

Capturing the Role of Temperature and Entropy in Crystal Polymorphs Through Molecular Dynamics Simulations

A Dissertation

Presented to
the faculty of the School of Engineering and Applied Science
University of Virginia

in partial fulfillment
of the requirements for the degree

Doctor of Philosophy

by

Eric Dybeck

December

2016

APPROVAL SHEET

The dissertation
is submitted in partial fulfillment of the requirements
for the degree of
Doctor of Philosophy


AUTHOR

The dissertation has been read and approved by the examining committee:

Michael Shirts

Advisor

David Green

Elizabeth Opila

Joshua Choi

Gaurav Giri

Accepted for the School of Engineering and Applied Science:



Dean, School of Engineering and Applied Science

December

2016

Acknowledgements

I would like to acknowledge the guidance and support provided by my advisor Michael Shirts. I would also like to thank Levi Naden, Natalie Schieber, Nate Abraham and the entire Shirts group for their help and technical support during my research. I also would like to acknowledge the UVA high performance computing team for the computational resources necessary to complete this work. Finally, I would like to acknowledge the endless love and support that my friends and family provided me during this challenging Ph.D. journey.

Abstract

The presence of multiple stable crystal structures in solid organic materials can limit their commercial viability when two or more observable polymorphs exhibit markedly different physical properties. Unintended restructuring events have hindered pharmaceutical solid form development in numerous therapeutic candidates and have led to costly market recalls. Conversely, intentional synthesis of a metastable form has led to significantly more favorable performance in numerous materials.

Current computational methods for predicting polymorphic behavior evaluate candidate crystal structures based on the minimized lattice energy. However, these static lattice energy-based approaches generate far more lattice energy minima than there are experimentally observed structures. Thermal motion of the crystals under working conditions has the potential to explain why many of these lattice minima are not observed experimentally. Lattice minima that are identified from a static crystal structure prediction can ultimately be unstable at experimental conditions through either temperature-mediated stability reranking, kinetic interconversion of multiple minima, or inaccuracies of the energy function in producing the real crystal ensemble.

In this work, we explore the role of thermal motion in eliminating candidate crystal structures using fully atomistic molecular dynamics simulations. Enthalpically favorable structures with low entropy will become unfavorable at high temperatures through temperature-mediated stability reranking. Our simulations correctly identify the high temperature solid form as having a larger entropy than the low temperature form in twelve small molecule organic systems with known temperature-mediated transformations. The estimated entropy differences in the classical point-charge potential are significantly closer to experimental measurements than estimated enthalpy differences. This result suggests that entropy difference estimates are less sensitive to the complexity of the simulation potential than the corresponding enthalpy estimates. We additionally find that a cheaper harmonic approximation provides a sufficient estimate of entropic contributions in small rigid molecules. However, entropies with

the harmonic approximation diverge from molecular dynamics-derived entropies in systems with multiple rotatable degrees of freedom or dynamically disordered crystalline structures.

We additionally probe the sensitivity of the stability estimates to the energy function by directly computing the added effects of a more accurate polarizable Hamiltonian on polymorph free energies using a novel Hamiltonian reweighting approach. We show that the change in free energy to the more complex potential comes predominately from enthalpic rather than entropic contributions in the system examined.

Finally, we demonstrate the utility of molecular dynamics in identifying lattice minima interconversion and order-disorder transitions in organic solids. A rapid conversion of multiple lattice minima into a single ambient temperature ensemble is presented in six of the systems examined. These kinetics events can significantly reduce the number of plausible candidate structures in future crystal structure prediction studies.

Contents

Acknowledgements	i
Abstract	iii
1 Introduction to Crystal Polymorphism	1
1.1 What is Polymorphism?	1
1.2 Case Studies in Industrial Polymorphism	3
1.2.1 Market Recalls due to Unexpected Transformation	3
1.2.2 Legal Controversy and Disappearing Polymorphs	5
1.2.3 Improving Material Performance Through Crystal Polymorphism	5
1.3 The Experimental Search for All Solid Forms	7
1.3.1 Characterizing Distinct Solid Forms	7
1.3.2 What Factors Lead to Different Polymorphs?	7
1.3.3 Screening for New Solid Forms	9
1.4 Computer-Aided Crystal Structure Prediction	10
1.5 Why Do Computer Models Find So Many Possible Structures?	12
2 Methods for Modeling Crystal Polymorphs	17
2.1 Searching the Crystal Landscape for Candidate Structures	17
2.1.1 Generating Candidate Structures	17
2.1.2 Refining Structures through Lattice Energy Minimization	20
2.1.3 The Need for Ambient Temperature Structures	22
2.2 Statistical Mechanics Formalism	23

