is possible that more sophisticated implicit solvent models (such as PBSA [83]) would produce superior starting weights, but the added cost may not be justified.

6.5.3 Complex weight initialization

In molecular design problems relating to drug discovery, another common phase is the complex phase—that is, protein and ligand interacting in solvent. This phase usually contains systems that have many more atoms, as well as much longer correlation times. Therefore, it would be extremely helpful to have an efficient weight initialization scheme. However, the complex phase poses several additional challenges. One challenge is that the relative free energies of different chemical states in the complex phase cannot be as easily estimated as in the solvent phase. Another is that the complex phase varies greatly between different projects; one cannot come up with a single, simple initialization scheme for all complex phase calculations as one can for the solvent phase. However, there are several ideas that may be implemented.

Initial acceptance probability

One simple idea is to use the initial $\ln P_{\text{accept}}$ values that are generated by the algorithm. This is appealing intuitively, because P_{accept} is a (very poor) estimate of the relative free energy of the two states [97], as well as being produced as a byproduct of the algorithm (requiring no additional effort). However, there are several pitfalls to this approach. In particular, due to the potentially high variance of the P_{accept} , setting a stochastic approximation weight with it can result in an initialization extremely far from the true value. This would actually make the simulation even less efficient. One approach would be to only initialize the weight with the initial acceptance probability if the log acceptance probability is within some bounds. This might prevent the SAMS weights from being initialized to quantities very far from the true value, while still allowing a decent initial guess.

Machine learning score

Another approach not explored here would involve using a machine learning algorithm that attempts to predict binding free energies [14, 94] to initialize the weights. A caveat here would be that a machine learning algorithm that is more accurate or otherwise has a different error pattern than the forcefield might start the weights farther from the desired true values (even if they are more accurate in reality). The area of weight initialization represents a fertile ground for future investigation.

Implicit ligand theory

Another interesting area of research is implicit ligand theory [61, 95, 54, 53], which uses one or more snapshots of a rigid protein and simulates only the small molecule, resulting in an impressive speedup. These rapidly-calculated approximations could be used to initialize the weights, providing a potentially profitable starting point.

CHAPTER 7

ALTERNATIVE TO SA: HIGHLY PARALLEL NONEQUILIBRIUM SWITCHING

7.1 Introduction

As discussed, there are several limitations of the stochastic approximation (SA) algorithm. However, that does not preclude the use of the reversible-jump (RJ) algorithm in other clever ways. For a long time, it has been known that rather than accepting or rejecting a proposal attempt, one can cache the resulting acceptance probability and use it to estimate the free energy difference offline. In molecular simulation, this technique is known as nonequilibrium switching [36, 1], and in statistics as annealed importance sampling (AIS) [60]. Examples of estimation techniques include the exponential averaging estimator [97], as well as the optimal Bennett Acceptance Ratio (BAR) [7, 76, 18, 51], which uses bidirectional data. Though this has not often been utilized for free energy calculations, it has several appealing advantages in the modern age. First, it is highly parallelizable—unlike SAMS-like expanded ensemble algorithm, it can easily run many switching trajectories or proposal attempts simultaneously. With hardware becoming increasingly parallel and the availability of cloud computing, this is an attractive feature. Second, there is no need to wait for a sampler to explore all the states—by definition of the algorithm, the user can choose which states. Third, it is still possible to employ a wide range of adaptive algorithms. For instance, one could use Bayesian experimental design [87, 20] to iteratively choose how to allocate effort to different “legs” of the calculation.

7.2 Advantage of RJ Nonequilibrium Switching

Since codes for nonequilibrium switching based relative free energy calculations already exist, one might wonder what advantage the complication of reversible jump brings. In order to understand why one would undertake this approach, consider what is required for a standard nonequilibrium switching based relative free energy calculation. In Figure 7.1, one can see that the typical setup of a nonequilibrium free energy calculation is to run simulations at the alchemical endpoints (that is, a hybrid system with the control parameters set to either 1.0 or 0.0), and periodically generate proposal attempts to reach the other endpoint. This requires simulation of the hybrid system at both endpoints for each pair for which one would like to calculate the free energy difference. Note that in this scheme, one cannot use the equilibrium simulation used in a calculation for $\Delta\Delta_{AB}$ to, for instance, estimate $\Delta\Delta_{AC}$, where A, B, C are molecule indices. This is because the alchemical system must contain all the degrees of freedom of both endpoints, which will not necessarily include the degrees of freedom of the rest. However, using the reversible jump scheme, one simulates from the *nonalchemical* endpoints for each molecule, and then in the course of the nonequilibrium switching trajectory, inserts the missing degrees of freedom and removes the superfluous degrees of freedom. This means that, to use the previous abstract example, one can simply run multiple equilibrium simulations, then adaptively choose which pairs should be emphasized. It would also allow new molecules to be added without re-simulating any of the existing molecules in the set; in principle, if one wanted a fast estimate, the exponential averaging estimator could even be used to without performing any additional equilibrium simulation. Additionally, the reversible jump scheme enables ring formation

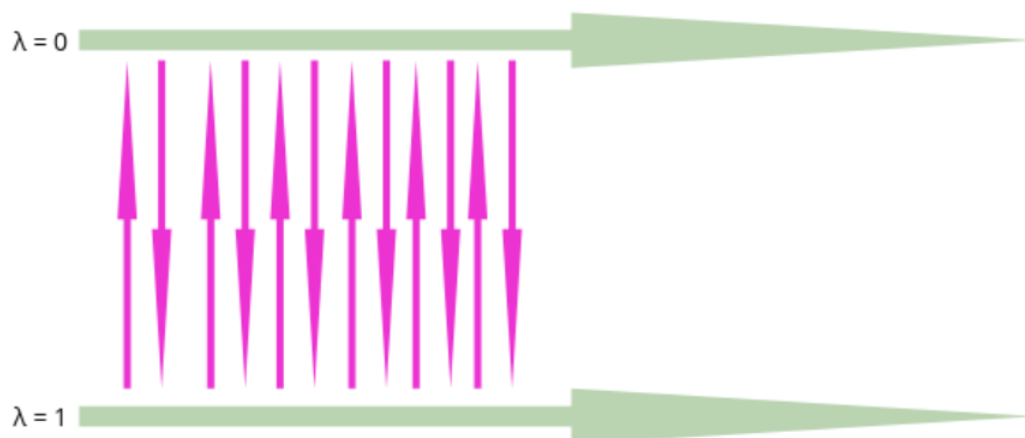


Figure 7.1: Illustration of a nonequilibrium switching free energy calculation. Note the requirement for simulation at the alchemical endpoints. The horizontal green lines represent the equilibrium simulations at each endpoint, while the magenta lines depict the switching trajectories from samples of each state to the other.

and breakage without resorting to more sophisticated tricks, which allows a much greater freedom for the user to choose different Finally, an adaptive design scheme based on Bayesian experimental design could easily incorporate other non-free-energy-based objectives, such as a model of synthetic accessibility or some other desirability criterion.

7.3 Use of cloud computing resources

This approach to the reversible-jump based algorithm presented in this thesis opens the door to ready use of cloud resources such as Amazon Web Services, which are already being employed in free energy simulations [16]. Offerings such as AWS not only enable users without in-house clusters to use high-performance

hardware, but also enhance the reproducibility of research. Entire stacks used to perform computation can be exported via services such as AWS CloudFormation, enabling rapid and straightforward reproduction of research. Another benefit of using on-demand cloud resources is that it enables the direct application of economics to the question of calculation efficiency. Rather than comparing algorithms in terms of asymptotic variance or other similar criteria, we can now compare them in terms of their dollar cost; how much do we have to spend in order to achieve the desirable result? This also allows us to put a direct price tag on speed improvements. While questions such as the variance of the acceptance probability once seemed abstract, they now strike our wallets. Since we are performing calculations on groups of molecules, we can also use this cost data to determine which edges are most profitable to explore, and which can be estimated by summing other edges. As in [89], the presence of cycles could also be used to correct for errors in the calculation.