2.3	Estimating Entropy with a Quasi-Harmonic Approximation	27
2.3.1	The Harmonic Approximation	27
2.3.2	The Quasi-Harmonic Approximation	28
2.3.3	Limitations of the Quasi-Harmonic Approximation	28
2.4	Introduction to Molecular Dynamics	29
2.4.1	Sampling the Configuration Space	29
2.4.2	Periodic Boundaries and Long-range Interactions	31
2.4.3	Temperature and Pressure Coupling	32
2.5	Computing Crystal Free Energies	34
2.5.1	Generating Crystal Ensembles with Molecular Dynamics	34
2.5.2	Defining a Free Energy Minima	35
2.5.3	Interconverting Crystal Polymorphs along a Pseudo-supercritical Path	37
2.5.4	Multistate Reweighting of MD Ensembles with MBAR	39
2.5.5	The Limits of Thermodynamics-based Crystal Structure Prediction	40
2.6	Eliminating Unobservable Crystals with Molecular Dynamics	42
3	Predicting Temperature-mediated Polymorphic Transformations	47
3.1	Introduction	47
3.2	Methods	50
3.2.1	Initial System Structures	50
3.2.2	Free Energy Estimation with MBAR	52
3.2.3	Exploring the Magnitude of Entropic Contributions	54
3.2.4	Quasi-harmonic Approximation	56
3.2.5	Simulation Details	57
3.3	Results and Discussion	60
3.3.1	Entropy and Enthalpy at Finite Temperature	60
3.3.2	Comparison of Molecular Dynamics and QHA	65
3.4	Conclusions	69

4	Exploring the Effects of Polarization Through Hamiltonian Reweighting	71
4.1	Introduction	71
4.2	Methods	78
4.2.1	Computing free energies between point charge potential polymorphs using a pseudo-supercritical path	78
4.2.2	Point-charge and Polarizable Potentials	80
4.3	Simulation Details	83
4.4	Results and Discussion	87
4.5	Conclusions	104
5	Lattice Minima Interconversion at Ambient Temperatures	107
5.1	Introduction	107
5.2	Methodology	110
5.3	Minima Conversion	111
5.4	Minima Interconversion	115
5.5	Conclusions and Future Outlook	116
6	Evaluating Methods to Reweight from MM to QM/MM Potentials	121
6.1	Introduction	121
6.2	Reweighting Methods	126
6.2.1	Non-Boltzmann Bennett	126
6.2.2	Multistate Bennett Acceptance Ratio	131
6.2.3	Effective Number of Samples	135
6.3	Simulation Details	136
6.3.1	Ethane-Methanol	136
6.3.2	SAMPL4 subset	138
6.3.3	Harmonic Oscillators	139
6.3.4	Uncertainty Analysis	140
6.4	Results and Discussion	141

6.4.1	Ethane-Methanol	141
6.4.2	SAMPL4 subset	148
6.4.3	Harmonic Oscillators	156
6.5	Conclusions	160
A	Appendix A	163
A.1	Derivations of Equations	163
A.1.1	Relating the Helmholtz Free Energy to the Gibbs Free Energy	163
A.1.2	Computing free energy differences at multiple temperatures with MBAR	164
B	Appendix B	167
B.1	Table of Experimental Polymorph Data	167
B.2	Extrapolation of Free Energy from Finite Temperature to 0 K	167
B.3	Complete PSCP Cycle Closure	168
C	Appendix C	171
C.1	End State Hamiltonian Reweighting for All Polymorphs	171
C.2	Indirect AMOEBA09 Free Energy vs Temperature Convergence	171
C.3	Potentials and Functional Forms	175
C.4	Pseudo-Supercritical Path Calculations	181
C.5	Benzene Polymorph Free Energies using 4- and 72- Benzene Model	182
C.6	Polymorph Free Energy is Insensitive to Cutoff Distance	187
D	Appendix D	189
D.1	Offset Harmonic Oscillators	189
D.2	Mathematical differences between different reweighting schemes	190

Chapter 1

Introduction to Crystal Polymorphism

1.1 What is Polymorphism?

Many solid materials can exist in multiple distinct and stable crystalline structures [1]. Essentially, a single compound can have multiple ways to pack together into a repeating crystalline lattice in the solid phase. The element carbon is a familiar example of a material that can adopt multiple crystal forms. The most stable form of carbon at ambient conditions is graphite sheets [2]. However, it is well known that carbon can also be observed as diamond [2], graphene [3], or even the more curious C₆₀ Buckminsterfullerene [4]. Other common examples of materials with multiple crystalline structures include close-packed metals like iron which can be seen in both a body-centered cubic (BCC) and face-centered cubic (FCC) crystal packing depending on the temperature and pressure [5].

When a material has multiple observable crystal structures, each form is referred to as a 'polymorph'. McCrone [6] provides the following definition: "A polymorph is a solid crystalline phase of a given compound resulting from the possibility of at least two different arrangements of the molecules of that compound in the solid state". The crystallographic use of the term 'polymorph' dates back to the 1820s with observations of different structures of arsenate and phosphate salts [1]. Since that time, polymorphism has seen steadily increasing interest from the scientific community from numerous high impact research groups [1,7–11].

The scientifically alluring aspect of polymorphic solids is that a material's properties can be greatly influenced by the specific crystal structure that it is currently adopting. In the example

of carbon mentioned above, graphite is a relatively soft and brittle material, whereas diamond is one of the hardest naturally occurring materials in the world. Although the hardness is arguably the most noticeable difference between graphite and diamond, the two forms also differ in optical transparency, electrical conductivity [12], and thermal conductivity [13] among other properties.

Many commercial products such as drugs [14–19], explosives [20–24], dyes [25], and electrochemical materials [26–30] involve crystallized organic molecules with multiple solid forms. Understanding polymorphism is important in these commercial materials because the structural and thermodynamic properties of the materials can change between different crystal structures. Properties that can change between polymorphs include density [11, 20, 31], compressibility [32], hardness [33, 34], gel strength [35], flexibility [36], conductivity [26–30, 37–39], explosivity [20, 22–24], catalytic activity [40, 41], transparency [42, 43], weather resistance [25, 44], heat capacity [20], melting point [32, 45, 46], vapor pressure [46], stability [25, 29], solubility [14, 47–49], bioavailability [14, 33, 48, 50–54], and dissolution rate [47, 55, 56].