7.4 Pricing and Economics of Cloud Computing

Due to familiarity, the remainder of this chapter will use Amazon Web Services as an example. Using AWS for molecular simulation, there are several cost considerations:

- Storage of input and output
- Compute time
- Data egress (ingress is free)

Here, I will primarily discuss the issue of compute time. There are several pricing models on AWS and its competitors for compute time. On-demand pricing refers to the use of resource at a fixed hourly price. So-called spot pricing allows users to bid on unused capacity at a significantly reduced price. This price fluctuates, however, and is not guaranteed to remain under the user's bid. In the present situation, this is tolerable: if an instance is killed because it becomes prohibitively expensive, we lose whatever nonequilibrium switching it has performed but not copied to storage, but nothing else. At the time of this writing and in the geographically nearest datacenter, the price of a single GPU instance hovered near \$0.27 per hour. In preliminary investigations on a GTX-Titan, a full proposal of 10ps from one alkane to another takes approximately 60 seconds. In other words, we can compute over 100 switching attempts for just 50 cents! Of course, this is a small system—scaling up will be pricier. These price signals can also guide the investigator into the most effective use of his or her time and money. However, this illustrates the low cost of the algorithm on modern infrastructure, and also highlights that pairs of molecules for which is advantageous to compute many switching attempts can be done cheaply and in parallel. One may be left wondering why one couldn't simply use many replicates of the stochastic approximation algorithm to take advantage of the parallelism. One could do this, of course, and whether this is efficient in terms of processor or wall clock time requires empirical investigation. However, one feature of the low spot pricing is its volatility—a user cannot rely on the price remaining at its current level indefinitely. In fact, this aspect has motivated interesting research in modeling the price of spot resources [37] that can be brought to bear on this application. Since the nonequilibrium switching trajectories are independent of each other, loss of some compute resources is not devastating. If necessary, the number of

nonequilibrium switching trajectories being performed in parallel can be reduced to save money when the cost increases. Effort can be adaptively reallocated by whatever model or heuristic the user desires. However, if one suspends a chain exploring chemical space due to cost, one loses the adaptation that that chain has already performed. There may exist powerful methods of coupling the stochastic approximation weights, but this is a topic for future research.

7.4.1 Factors affecting performance

One factor affecting performance is the nonequilibrium protocol itself. Another factor that could potentially affect performance of this algorithm is the correlation time of the equilibrium simulation. Unlike the SAMS case, where the algorithm accepts or rejects moves to different chemical states, this approach does not benefit from the potential decorrelating effect of changing chemical states. Additionally, the number of parallel computing devices available will profoundly affect whether this choice is a feasible one. With a large number of processing devices, this approach can reduce wall clock time, but empirical evaluation is necessary to determine whether for relative free energy calculations this approach is comparable in efficiency to equilibrium staging.

CHAPTER 8

HYDRATION FREE ENERGY

8.1 Introduction

In this example, we seek to compute the relative hydration free energies of different ligands in explicit solvent. We will use the approach of transdimensional nonequilibrium switching to compute relative hydration free energies for a set of molecules in the FreeSolv dataset [57, 49], a set of experimentally and computationally characterized hydration free energies for small molecules. Although this is a simple test, it demonstrates that the algorithm is capable of efficiently producing low-variance work values in explicit solvent, and among diverse sets of ligands. For the example presented in this chapter, we will be using a randomly selected subset of substituted benzene molecules. For each pair of ligands, we will run 1 nanosecond of equilibrium simulation and a nonequilibrium switching protocol 5 picoseconds long. We will use 1000 switching trajectories in the forward and reverse direction for each pair, and then use the work values to estimate the relative free energies.

8.2 Thermodynamic Cycle

In order to compute the relative hydration free energy, we employ a thermodynamic cycle wherein we transform the small molecule from one molecule to another in solvent and in vacuum, and subtract the resulting free energies, depicted in 8.1. We compare to FreeSolv [57], which is a database of computed

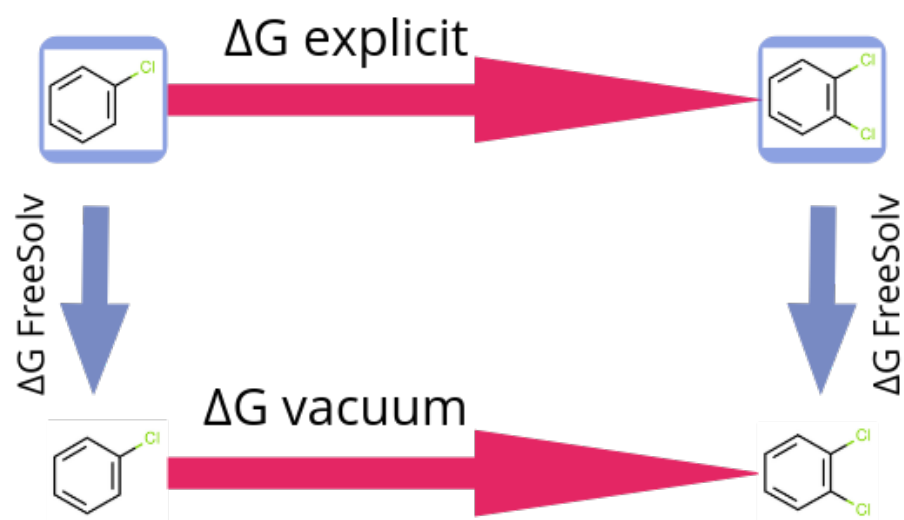


Figure 8.1: The thermodynamic cycle of the hydration free energy calculation. Blue arrows denote the data available in FreeSolv; magenta arrows denote the legs performed in this work.

absolute hydration free energies; therefore, in order to compare our results to FreeSolv, we will need to compute relative hydration free energies from the absolute data.

8.3 Vacuum Phase

For the vacuum phase of the simulation, we performed 1 ns of equilibrium simulation for each molecule, and 100 switches between each pair in both directions. For the nonequilibrium switching length, we used 5ps, as it was sufficient for explicit solvent, and so should be sufficient for vacuum.

8.3.1 Explicit Phase

8.3.2 Nonequilibrium switching protocol length

Prior to performing this calculation, we performed a quick check of the performance of various nonequilibrium switching times vs. the amount of work that was performed. In this check, we simulated benzene and naphthalene in solvent for 1 nanosecond, and then performed 100 switching trajectories in each direction. We collected the data and analyzed the standard deviation of the work performed by the protocol, shown in Figure 8.2 Given that the transformation from benzene to naphthalene is rather difficult without extra bond and angle softening, we decided to use a switching time of 5 picoseconds, and 1000 switching attempts between each pair in each direction.

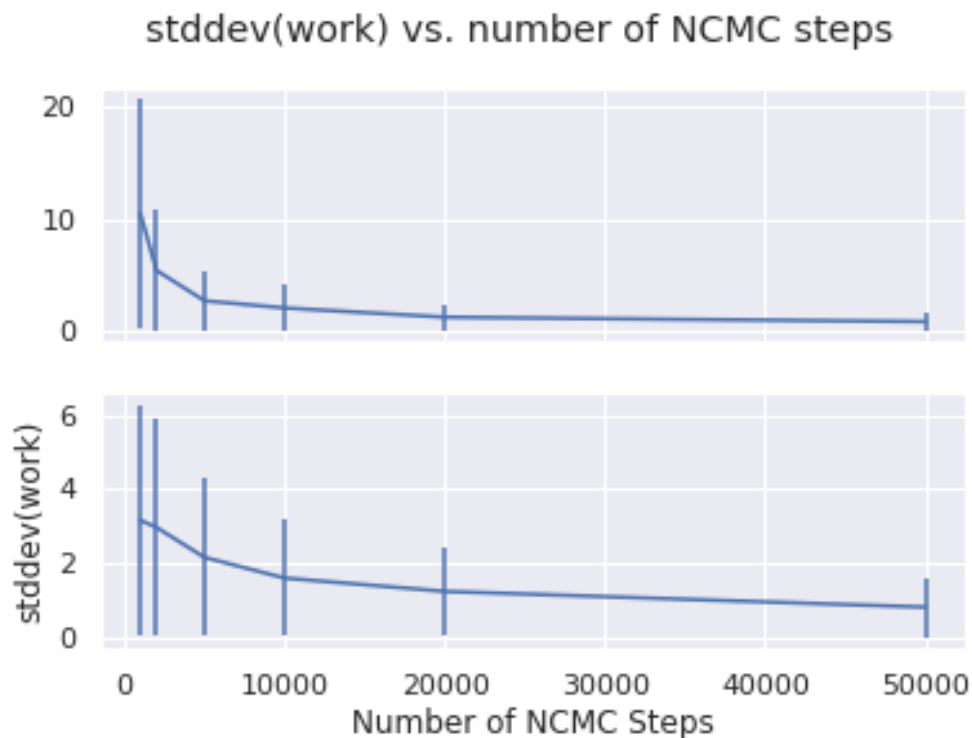


Figure 8.2: The standard deviation of the work for an NCMC switching trajectory between benzene and naphthalene vs. number of steps, with a 1fs timestep. Top: work for transforming benzene into naphthalene. Bottom: work for transforming naphthalene into benzene.

8.4 Analysis of the work of the entire move

In addition to the protocol work, we can also examine the distributions of the work of the entire move in both directions. Since this is the quantity that is fed into BAR, it is critical to examine. In 8.3, there is a representative pair of work distributions. Once again, there appears to be sufficient overlap for proceeding with the estimation of free energy differences. Surprisingly, even in 8.4, the overlap is quite good.

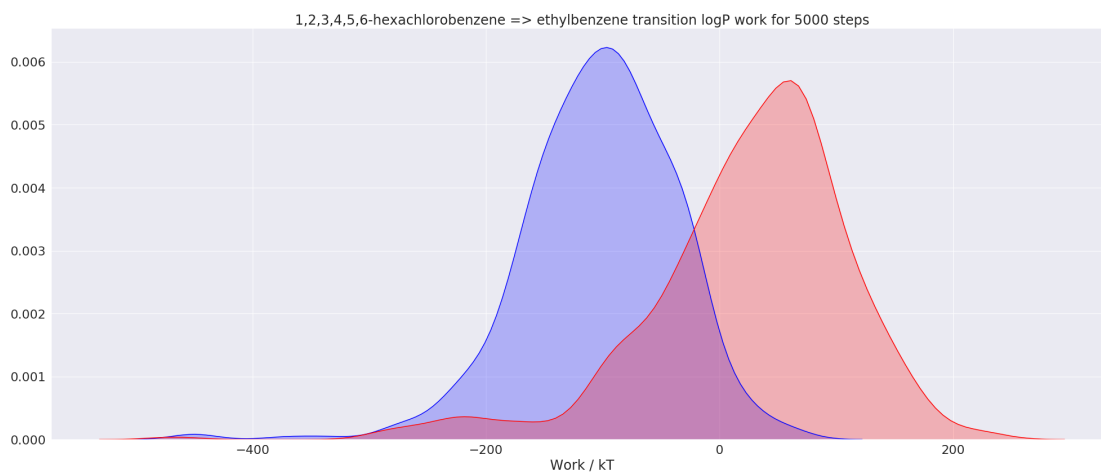


Figure 8.3: The forward and reverse (blue and red, respectively) log acceptance probability distributions for transitions between hexachlorobenzene and ethylbenzene. Although the variances of the distributions appear high, the overlap is very good, allowing a low-error free energy estimate.

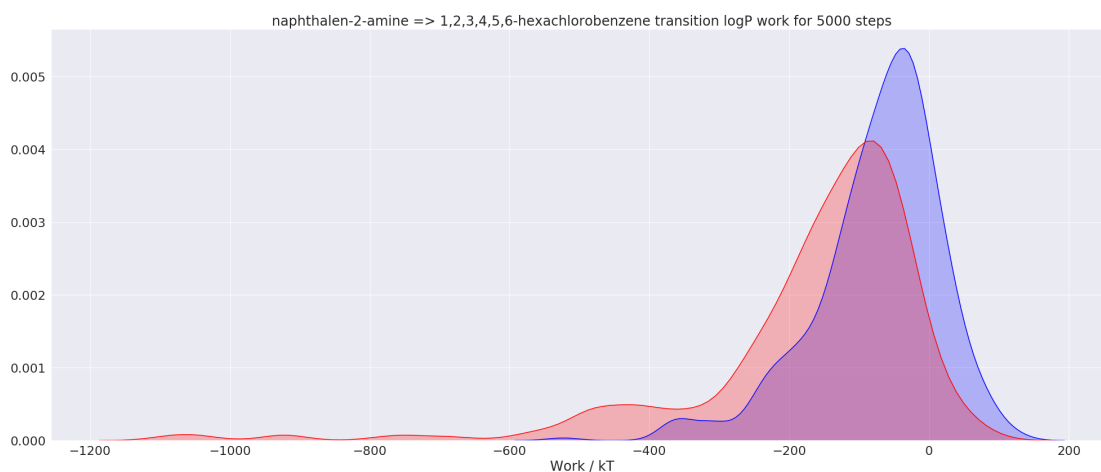


Figure 8.4: The forward and reverse (blue and red, respectively) log acceptance probability distributions for transformations between naphthalene-2-amine and hexachlorobenzene. Despite having to create or break a fused ring, the overlap is still quite good, allowing low variance free energy estimates

CHAPTER 9

CONCLUSION

9.1 Contributions of work

This work makes several novel contributions to the field. First, this work develops a theoretically rigorous algorithm for simulations that can explore multiple chemical states in a single simulation. This relieves the restriction that most free energy calculation methods have, which is that they only compute relative free energy differences between one pair of ligands, or, in certain cases, a small set. Second, this work provides a theoretically rigorous approach to the formation of ring formation and breaking. Finally, the algorithm presented herein allows for adaptive approaches that have heretofore been impossible.

9.2 Exploration of Chemical Space

The first goal that this work accomplished was the introduction of a formal and rigorous basis for approaches such as Chemical Monte Carlo (CMC) [67]. While the concept of CMC is attractive, there was no basis to guarantee that the sampling algorithm would be asymptotically correct. By resorting to Reversible Jump MCMC [30], I am able to provide a basis for jumps in chemical space that provably preserves the correct invariant distribution, allowing the application of additional techniques such as stochastic approximation.