One or more of the above properties may be the aspect that makes a material commercially viable. Crystallization of the material into a different form with new properties could ultimately compromise the intended performance. Therefore, commercial interest in a material can be as much dependent on the crystal structure as on the chemical composition. Commercial manufacturers will typically send candidate materials through extensive polymorphic screening experiments to identify all possible solid forms, either to improve the material's performance [27, 33, 35, 38] or to avoid an unexpected transformation in the future [14, 48, 50]. In the following section, a number of industrially relevant case studies are described involving polymorphic materials. In each case, the presence of multiple stable crystal structures had a significant impact on the effectiveness and/or marketability of the product.

1.2 Case Studies in Industrial Polymorphism

1.2.1 Market Recalls due to Unexpected Transformation

One of the more salient concerns in pharmaceutical manufacturing is the appearance of a new unknown crystal form late in a drug's development process. In many cases, the new forms appear during the scale-up from the lab or pilot plant to full scale production. A sharp change in the materials' hardness, flowability, or compressibility as a result of this new form can render the original manufacturing process unviable and require costly reengineering of the production line.

The most serious cases of late crystal form appearance occur when the product has already entered the commercial market. The most famous case of this late appearance was Abbott's HIV drug ritonavir. The drug was originally formulated in a semisolid gel capsule based on the only known solid form at the time (Form I) [50, 57]. However, the globally stable Form II crystal ultimately precipitated out of the capsule due to a lower solubility in the hydroalcoholic solution and led to failed dissolution in patients. Furthermore, the labs which were now seeded with Form II ritonavir could no longer produce Form I even when following the original crystallization procedure. The drug was eventually reformulated and reintroduced to the market in 1999, three years after the original release of the product. During this time, it is estimated that Abbott lost \$250 million in potential sales and delayed the treatment of thousands of patients [14].

Another example of late crystal form appearance forcing a market recall is the drug rotigotine. This drug is used to treat Parkinson's disease as well as restless leg syndrome [57, 58]. Rotigotine was originally formulated in a transdermal patch with the known Form I solubilized in a polymer matrix [59]. However, the unknown Form II began crystallizing out of the patch into snowflake-like crystals (shown in Figure 1.1) shortly after production due to a higher stability and lower solubility [15]. These solid crystals prevent the necessary transdermal diffusion of the active pharmaceutical ingredient (API) when the patch is administered. The drug was eventually reformulated and launched in 2012. However, between 2008 and

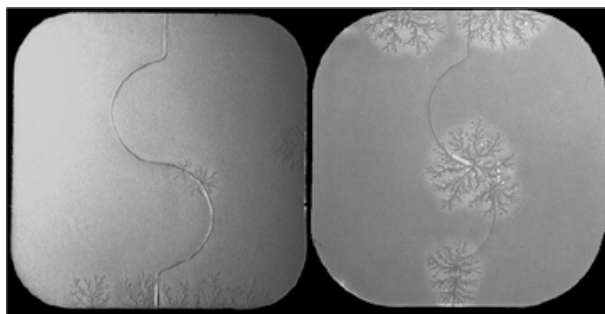


Figure 1.1: Rotigotine is administered using a polymer matrix transdermal patch. Form I rotigotine is soluble in the polymer matrix (Left). Form II rotigotine is insoluble in the patch and precipitated out into snowflake-like crystal which forced a market recall.

2012 the US experienced a complete unavailability of the drug [57].

The FDA now mandates that manufacturers screen for polymorphs well ahead of mainstream production to demonstrate that potential materials does not have unknown and unwanted crystal structures. The details of the experimental procedure for seeking out new polymorphs is discussed further in the next section. If polymorphism does exist in the candidate product, manufacturers must demonstrate that the multiple crystals have nearly equivalent bioavailability between crystal structures or that the crystal structures can be reliably stabilized in the desired form.

Commercial manufacturers must also demonstrate that the correct crystal structures can be reliably produced by following the given synthesis procedures. The FDA certifies a drug in a specific solid form, meaning that the drug must always be created and distributed in that specific form. Furthermore, any crystallographic transformation of the solid drug form during transport or storage necessarily leads to the drug being considered expired. An expiration date therefore represents the amount of time in which a drug is guaranteed to remain in its manufactured solid form as long as the storage instructions are properly followed (such as being stored in a cool and dry environment). Further details of how a crystal can be unintentionally synthesized in or subsequently transform into an undesired metastable form are provided in the following section.

1.2.2 Legal Controversy and Disappearing Polymorphs

The phenomena of one crystal being synthesized in the lab-scale, but then no longer reproducible at the industrial scale is known as ‘disappearing polymorphs’ [7]. These disappearing polymorphs have led to legal controversy regarding how broadly a patent’s protection extends when multiple crystal structures exist [57]. In general, one patents a particular crystallization process which therefore patents a particular crystal form of an API. When a new crystal structure is formed during the industrial scaleup, this opens the door for the new crystallization process and crystal structure to be repatented, which can increase the revenue of a drug considerably.