9.2.1 Chemical State Proposals

The first part of the work addresses the requirements as well as the possibilities for making proposals to jump to different chemical states. Unlike in previous algorithms such as Multisite Lambda Dynamics [39], the set of ligands that one chooses for a calculation can be quite diverse. As such, a mechanism of efficiently proposing jumps from one state to another must be developed. In addition to setting forth the conditions for correctness, one must also develop a method that is reasonably efficient. Although formal statements exist to define which chemical states would be most "near" each other in a sense, known as thermodynamic length [17], these expressions are intractable to compute. Therefore, one must resort to heuristics. As a first attempt, I used the number of atoms deemed to be in common between the current and proposed molecule by a maximum common substructure search (MCSS), itself an algorithm with several hyperparameters. I conducted an exploration of various parameters for the MCS search, determining that for some pairs this is a highly influential setting, while for others it matters far less. Finally, I presented groundwork for future work in the realm of chemical state proposals, especially techniques that would allow for a chemical space that is not predefined in the beginning of the calculation. I also described aspects of the algorithm that can be combined with the power of machine learning to leverage the best of both techniques.

9.2.2 Dimension Matching

Once it is decided to attempt a jump to a new chemical state, with atoms assigned to one molecule, the other, or both, the task becomes to insert the missing atoms

and delete the ones that do not belong to the following state. This procedure, known as dimension matching, is a critical component of the algorithm. First of all, there are several important requirements that must be imposed. The proposal must be accompanied by a normalized proposal probability to compute the requisite acceptance test. Second, the proposal must be exact—that is, it cannot be drawn by an approximate scheme (such as most MCMC algorithms). Third, the probability of proposing atoms that are being deleted must also be computable exactly. These restrictions combine to eliminate several straightforward approaches to inserting missing atoms. In this work, I developed a scheme based on Configurational Bias Monte Carlo (CBMC) [79] that meets all of these criteria. This scheme inserts atoms one at a time. First, it proposes a valid order of proposal stochastically. Then, for each atom to be proposed, a dihedral angle with 3 position-bearing atoms is chosen stochastically to establish a reference frame. The proposal then proceeds by drawing a bond length r and angle θ from normal distributions, as both are represented harmonically in common forcefields such as AMBER [68]. The dihedral distribution, however, is not integrable in closed form. To remedy this, the algorithm drives the atom around its dihedral angle, computing a potential energy at each point. This dihedral scan is normalized and used as a proposal probability for the dihedral angle ϕ . This scheme is computationally efficient and meets the theoretical requirements imposed by the proof of correctness of the entire algorithm. An interesting challenge with the CBMC method is to allow the closure of rings—something very difficult for other approaches. In this work, I employed a method inspired by [93]. In this method, I add additional guide force terms so that an atom is almost always placed in a position that will ultimately close the ring. This adds a significant degree of power to this algorithm, as it is now capable of very complex transfor-

mations that would otherwise have been prohibitively expensive or impossible. By setting forth these requirements, I have also left open considerable ground for future development. One interesting avenue to follow would be to use one of a family of algorithms known as particle filtering [3, 15]. This family of algorithms constructs an ensemble of proposals, and at various points during the proposal sequence (in this case, a sequence of atoms), allows the proposals to be resampled from the current set according to a weight. Practically speaking, this approach permits the algorithm to discard "bad choices" that did not become apparent until later in the proposal due to long-range dependencies. This can especially help with challenges such as ring closure, where an apparently reasonable atom proposal can result in a disastrous geometry once a ring cannot be closed. In order to implement a scheme such as this, care will need to be taken to compute the weights at each step in a computationally-efficient manner. This is likely to be fertile ground for future research. Finally, another fascinating potential improvement of this algorithm is to simply resort to a deep learning-based conformation generator. While some of these conformation generators do not meet the above requirements, there have been fascinating recent advances [26] that could enable very rapid and accurate conformer generation in compliance with this algorithm's requirements. It should be noted that such an algorithm used in this case should be trained on *simulated* data, not experimental or quantum mechanical data, as the other sources would not appropriately match the forcefield.

9.2.3 NCMC Switching

Once the new atoms can be given positions, the task of the algorithm is of course not finished. Although in vacuum this may be all that is required, most

calculations are carried out in explicit solvent. This means that each new atom will likely clash not only with a receptor atom, but also potentially with a solvent atom. These clashes will result in very unfavorable acceptance probabilities and virtually guaranteed rejection. To resolve this, I integrated the approach of [62, 60], known either as nonequilibrium switching or annealed importance sampling, into this algorithm. This approach works as follows. First, the new atoms are introduced in a decoupled state. We then take one step of dynamics, followed by a step changing the control parameters of the simulation. These parameters, canonically called λ , control whether the parameters of the system behave more closely like the initial molecule or the final molecule. We alternate between taking a step of dynamics and a step in control parameter space, accumulating a work—a change in potential energy—each time we change the control parameters. At the conclusion of the protocol, the contribution to the overall acceptance probability becomes simply the exponentiated negative work. In this way, we can insert atoms into densely packed systems, allowing the system to reorganize and dramatically improving acceptance probabilities. This algorithm is attractive in its simplicity but also in its potential for future work. A significant body of work exists regarding the so-called thermodynamic length [17]. Through this formalism, the path through control parameter space that we take can be optimized. As such, the details of the rest of the algorithm can be simplified if this component can “clean up” most unfavorable configurations. There is ample room to design not only on-average better protocols, but even protocols that are specific to each pair of ligands. Such information is useful not only for this algorithm, but also for others using nonequilibrium switching for free energy calculations [1].

9.2.4 Stochastic Approximation

Having collected the necessary components to maintain the invariant distribution over configurations and chemicals, we now turn to achieving an appropriate weighting of these states. Traditionally, in such expanded ensemble simulations, one might resort to a stochastic approximation method to achieve even sampling [88, 84, 42]. This not only results in the individual state weights converging to the relative free energies, but more importantly forces sampling of various states that other otherwise made unfavorable by even moderate free energy differences. However, in this case, we have defined a distribution over many different chemical species. We thus have the opportunity to not merely achieve even sampling, but to adaptively direct sampling toward states that are more likely to be favorable. This is an approach that is very difficult or impossible to achieve in other schemes. Here, I perform it by adding extra recursion steps to the stochastic approximation algorithm, allowing the *target* weights of various states to depend on the current stochastic approximation weights, where before they were fixed. This allows us to adaptively reweight chemical states so that sampling focuses on those that achieve some desirable free-energy based property, such as binding, selectivity, multiple targeting, and more. Much future work can be accomplished in this realm. One of the most pressing issues is a reasonable initialization of the stochastic approximation weights, since the adaptive algorithm performs best when initialized near its optimum. One interesting approach that future work can take is to use a machine learning algorithm, which may be less accurate than the free energy calculation, to initialize the weights. This may provide a "good guess" that will then be refined by the calculation. Another fascinating future approach might involve adapting the weights of a function approximator such as a neural net rather than the individual chemical state weights. This may have

several advantages: it can accommodate the insertion of new chemical states easily, it can then be used as a trained model (or fine-tuned on experiment) after the fact, and it can be initialized to a reasonable starting position. One potential challenge of this approach is to ensure that the convergence of the function approximator’s parameters occurs in a reasonable amount of time. Finally, one might imagine using a Bayesian scheme to determine the weights, a method not explored here but discussed with several colleagues. All of these are promising extensions to the present algorithm.