For example, ranitidine hydrochloride (known commonly as its commercial name Zantac[®]) was first patented by Glaxo in 1978 in the form 1 crystal. Form 2 was immediately found upon scaleup and was subsequently repatented by Glaxo in 1985. When the form 1 crystal patent expired, generic companies following the original patent always produced form 2, which was still patented by Glaxo. As a result, generic production was prevented until the second patent expired in 2002 [57]. Glaxo effectively extended the patent on their lucrative drug Zantac by 7 years through an appropriate understanding of crystal polymorphism.

Similar examples of patent extension through discovering a new crystal form have been reported for Paroxetine Hydrochloride (Paxil[®]) [60] and Cefdinir [61]. Large pharmaceutical companies are actively looking for new crystal structures to extend patents for this reason. Conversely, a generic company could theoretically circumvent active patents or block a patent extension through preemptively patenting a new crystalline form. Therefore larger pharmaceutical companies must patent all possible crystal forms to fully protect their intellectual property.

1.2.3 Improving Material Performance Through Crystal Polymorphism

A full understanding of alternate crystal structures can in some cases lead to improved material performance. One clear example of improved performance through crystal rearrangement

appears in the organic semiconductor industry. Organic electronics exhibit many potential advantages over traditional silicon based materials in products such as flexible light sources, portable solar cells, and curved television screens [27, 38, 62, 63]. Other advantages include cheaper manufacturing procedures like large-area deposition on flexible substrates [38, 64].

A key factor in the commercial viability of potential organic semiconductors is having a charge-carrier mobility within 1-2 orders of magnitude of standard silicon devices. Polymorphism has the ability to transform a material into a commercially viable product because charge-carrier mobility is highly sensitive to polymorph-dependent structural details like $\pi-\pi$ stacking [27, 65]. The material 7,14-bis((trimethylsilyl)ethynyl)dibenzo[b,def]-chrysene (TMS-DBC) was recently shown to have a 75x larger hole mobility when adopting a 2D brickwork configuration rather than a 1D slipped stack configuration [29]. Similar examples of improved charge transport when adopting a new crystal structure have been observed in Rubrene [26], TIPS-pentacene [27], and C8-BTBT [28].

There are a number of other industries where adopting a particular crystal polymorph can improve the materials performance. The common explosive 1,3,5,7-tetranitro-1,3,5,7-tetrazacyclooctane (HMX) has four polymorphs α , β , γ , δ . All four forms can be detonated through impact. However, the β form is more stable and therefore less prone to accidental detonation [20, 21]. In the pharmaceutical industry, acetaminophen has a metastable form with significantly better tabletability than the globally stable form [33]. Disordered crystals like fenoprofen also have a number of promising uses for drug manufacturing due to their higher stability than amorphous crystals and their higher solubility/dissolution rates compared to crystalline structures [66, 67]. Finally, the hydrogel 4,6-O-Benzylidene--D-galactosyl azide becomes stronger after a spontaneous temperature induced polymorphic transformation from the α to β form [35, 68].

The case studies highlighted above demonstrate the commercial imperative of knowing all stable solid forms of a given material. A detailed understanding of the possible solid forms can ultimately avoid costly market recalls and/or improve a material's performance. In the next section, the process of experimentally searching for new solid forms is discussed.

1.3 The Experimental Search for All Solid Forms

1.3.1 Characterizing Distinct Solid Forms

In order to identify a new solid form, one must be able to discriminate between different polymorphs. The precise orientation, composition, and neighbor distances of synthesized crystal structures are most often characterized experimentally using single-crystal and powder x-ray diffraction [45, 49, 68, 69] as well as FTIR, RAMAN, and Stokes Ellipsometric spectroscopic analysis [32, 34, 45, 70, 71].

Solid forms can also be distinguished using thermogravimetric analysis [22, 45, 49], differential scanning calorimetry [14, 45, 48, 49, 53], as well as by examining the physical appearance such as color or crystalline shape [14, 25, 29, 49, 72]. Differential scanning calorimetry has the benefit of also providing key thermodynamic quantities such as heat capacities and enthalpies of transition between different forms. These thermodynamic quantities can be compared directly against computational estimates and validate model accuracy, which will be described more in Chapter 3.

1.3.2 What Factors Lead to Different Polymorphs?

When searching for alternate polymorph structures, it is useful to understand the factors which can stabilize a new solid form, either thermodynamically or kinetically. In any given thermodynamic state, there will always be one globally stable structure. However, metastable solid forms can be synthesized that are kinetically stable for days, months, or even years. Therefore, metastable forms must be considered as possibly observable structures in addition to the globally stable one. Indeed, Ostwald's famous step rule [73] states that when crystallizing a material from solution, the metastable phases will be discovered first before the thermodynamically most stable form.

There are many factors that can cause a solid to crystallize into a metastable conformation rather than the thermodynamically preferred structure. One common example is the choice of

solvent during the crystallization process. The cancer drug 5-fluorouracil is known to nucleate in a metastable form when crystallized in polar water. Computational and experimental studies later revealed that the thermodynamically preferred form is reached by synthesizing in non-polar solvent [53, 54, 74]. In this instance, polarity of the solvent kinetically biases the molecule into one specific molecular packing. Other environmental factors that can influence the synthesized crystal structure include temperature [26, 29, 44, 75], pressure [76–81], humidity [47, 56, 82, 83], impurities [84–86], impeller speed, and other shear forces [27, 52, 75, 85, 87].