9.3 Transdimensional Nonequilibrium Switching

Drawing from the work on the chemical space sampling algorithm, it becomes apparent that there is another algorithm lurking within it. This algorithm involves so-called transdimensional nonequilibrium switching. Unlike standard nonequilibrium switching, where a hybrid system consisting of a union of atoms of both endpoints must be simulated at equilibrium, transdimensional nonequilibrium switching simulates only the non-alchemical endpoints. By pushing the atom mapping and dimension jump into the nonequilibrium switching protocol, this algorithm enables the above techniques to be immediately applied on a massive scale with parallel resources. Once equilibrium simulations have begun, in parallel, transdimensional nonequilibrium switching trajectories can be initiated in parallel. This approach offers two major advantages over the previous algorithms. One, it enables only a single equilibrium simulation for each molecule in the set (as opposed to a pair of simulations for every pair), significantly cutting down on equilibrium simulation time. More fascinatingly, it enables the atom mapping parameters to be adapted online, as the equilibrium samples that are

used to start the nonequilibrium trajectories are no longer dependent on the atom map. This is a very powerful approach that lends itself to sophisticated Bayesian design schemes, where molecules can be optimized not only for free-energy based objectives, but also for other heuristic and cost objectives. Future work in this direction will bring the field closer to very large-scale use of free energy calculations earlier in the drug discovery pipeline.

9.4 Conclusion

This work has made three main contributions to the field of free energy computation. First of all, it has developed a scheme which can adaptively sample regions of chemical space according to free-energy based design criteria. This algorithm also leaves open the theoretical requirements for any additional extensions; as such, it can be treated in a "plug and play" manner if the requirements are met. Second, in the course of developing this algorithm, I have also enabled the forming and breaking of fused ring systems, which exist in many molecules relevant to drug discovery and materials science. Finally, I extend the chemical space sampling algorithm to be applicable to the nonequilibrium switching setting, a method I call transdimensional nonequilibrium switching. This enables the above algorithm to be immediately adapted to highly-parallel heterogeneous compute environments like the ones used commonly today. It is my hope that these algorithmic advances lay the groundwork for significantly more work on scaling free energy calculations to situations both with resource constraints and very large sets of molecules.

APPENDIX A

APPENDIX A

A.1 Derivation of the Acceptance Probability

Having assembled all the components for the algorithm to leave the expanded ensemble of configurations and chemical states invariant, I now derive the formula for the acceptance probability. In order to do this, I will impose pathwise detailed balance, ensuring that First, let us review the notation that we will be using.

A.1.1 Proposal of Chemical State Jump

Beginning at a state $(x_{core}, x_{old}, M_{old})$, where we denote the unique degrees of freedom not common to the new molecule x_{old} , the degrees of freedom in common x_{core} , the old molecule identity M_{old} , and denoting \mathcal{T} as the forward transition and $\tilde{\mathcal{T}}$ as the reverse transition:

$$\mathcal{T} = (M_{old} \rightarrow M_{new}) \sim q(\cdot | M_{old}) \quad (\text{A.1})$$

$$order_{old \rightarrow new} \sim order(n_{x_{new}}) \quad (\text{A.2})$$

$$x_{new} \sim \phi(x_{new} | x_{core}, M_{old}, M_{new}) \quad (\text{A.3})$$

$$x_{new}^*, x_{core}^* \sim \Phi(x_{core}, x_{new}, x_{old} \rightarrow x_{core}^*, x_{new}^*, x_{old}^*) \quad (\text{A.4})$$

$$(\text{A.5})$$

Here, x_{new} are the new degrees of freedom in the new molecule that are not shared with the old molecule, and x^* denotes coordinates that are proposed as part of the probabilistic proposal.

A.1.2 Derivation of Acceptance Probability

Now, using $\pi(x, \mathcal{M})$ as the equilibrium distribution, we impose super-detailed balance on the specific proposed transition path \mathcal{T} :

$$\pi(x_{\text{core}}, x_{\text{old}}, \mathcal{M}_{\text{old}}) q[\mathcal{T}] \mathcal{A}[\mathcal{T}] = \pi(x'_{\text{core}}, x'_{\text{new}}, \mathcal{M}_{\text{new}}) q[\tilde{\mathcal{T}}] \mathcal{A}[\tilde{\mathcal{T}}] \quad (\text{A.6})$$

$$\begin{aligned} & \pi(x_{\text{core}}, x_{\text{old}}, \mathcal{M}_{\text{old}}) q(\mathcal{M}_{\text{new}}|\mathcal{M}_{\text{old}}) \phi(x_{\text{new}}|x_{\text{core}}, \mathcal{M}_{\text{old}}, \mathcal{M}_{\text{new}}) \Phi(x \rightarrow x'|\mathcal{M}_{\text{old}} \rightarrow \mathcal{M}_{\text{new}}) \mathcal{A}[\mathcal{T}] \\ = & \pi(x'_{\text{core}}, x'_{\text{new}}, \mathcal{M}_{\text{new}}) P(\mathcal{M}_{\text{old}}|\mathcal{M}_{\text{new}}) \phi(x'_{\text{old}}|x'_{\text{core}}, \mathcal{M}_{\text{new}}, \mathcal{M}_{\text{old}}) \Phi(x' \rightarrow x|\mathcal{M}_{\text{new}} \rightarrow \mathcal{M}_{\text{old}}) \mathcal{A}[\tilde{\mathcal{T}}] \end{aligned} \quad (\text{A.7})$$

Collecting terms, we can arrive at a condition on the acceptance criteria:

$$\begin{aligned} \frac{\mathcal{A}[\mathcal{T}]}{\mathcal{A}[\tilde{\mathcal{T}}]} &= \frac{\pi(x'_{\text{core}}, x'_{\text{new}}, \mathcal{M}_{\text{new}})}{\pi(x_{\text{core}}, x_{\text{old}}, \mathcal{M}_{\text{old}})} \frac{P(\mathcal{M}_{\text{old}}|\mathcal{M}_{\text{new}})}{P(\mathcal{M}_{\text{new}}|\mathcal{M}_{\text{old}})} \frac{\phi(x'_{\text{old}}|x'_{\text{core}}, \mathcal{M}_{\text{new}}, \mathcal{M}_{\text{old}})}{\phi(x_{\text{new}}|x_{\text{core}}, \mathcal{M}_{\text{old}}, \mathcal{M}_{\text{new}})} \\ & \frac{\Phi(x' \rightarrow x|\mathcal{M}_{\text{new}} \rightarrow \mathcal{M}_{\text{old}})}{\Phi(x \rightarrow x'|\mathcal{M}_{\text{old}} \rightarrow \mathcal{M}_{\text{new}})} \\ &= \frac{e^{-u(x'_{\text{core}}, x'_{\text{new}}, \mathcal{M}_{\text{new}}) + g(\mathcal{M}_{\text{new}})}}{e^{-u(x_{\text{core}}, x_{\text{old}}, \mathcal{M}_{\text{old}}) + g(\mathcal{M}_{\text{old}})}} \frac{P(\mathcal{M}_{\text{old}}|\mathcal{M}_{\text{new}})}{P(\mathcal{M}_{\text{new}}|\mathcal{M}_{\text{old}})} \frac{\phi(x'_{\text{old}}|x'_{\text{core}}, \mathcal{M}_{\text{new}}, \mathcal{M}_{\text{old}})}{\phi(x_{\text{new}}|x_{\text{core}}, \mathcal{M}_{\text{old}}, \mathcal{M}_{\text{new}})} \\ & e^{-\Delta S[x \rightarrow x'|\mathcal{M}_{\text{old}} \rightarrow \mathcal{M}_{\text{new}}]} \\ &= \frac{e^{-u(x'_{\text{core}}, x'_{\text{new}}, \mathcal{M}_{\text{new}}) + g(\mathcal{M}_{\text{new}})}}{e^{-u(x_{\text{core}}, x_{\text{old}}, \mathcal{M}_{\text{old}}) + g(\mathcal{M}_{\text{old}})}} \frac{P(\mathcal{M}_{\text{old}}|\mathcal{M}_{\text{new}})}{P(\mathcal{M}_{\text{new}}|\mathcal{M}_{\text{old}})} \frac{\phi(x'_{\text{old}}|x'_{\text{core}}, \mathcal{M}_{\text{new}}, \mathcal{M}_{\text{old}})}{\phi(x_{\text{new}}|x_{\text{core}}, \mathcal{M}_{\text{old}}, \mathcal{M}_{\text{new}})} \\ & e^{-w[x \rightarrow x'|\lambda=0 \rightarrow 1]} \frac{e^{-u(x, \lambda=0)}}{e^{-u(x', \lambda=1)}} \end{aligned} \quad (\text{A.8})$$