A supersaturated solution can also be intentionally directed towards a metastable form by introducing a seed crystal [7, 85, 88], using a templated monolayer [38, 89], or using laser-induced nucleation [75]. Co-crystallizing agents are also well known to influence crystal structure formation [90], and it is estimated that roughly half of all commercially available drugs use salt as a co-crystallizer to control the synthesized crystal form [91].

The above paragraphs describe methods to trap a crystal in a metastable form rather than the globally stable structure. However, the kinetically stable forms will themselves depend on the specific thermodynamic state. Structures that can be isolated at one temperature may become kinetically unstable at higher temperatures. Furthermore, the thermodynamically most stable structure can also change with the temperature [18, 46, 92–103] or pressure [76–81] leading to new observable forms.

The complete set of all ‘observable polymorphs’ for a given material therefore consists of all structures that can be produced from any crystallization method in any thermodynamic state of practical interest to the material. The challenge for experimental polymorph screening procedures is to cover a large enough range within this multidimensional space to achieve high confidence that another undiscovered form, which could compromise the performance of the product, will not be observed. The following section details the current best practices in polymorphic screening procedures as well as their limitations.

1.3.3 Screening for New Solid Forms

Pharmaceutical researchers will typically search for alternate more-stable polymorph structures before mainstream production using a number of different screening techniques include slow cooling of the crystal in solution [29], slow evaporation of solvent [29], thermal annealing of known solid forms [29,75], as well as slurring experiments [20,49,50,75,84] and Polymer-induced Heteronucleation [104].

Each of the above approaches is designed to change the crystallization conditions such that a new solid form will emerge, if such a form exists. Generally the longer that a crystallization is allowed to occur, and the higher the temperature of crystallization, the more likely the material is to find the most stable solid form. By changing the time and/or temperature of the crystallization one can tune the synthesis toward metastable or globally stable forms and in theory find all possible crystallizable forms.

However, each of the crystal screening methods above introduces a new variable into the synthesis process. Any combination of time, temperature, pressure, humidity, solvent, and stirring rate describes a unique crystallization condition that could, in theory, lead to a new unknown solid form. Given the vast space of possible synthesis conditions, it is understandable that unexpected transition into more stable crystal structures continues to plague commercial products, perhaps most famously illustrated by the market recall of the HIV drug ritonavir described earlier. In the particularly striking case of thymine, slurring experiments in 27 organic solvents all yielded the known forms [105,106]. However, subsequent sublimation experiments almost a decade later produced a new stable polymorph, and further dehydration experiments found two more previously unknown forms [106]. Examples like these demonstrate that new stable forms can be missed even after thorough polymorphic screenings have already been conducted.

New unknown forms could, in theory, be identified early in the drug development process by using the precise combination of crystallization conditions leading to the new form. However, there is no practical method to rigorously test all crystallization procedures even with

state-of-the-art high-throughout screening techniques. This dilemma of too many possibilities directly motivates the use of cheaper and faster computer models that can help predict all possible solid forms of a material.

1.4 Computer-Aided Crystal Structure Prediction

Computational crystal structure prediction (CSP) methods represent an alternative approach to traditional experimental screenings for finding stable crystal forms. Computer-based structure prediction models calculate the thermodynamic stability of crystals directly and thus avoid many of the issues in experimental polymorph screenings such as nucleation source, choice of solvent, and high kinetic barriers between forms.

Computer-based structure prediction models can, in theory, be significantly less expensive than full experimental screening of polymorphs, which costs roughly \$40,000-\$50,000 for a comprehensive screening of a single compound [107]. This same money could be used to purchase 1,250,000 CPU-hours on publicly available cloud computing servers [108] translating to thousands of independent molecular simulations. Computer-based polymorph screening methods will therefore be a cost-effective approach if accurate *in silico* structure predictions or eliminating unnecessary experiments can be achieved using less than this target amount of computing time.

This raises the question: are current computer models accurate enough to reliably predict correct crystal structures? The current state of computational crystal structure predictions in finding real structures is measured every 3-4 years through blind crystal prediction challenges. As of 2016, there have been six blind challenges for predicting crystal structures [109–114]. In these blind challenges, entrants are asked to predict the crystal structures of a numbers of molecules which have not had experimentally published structures to date. The participants are provided with only a two-dimensional schematic of the target systems as well as the crystallization conditions. Each participant in the challenge can use any computational methods necessary to predict the correct structure and are allowed to submit 3 proposed structures for

each molecule. A more detailed review of the procedure for generating predicted structures is provided in Chapter 2. At the end of the challenge, the submitted structures are compared against the true experimentally synthesized crystals.

The rate of successful structure prediction has increased in each subsequent blind challenge due to a combination of advances in sampling the full crystal landscape [112], better energy functions [115], and treatment of flexible degrees of freedom [114]. Generally, one group will bring forward a particularly effective new strategy which will then be adopted by all other groups in future challenges. This directed evolution of CSP methods has led to considerably reliable predictions in recent challenges. Indeed, in the last three blind challenges all but one of the target systems had the correct structure found by at least one submission [112–114].

The recent blind crystal prediction studies demonstrate that current computational methods are capable of identifying the true experimentally preferred structure. Additionally, computer-aided crystal structure prediction methods can in some cases identify stable structures that were not otherwise found from experimental polymorph screenings. A new stable form was found computationally before experimental synthesis in aspirin [16, 116], thymine [106], 5-Fluorouracil [53], carbamazepine [89], cyheptamide [117], and malic acid. In the case of 5-Fluorouracil, the computer models also provided the key insight that synthesis in *aprotic* solvent leads to the globally stable form, explaining why previous experiments in protic solvent had only crystallized the metastable known form [54].