The Metropolis-like criteria satisfies this requirement:

$$P_{\text{accept}}[\mathcal{T}] = \min \left\{ 1, e^{-w[x \rightarrow x'|\lambda=0 \rightarrow 1]} \frac{e^{-u(x, \lambda=0)}}{e^{-u(x', \lambda=1)}} \right\} \quad (\text{A.9})$$

BIBLIOGRAPHY

- [1] Matteo Aldeghi, Bert L. de Groot, and Vytautas Gapsys. Accurate calculation of free energy changes upon amino acid mutation. In *Methods in Molecular Biology*, pages 19–47. Springer New York, sep 2018.
- [2] Christophe Andrieu and Yves Atchade. On the efficiency of adaptive MCMC algorithms. *Electronic Communications in Probability*, 12(0):336–349, 2007.
- [3] Christophe Andrieu, Arnaud Doucet, and Roman Holenstein. Particle markov chain monte carlo methods. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 72(3):269–342, jun 2010.
- [4] Christophe Andrieu, Éric Moulines, and Pierre Priouret. Stability of stochastic approximation under verifiable conditions. *SIAM Journal on Control and Optimization*, 44(1):283–312, jan 2005.
- [5] Christophe Andrieu and Johannes Thoms. A tutorial on adaptive mcmc. *Statistics and Computing*, 18(4):343–373, Dec 2008.
- [6] Johan qvist, Carmen Medina, and Jan-Erik Samuelsson. A new method for predicting binding affinity in computer-aided drug design. *"Protein Engineering, Design and Selection"*, 7(3):385–391, 1994.
- [7] Charles H Bennett. Efficient estimation of free energy differences from monte carlo data. *Journal of Computational Physics*, 22(2):245–268, 1976.
- [8] Matko Bošnjak, Tim Rocktäschel, Jason Naradowsky, and Sebastian Riedel. Programming with a differentiable forth interpreter. In *Proceedings of the 34th International Conference on Machine Learning*, 2017.
- [9] Steve Brooks, Andrew Gelman, Galin Jones, and Xiao-Li Meng. *Handbook of Markov Chain Monte Carlo*. CRC press, 2011.
- [10] Hongming Chen, Ola Engkvist, Yinhai Wang, Marcus Olivecrona, and Thomas Blaschke. The rise of deep learning in drug discovery. *Drug Discovery Today*, 23(6):1241–1250, jun 2018.
- [11] Tiejun Cheng, Qingliang Li, Zhigang Zhou, Yanli Wang, and Stephen H. Bryant. Structure-based virtual screening for drug discovery: a problem-centric review. *The AAPS Journal*, 14(1):133–141, Mar 2012.

- [12] John D. Chodera and Michael R. Shirts. Replica exchange and expanded ensemble simulations as gibbs sampling: Simple improvements for enhanced mixing. *J. Chem. Phys.*, 135(19):194110, 2011.
- [13] Ian Churcher. Protac-induced protein degradation in drug discovery: Breaking the rules or just making new ones? *Journal of Medicinal Chemistry*, 61(2):444–452, dec 2017.
- [14] Lucy J Colwell. Statistical and machine learning approaches to predicting protein–ligand interactions. *Current Opinion in Structural Biology*, 49:123–128, apr 2018.
- [15] NICOLAS COMBE, THIJS J. H. VLUGT, PIETER REIN TEN WOLDE, and DAAN FRENKEL. Dynamic pruned-enriched rosenbluth method. *Molecular Physics*, 101(11):1675–1682, jun 2003.
- [16] Zoe Cournia, Bryce Allen, and Woody Sherman. Relative binding free energy calculations in drug discovery: Recent advances and practical considerations. *Journal of Chemical Information and Modeling*, 57(12):2911–2937, dec 2017.
- [17] Gavin E. Crooks. Measuring thermodynamic length. *Physical Review Letters*, 99(10), sep 2007.
- [18] Gavin Earl Crooks. *Excursions in statistical dynamics*. PhD thesis, University of California, Berkeley, 1999.
- [19] Tom Darden, Darrin York, and Lee Pedersen. Particle mesh ewald: An $n\log(n)$ method for ewald sums in large systems. *The Journal of Chemical Physics*, 98(12):10089–10092, jun 1993.
- [20] Anirban DasGupta. 29 review of optimal bayes designs. In *Handbook of Statistics*, pages 1099–1147. Elsevier, 1996.
- [21] Christoph Dellago and Gerhard Hummer. Computing equilibrium free energies using non-equilibrium molecular dynamics. *Entropy*, 16(1):41–61, 2014.
- [22] Xinqiang Ding, Jonah Z. Vilseck, Ryan L. Hayes, and Charles L. Brooks. Gibbs sampler-based -dynamics and rao–blackwell estimator for alchemical free energy calculation. *Journal of Chemical Theory and Computation*, 13(6):2501–2510, may 2017.

- [23] R. Douc and O. Cappe. Comparison of resampling schemes for particle filtering. In *ISPA 2005. Proceedings of the 4th International Symposium on Image and Signal Processing and Analysis*, 2005. IEEE, 2005.
- [24] Peter Eastman, Jason Swails, John D. Chodera, Robert T. McGibbon, Yutong Zhao, Kyle A. Beauchamp, Lee-Ping Wang, Andrew C. Simmonett, Matthew P. Harrigan, Chaya D. Stern, Rafal P. Wiewiora, Bernard R. Brooks, and Vijay S. Pande. OpenMM 7: Rapid development of high performance algorithms for molecular dynamics. *PLOS Computational Biology*, 13(7):e1005659, jul 2017.
- [25] Péter Englert and Péter Kovács. Efficient heuristics for maximum common substructure search. *Journal of Chemical Information and Modeling*, 55(5):941–955, may 2015.
- [26] Niklas W. A. Gebauer, Michael Gastegger, and Kristof T. Schütt. Generating equilibrium molecules with deep neural networks. *CoRR*, abs/1810.11347, 2018.
- [27] Andrew Gelman and Xiao-Li Meng. Simulating normalizing constants: from importance sampling to bridge sampling to path sampling. *Statistical Science*, 13(2):163–185, may 1998.
- [28] M.K. Gilson, J.A. Given, B.L. Bush, and J.A. McCammon. The statistical-thermodynamic basis for computation of binding affinities: a critical review. *Biophysical Journal*, 72(3):1047–1069, mar 1997.
- [29] Todd R. Gingrich, Grant M. Rotskoff, Gavin E. Crooks, and Phillip L. Geissler. Near-optimal protocols in complex nonequilibrium transformations. *Proceedings of the National Academy of Sciences*, 113(37):10263–10268, aug 2016.
- [30] PETER J. GREEN. Reversible jump markov chain monte carlo computation and bayesian model determination. *Biometrika*, 82(4):711–732, 1995.
- [31] Anvita Gupta, Alex T. Müller, Berend J. H. Huisman, Jens A. Fuchs, Petra Schneider, and Gisbert Schneider. Generative recurrent networks for de novo drug design. *Molecular Informatics*, 37(1-2):1700111, nov 2017.
- [32] Aric A. Hagberg, Daniel A. Schult, and Pieter J. Swart. Exploring network structure, dynamics, and function using networkx. In Gaël Varoquaux, Travis Vaught, and Jarrod Millman, editors, *Proceedings of the 7th Python in Science Conference*, pages 11 – 15, Pasadena, CA USA, 2008.

- [33] W. K. Hastings. Monte carlo sampling methods using markov chains and their applications. *Biometrika*, 57(1):97–109, apr 1970.
- [34] Soren Henriksen, Adrian Wills, Thomas B. Schön, and Brett Ninness. Parallel implementation of particle MCMC methods on a GPU*. *IFACProceedingsVolumes*, 45(16) : 1143 – –1148, jul2012.
- [35] Viktor Hornak and Carlos Simmerling. Development of softcore potential functions for overcoming steric barriers in molecular dynamics simulations. *Journal of Molecular Graphics and Modelling*, 22(5):405–413, may 2004.
- [36] Gerhard Hummer. Fast-growth thermodynamic integration: Error and efficiency analysis. *The Journal of Chemical Physics*, 114(17):7330–7337, 2001.
- [37] B. Javadi, R. K. Thulasiramy, and R. Buyya. Statistical modeling of spot instance prices in public cloud environments. In *2011 Fourth IEEE International Conference on Utility and Cloud Computing*. IEEE, dec 2011.
- [38] Jennifer L. Knight and Charles L. Brooks. -dynamics free energy simulation methods. *Journal of Computational Chemistry*, 30(11):1692–1700, aug 2009.
- [39] Jennifer L. Knight and Charles L. Brooks. Multisite dynamics for simulated structure–activity relationship studies. *Journal of Chemical Theory and Computation*, 7(9):2728–2739, sep 2011.
- [40] A. Lee, C. Yau, M.B. Giles, A. Doucet, and C.C. Holmes. On the utility of graphics cards to perform massively parallel simulation of advanced monte carlo methods. *Journal of Computational and Graphical Statistics*, 19(4):769–789, 2010.
- [41] G. M. Lee and C. S. Craik. Trapping moving targets with small molecules. *Science*, 324(5924):213–215, apr 2009.
- [42] Faming Liang, Chuanhai Liu, and Raymond J Carroll. Stochastic approximation in monte carlo computation. *Journal of the American Statistical Association*, 102(477):305–320, mar 2007.
- [43] Jun S. Liu. *Monte Carlo Strategies in Scientific Computing*. Springer New York, 2004.

- [44] Shuai Liu, Lingle Wang, and David L. Mobley. Is ring breaking feasible in relative binding free energy calculations? *Journal of Chemical Information and Modeling*, 55(4):727–735, apr 2015.
- [45] Yu-Chen Lo, Stefano E. Rensi, Wen Torng, and Russ B. Altman. Machine learning in chemoinformatics and drug discovery. *Drug Discovery Today*, 23(8):1538–1546, aug 2018.
- [46] Pinyi Lu, David R. Bevan, Andrew Leber, Raquel Hontecillas, Nuria Tubau-Juni, and Josep Bassaganya-Riera. *Computer-Aided Drug Discovery*, pages 7–24. Springer International Publishing, Cham, 2018.
- [47] A. P. Lyubartsev, A. A. Martsinovski, S. V. Shevkunov, and P. N. Vorontsov-Velyaminov. New approach to monte carlo calculation of the free energy: Method of expanded ensembles. *J. Chem. Phys.*, 96(3):1776, 1992.
- [48] Soma Mandal, Meenal Moudgil, and Sanat K. Mandal. Rational drug design. *European Journal of Pharmacology*, 625(1-3):90–100, dec 2009.
- [49] Guilherme Duarte Ramos Matos, Daisy Y. Kyu, Hannes H. Loeffler, John D. Chodera, Michael R. Shirts, and David L. Mobley. Approaches for calculating solvation free energies and enthalpies demonstrated with an update of the FreeSolv database. *Journal of Chemical & Engineering Data*, 62(5):1559–1569, apr 2017.
- [50] Xiao-Li Meng and Wing Hung Wong. Simulating ratios of normalizing constants via a simple identity : a theoretical exploration. 1996.
- [51] Xiao-Li Meng and Wing Hung Wong. Simulating ratios of normalizing constants via a simple identity: A theoretical exploration. *Statistica Sinica*, 6(4):831–860, 1996.
- [52] Julien Michel, Nicolas Foloppe, and Jonathan W. Essex. Rigorous free energy calculations in structure-based drug design. *Molecular Informatics*, 29(8-9):570–578, jul 2010.
- [53] David Minh. Protein-ligand binding potential of mean force calculations with hamiltonian replica exchange on alchemical interaction grids. 07 2015.
- [54] David D. L. Minh. Implicit ligand theory: Rigorous binding free energies and thermodynamic expectations from molecular docking. *The Journal of Chemical Physics*, 137(10):104106, sep 2012.

- [55] David D. L. Minh and John D. Chodera. Estimating equilibrium ensemble averages using multiple time slices from driven nonequilibrium processes: Theory and application to free energies, moments, and thermodynamic length in single-molecule pulling experiments. *The Journal of Chemical Physics*, 134(2):024111, jan 2011.
- [56] David L. Mobley, Ken A. Dill, and John D. Chodera. Treating entropy and conformational changes in implicit solvent simulations of small molecules. *J. Phys. Chem. B*, 112(3):938–946, jan 2008.
- [57] David L. Mobley and J. Peter Guthrie. FreeSolv: a database of experimental and calculated hydration free energies, with input files. *Journal of Computer-Aided Molecular Design*, 28(7):711–720, jun 2014.
- [58] Pierre Del Moral. Nonlinear filtering: Interacting particle resolution. *Comptes Rendus de l’Académie des Sciences - Series I - Mathematics*, 325(6):653–658, sep 1997.
- [59] Lawrence M. Murray, Anthony Lee, and Pierre E. Jacob. Parallel resampling in the particle filter. *Journal of Computational and Graphical Statistics*, 25(3):789–805, jul 2016.
- [60] Radford M Neal. Annealed importance sampling. *Statistics and computing*, 11(2):125–139, 2001.
- [61] Trung Hai Nguyen and David D. L. Minh. Implicit ligand theory for relative binding free energies. *The Journal of Chemical Physics*, 148(10):104114, mar 2018.
- [62] J. P. Nilmeier, G. E. Crooks, D. D. L. Minh, and J. D. Chodera. Nonequilibrium candidate monte carlo is an efficient tool for equilibrium simulation. *Proceedings of the National Academy of Sciences*, 108(45):E1009–E1018, oct 2011.
- [63] Ruth Nussinov and Chung-Jung Tsai. Allostery in disease and in drug discovery. *Cell*, 153(2):293–305, apr 2013.
- [64] Sanghyun Park, Daniel L. Ensign, and Vijay S. Pande. Bayesian update method for adaptive weighted sampling. *Physical Review E*, 74(6), dec 2006.
- [65] Keith Paton. An algorithm for finding a fundamental set of cycles of a graph. *Communications of the ACM*, 12(9):514–518, 1969.