Ultimately, the recent blind crystal prediction studies demonstrate that current computational methods are capable of identifying the true experimentally preferred structure. Furthermore, the insight from computer models can occasionally guide the experimental synthesis procedure to crystallize the new stable solid forms [54, 89, 106]. These computational successes reveal the potential for computer-aided crystal structure prediction in screening for all solid forms of a given material.

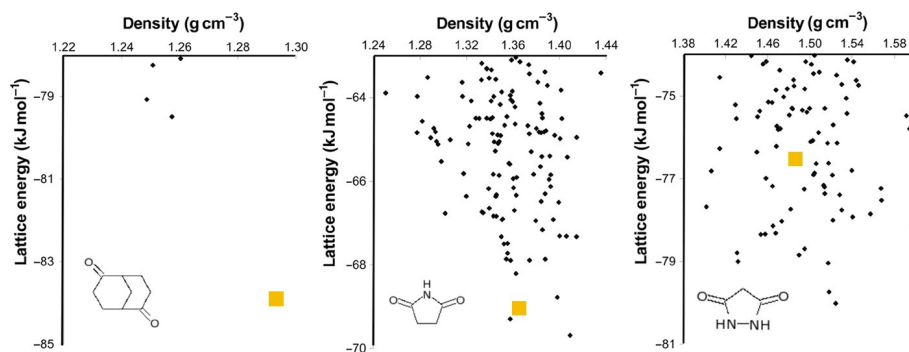


Figure 1.2: The crystal structure prediction of (a) bicyclo(3.3.1)nonane-2,6-dione finds the experimental form as the clear unambiguous most stable minima. The crystal prediction of (b) succinimide finds the experimental structure as one of many low energy lattice minima. The crystal prediction for (c) 3,5-pyrazolidinedione has numerous structures more stable than the experimental form. All three CSP are reproduced from ref 11.

1.5 Why Do Computer Models Find So Many Possible Structures?

The observation that computer-aided crystal structure prediction models generally find the true experimental structure initially gives optimism that *ab initio* crystal structure prediction is close. However, in many cases these models find tens or hundreds of other structures in addition to the true structure with similar relative stabilities.

Figure 1.2 depicts three real crystal structure prediction results reproduced from ref 11. Each point in Figure 1.2 represents a lattice minima identified by the model, and the true experimental structure is highlighted in orange. The example of bicyclo(3.3.1)nonane-2,6-dione in Figure 1.2a represents an ideal result from a crystal structure prediction study, with the experimental structure identified as the clear lowest energy form. There are multiple other candidate lattice minima identified for this molecule which are not observed. However, the lack of experimental observation for these structures could be explained away by a significantly lower stability, and once the globally stable polymorph is synthesized it is unlikely to revert into another metastable form.

If every crystal structure prediction study resulted in a figure similar to that of bicyclo(3.3.1)nonane-2,6-dione, computer-aided crystal structure prediction would be a mature science that was well

incorporated into the commercial material synthesis pipeline. However, the more common result from a CSP is that the experimental structure is identified as one of multiple low energy structures (as in Figure 1.2b) or as a significantly metastable structure with many lower energy forms (as in Figure 1.2c).

The existence of so many other predicted structures in these cases significantly inhibits the ability of the model to direct experimental synthesis of new forms. It's entirely plausible that one or more of these predicted structures does correspond with a real observable form that could improve (or compromise) a commercial material. However it is difficult to determine which predicted more stable form in this space is worth searching for with targeted experimental synthesis procedures. An urgent need in the computational crystallographic community is to find ways to eliminate the dead-end lattice minima that do not correspond with structures which will eventually be observed experimentally.

There are three primary explanations for why a predicted structure from standard CSP approaches may ultimately be unrealizable in the physical world.

1. Temperature-mediated Stability Reranking

Large entropy differences between crystal structures could explain why certain predicted structures are not favorable under working conditions. Current crystal structure prediction methods identify all minima on the *lattice energy* surface as a potentially observable structure. However, experimental structures are minima on the *free energy* surface at ambient conditions. Structures that have low minimized energy may nevertheless have significantly lower entropies which would increase the free energy at room temperature relative to more entropically favorable forms. Said another way, if both energy *and* entropy were used to rank the relative stabilities of each minima, the ambiguous CSP results like those in Figure 1.2 b and c may become as clear as that of Figure 1.2a.

2. Sensitivity to the Hamiltonian

Subtle intermolecular interactions that are neglected in the energy model can also explain why predicted structures are not as favorable as initially estimated. CSP studies are often

done with off-the-shelf classical potentials, or with classical potentials fitted to select quantum-mechanical (QM) calculations [118]. The studies are rarely done with full *ab initio* calculations at every step of the process. Including more advanced effects such as anisotropic multipoles, intermolecular polarization, and QM electron correlations can significantly change the lattice energy landscape and therefore the relative ranking of different predicted structures. In some cases, the lattice minima in a cheaper potential can completely restructure when more accurate interaction effects are included [119]. It is therefore possible that with a sufficiently accurate potential, all crystal structure prediction studies would appear like that of bicyclo(3.3.1)nonane-2,6-dione with a single clear prediction of the experimental structure and all unobservable structures being far less energetically stable.

3. Lattice Minima Interconversion at Ambient Temperatures

A final hypothesis for why there are more predicted structures than experimental forms is that many of the lattice energy minima rapidly (inter)convert to the same ensemble when thermal motions are present. In the limit of zero Kelvin, all local minima on the lattice energy surface represent stable and observable structures. However, at higher temperatures, many lattice minima will have low kinetic barriers to escape into other nearby minima. In this hypothesis, many of the predicted minima on the lattice energy surface will actually lead to the same ambient temperature experimental structure once thermal motions are added to the system. In such cases, the plethora of close lattice minima seen at 0K (such as in 1.2 b and c) would ultimately appear as just a small set of unique forms at room temperature.

Ultimately, each of the three hypotheses above can be probed with the appropriate use of molecular simulation. However, each requires a significant improvement on the current state-of-the-art in how organic crystals are modeled computationally. Specifically, each hypothesis can be explored through the use of molecular dynamics (MD) simulations. In chapter 2, the current models for crystal structure prediction are reviewed in detail. Molecular dynamics is

then introduced as a general simulation method followed by the specific details of how MD is applied to organic crystal polymorphs. Chapters 3, 4, and 5 are dedicated to exploring each of the main hypotheses for lattice minima elimination including temperature-mediated reranking, corrections for more accurate potentials, and thermal interconversion at ambient temperatures.

Chapter 2

Methods for Modeling Crystal Polymorphs

2.1 Searching the Crystal Landscape for Candidate Structures

2.1.1 Generating Candidate Structures

The task of computationally predicting crystal structures requires generating a set of candidate structures followed by estimating the relative free energy to identify the most stable packing configurations. A full review of crystal structure prediction methodology can be found in ref 11. The general process is briefly summarized here. The starting point of a crystal structure prediction is generally a 2-dimensional atomic connectivity chart of the molecule of interest. The 2-dimensional schematic is transformed into a 3-dimensional molecule through vacuum phase geometry optimization, either classically or quantum mechanically. This 3-dimensional molecule and multiple replicas are then packed together into a number of different possible unit cells. Finally, all of these possible unit cells are evaluated to determine the lowest energy and highest density structures. A schematic of this process with all three steps described above is shown in Figure 2.1 taken from ref 11.

Once an isolated molecular geometry is determined, one could in theory scan through all possible unit cells and find the absolute most favorable packing configuration. However, every lattice has at least 6 degrees of freedom for the unit cell (the three cell lengths a , b , c and three

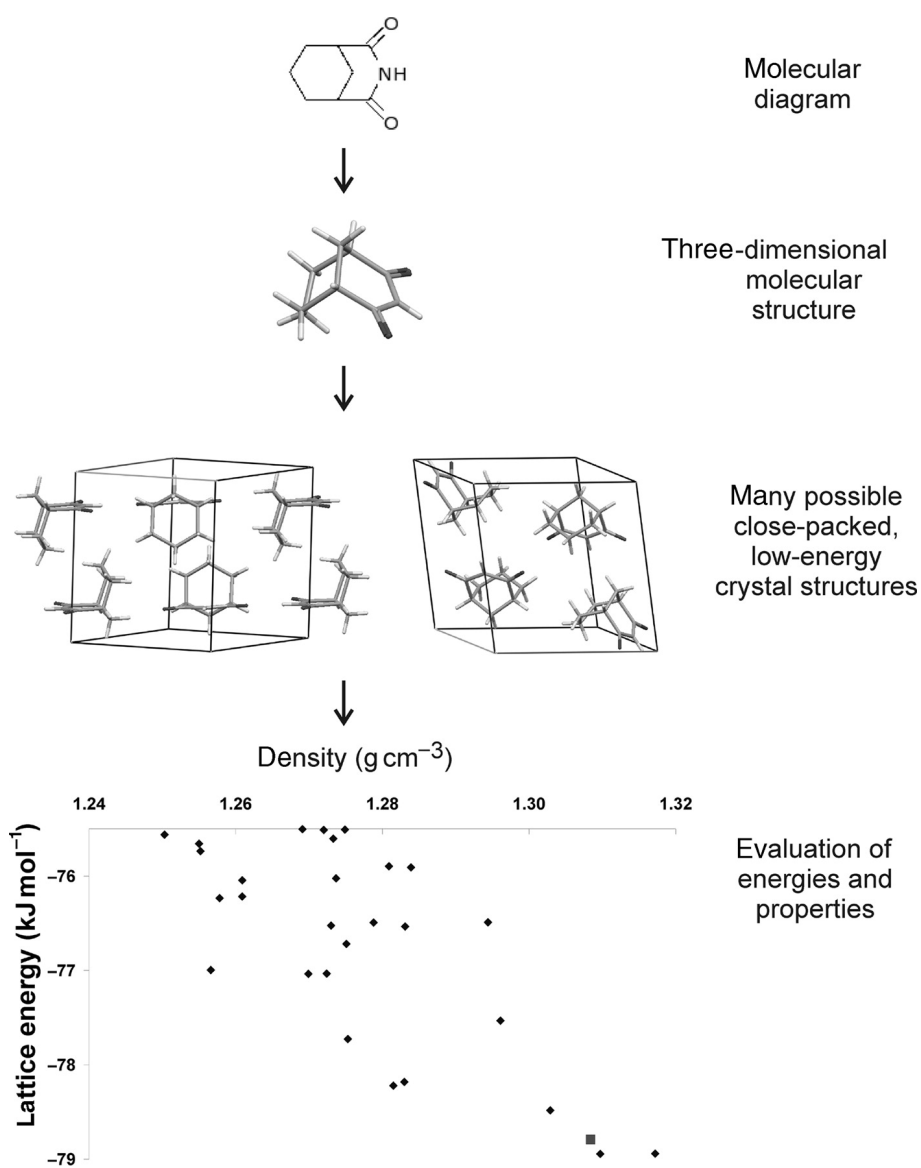


Figure 2.1: Workflow depicting how lattice energy-based crystal structure predictions are carried out (reproduce from ref 11).