- [66] Steven M. Paul, Daniel S. Mytelka, Christopher T. Dunwiddie, Charles C. Persinger, Bernard H. Munos, Stacy R. Lindborg, and Aaron L. Schacht. How to improve r&d productivity: the pharmaceutical industrys grand challenge. *Nature Reviews Drug Discovery*, 9(3):203–214, feb 2010.
- [67] Jed Pitera and Peter Kollman. Designing an optimum guest for a host using multimolecule free energy calculations: predicting the best ligand for rebeks “tennis ball”. *Journal of the American Chemical Society*, 120(30):7557–7567, aug 1998.
- [68] Jay W. Ponder and David A. Case. Force fields for protein simulations. In *Protein Simulations*, pages 27–85. Elsevier, 2003.
- [69] C. Radhakrishna Rao. Information and the accuracy attainable in the estimation of statistical parameters. In *Springer Series in Statistics*, pages 235–247. Springer New York, 1992.
- [70] H. Robbins and S. Monro. A stochastic approximation method. *Annals of Mathematical Statistics*, 22:400–407, 1951.
- [71] Gareth O. Roberts and Jeffrey S. Rosenthal. Examples of adaptive MCMC. *Journal of Computational and Graphical Statistics*, 18(2):349–367, jan 2009.
- [72] Lars Ruddigkeit, Ruud van Deursen, Lorenz C. Blum, and Jean-Louis Reymond. Enumeration of 166 billion organic small molecules in the chemical universe database GDB-17. *Journal of Chemical Information and Modeling*, 52(11):2864–2875, nov 2012.
- [73] Rita Santos, Oleg Ursu, Anna Gaulton, A. Patrícia Bento, Ramesh S. Donadi, Cristian G. Bologa, Anneli Karlsson, Bissan Al-Lazikani, Anne Hersey, Tudor I. Oprea, and John P. Overington. A comprehensive map of molecular drug targets. *Nature Reviews Drug Discovery*, 16(1):19–34, dec 2016.
- [74] Michael R Shirts and John D Chodera. Statistically optimal analysis of samples from multiple equilibrium states. *The Journal of chemical physics*, 129(12):124105, 2008.
- [75] Michael R. Shirts, David L. Mobley, and John D. Chodera. Chapter 4 alchemical free energy calculations: Ready for prime time? In *Annual Reports in Computational Chemistry*, pages 41–59. Elsevier, 2007.

- [76] Michael R Shirts and Vijay S Pande. Comparison of efficiency and bias of free energies computed by exponential averaging, the bennett acceptance ratio, and thermodynamic integration. *The Journal of chemical physics*, 122(14):144107, 2005.
- [77] Brian K. Shoichet. Virtual screening of chemical libraries. *Nature*, 432(7019):862–865, dec 2004.
- [78] Brian K Shoichet, Susan L McGovern, Binqing Wei, and John J Irwin. Lead discovery using molecular docking. *Current Opinion in Chemical Biology*, 6(4):439–446, aug 2002.
- [79] Jörn Ilja Siepmann and Daan Frenkel. Configurational bias monte carlo: a new sampling scheme for flexible chains. *Molecular Physics*, 75(1):59–70, jan 1992.
- [80] David A. Sivak and Gavin E. Crooks. Thermodynamic metrics and optimal paths. *Physical Review Letters*, 108(19), may 2012.
- [81] Richard S. Sutton. Learning to predict by the methods of temporal differences. *Machine Learning*, 3(1):9–44, aug 1988.
- [82] 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. *Biophysical Journal*, 86(1):67–74, jan 2004.
- [83] Chunhu Tan, Lijiang Yang, and Ray Luo. How well does poisson-boltzmann implicit solvent agree with explicit solvent? a quantitative analysis. *The Journal of Physical Chemistry B*, 110(37):18680–18687, sep 2006.
- [84] Zhiqiang Tan. Optimally adjusted mixture sampling and locally weighted histogram analysis. *Journal of Computational and Graphical Statistics*, 26(1):54–65, jan 2017.
- [85] Christopher P. Tinworth, Hannah Lithgow, and Ian Churcher. Small molecule-mediated protein knockdown as a new approach to drug discovery. *MedChemComm*, 7(12):2206–2216, 2016.
- [86] Don van Ravenzwaaij, Pete Cassey, and Scott D. Brown. A simple introduction to markov chain monte-carlo sampling. *Psychonomic Bulletin & Review*, 25(1):143–154, Feb 2018.

- [87] J. Vanlier, C. A. Tiemann, P. A. J. Hilbers, and N. A. W. van Riel. A bayesian approach to targeted experiment design. *Bioinformatics*, 28(8):1136–1142, feb 2012.
- [88] Fugao Wang and D. P. Landau. Efficient, multiple-range random walk algorithm to calculate the density of states. *Physical Review Letters*, 86(10):2050–2053, mar 2001.
- [89] Lingle Wang, Yujie Wu, Yuqing Deng, Byungchan Kim, Levi Pierce, Goran Krilov, Dmitry Lupyan, Shaughnessy Robinson, Markus K. Dahlgren, Jeremy Greenwood, Donna L. Romero, Craig Masse, Jennifer L. Knight, Thomas Steinbrecher, Thijs Beuming, Wolfgang Damm, Ed Harder, Woody Sherman, Mark Brewer, Ron Wester, Mark Murcko, Leah Frye, Ramy Farid, Teng Lin, David L. Mobley, William L. Jorgensen, Bruce J. Berne, Richard A. Friesner, and Robert Abel. Accurate and reliable prediction of relative ligand binding potency in prospective drug discovery by way of a modern free-energy calculation protocol and force field. *Journal of the American Chemical Society*, 137(7):2695–2703, feb 2015.
- [90] Wendy A. Warr. A short review of chemical reaction database systems, computer-aided synthesis design, reaction prediction and synthetic feasibility. *Molecular Informatics*, 33(6-7):469–476, jun 2014.
- [91] Gregory L. Warren, C. Webster Andrews, Anna-Maria Capelli, Brian Clarke, Judith LaLonde, Millard H. Lambert, Mika Lindvall, Neysa Nevins, Simon F. Semus, Stefan Senger, Giovanna Tedesco, Ian D. Wall, James M. Woolven, Catherine E. Peishoff, and Martha S. Head. A critical assessment of docking programs and scoring functions. *Journal of Medicinal Chemistry*, 49(20):5912–5931, oct 2006.
- [92] Christopher J. C. H. Watkins and Peter Dayan. Q-learning. *Machine Learning*, 8(3):279–292, May 1992.
- [93] Collin D. Wick and J. Ilja Siepmann. Self-adapting fixed-end-point configurational-bias monte carlo method for the regrowth of interior segments of chain molecules with strong intramolecular interactions. *Macromolecules*, 33(19):7207–7218, sep 2000.
- [94] Zhenqin Wu, Bharath Ramsundar, Evan N. Feinberg, Joseph Gomes, Caleb Geniesse, Aneesh S. Pappu, Karl Leswing, and Vijay Pande. MoleculeNet: a benchmark for molecular machine learning. *Chemical Science*, 9(2):513–530, 2018.

- [95] Bing Xie, Trung Hai Nguyen, and David D. L. Minh. Absolute binding free energies between t4 lysozyme and 141 small molecules: Calculations based on multiple rigid receptor configurations. *Journal of Chemical Theory and Computation*, 13(6):2930–2944, may 2017.
- [96] Patrick R. Zulkowski, David A. Sivak, Gavin E. Crooks, and Michael R. DeWeese. Geometry of thermodynamic control. *Physical Review E*, 86(4), oct 2012.
- [97] Robert W. Zwanzig. High-temperature equation of state by a perturbation method. i. nonpolar gases. *The Journal of Chemical Physics*, 22(8):1420–1426, aug 1954.