angles α , β , γ) plus 3 additional degrees of freedom for every atom of every molecule in the system. This astronomically large number of degrees of freedom cannot be rigorously searched with any practical brute force technique. Furthermore, an overwhelming majority of structures produced from a blind scan through all degrees of freedom will involve obviously unfavorable conformations with overlapping atoms or unphysically distorted molecules. For this reason, a number of simplifying assumptions are usually employed to facilitate more efficient unit cell generation.

First, the molecules are typically held rigid at the internal geometry from the vacuum optimized structure. The justification for this simplification is that intramolecular interactions are typically much stronger than intermolecular interactions, and will force a molecule to maintain roughly its vacuum-phase geometry even in a packed crystal [120,121]. However, one should take care to recognize that dihedral rotations can vary significantly between vacuum and crystallized molecules. As an extreme example, biphenyl slightly favors a 45 degree out-of-plane tilt in an isolated environment [122]. However, the molecule adopts a planar orientation in its room temperature crystal structure to reduce the free volume [123]. Thus, some chemical intuition is generally required to determine the correct rigid molecule orientation to use in the remainder of the search.

The second simplifying assumption is to restrict the lattice parameters to one of the 230 crystallographic space groups. Restricting the cell to a spacegroup symmetry significantly reduces the number of possible cells to scan through. For example, an unconstrained cell with 4 rigid molecules has 30 degrees of freedom. However, a cell constrained to the P21/c space-group with $Z'=1$ has only 10 degrees of freedom [11]. An obvious drawback of this approach is that crystals adopting an amorphous or disordered structure that cannot be classified in any individual space group will be missed. However, roughly 80% of homomolecular crystals are classified in one of the six most common space groups [11].

The final simplification is assuming that the asymmetric unit cell contains only a small number of independent molecules in the cell. It is estimated that 90% of crystals have one

or less molecules in the asymmetric unit cell ($Z' \leq 1$) [124]. However, one can in theory include any number independent molecules above one in the structure generation process. The number of independent molecules is limited only by the computational power, but is often restricted to single digits.

Once these simplifications are employed, reasonable candidate crystal structures can be produced through brute force approaches such as sampling a grid of all degrees of freedom or using random Monte Carlo methods. A number of different programs exist for generating candidate structures in this manner including Crystal Structure Predictor [125, 126], GRACE [118, 127], MGAC [128–130], and USPEX [131, 132] among many others [133–136]. Each of these programs have subtle differences in the implementation and parameters of the search process. However, the end result from all of these programs is the same: a collection of plausible crystal configurations. From this point, the lattice energies of these configurations can be directly evaluated and compared. Typically, the raw configurations from a crystal structure generation algorithm are further refined through a final lattice energy minimization step before being assessed for stability.

2.1.2 Refining Structures through Lattice Energy Minimization

The goal of an energy minimization process is to find the precise local energy minima closest to an initial configuration. This essentially condenses to a mathematical problem of minimizing a function (the potential energy) by varying the spatial positions of all the atoms in the system. There are numerous well-established methods to minimize a multivariate function including steepest descent, conjugate-gradient, and Newton-Rhapson.

In the context of a crystal energy minimization, the lattice parameters and angles must also be varied to find the true lattice energy minimum. The minimization of the box vectors and minimization of the atoms within the box are generally done in alternating fashion until a sufficient tolerance is reached for both. There are a number of possible methods for implementing this crystal energy minimization process. The two crystal minimizations used in this

work are the XTALMIN subroutine within the TINKER [137] simulation package as well as a self-developed crystal minimization scheme in GROMACS [138].

The XTALMIN subroutine takes as input the structure coordinates and unit cell lattice parameters. It then alternates cycles of Newton-style optimization of the structure and conjugate gradient optimization of the crystal lattice parameters [137]. The crystal minimization process in GROMACS also takes a structure file as input. It then alternates minimizing the energies within a cell using the BFGS algorithm [139] to a desired force tolerance along with a Nelder-Mead minimization algorithm (as implemented in numpy version 1.8.2) using the box vectors as optimization variables to find the minimum energy crystal structure.

The above minimization processes can be done with rigid/constrained molecules or with fully flexible molecules. In many cases, the energy minimization is done with fully flexible molecules even if the structures were generated with rigid bonds. Allowing the internal molecular geometry to relax in the final energy minimization step alleviates a portion of the error associated with generating candidate structures using a rigid molecule.

Once the lattice energy minimization is complete, the resulting structures can be compared for relative stability. Typically, the stability ranking is presented through a plot of lattice energy against density for each minima, similar to that shown in Figure 1.2. The lattice energy landscape will be highly dependent on the energy function used to score candidate structures. There are numerous methods for modeling intramolecular and intermolecular interactions in molecular crystals through classical or quantum mechanics based approaches. In classical Hamiltonians, particular emphasis is focused on describing the Coulombic interactions between the atoms and molecules. More details of possible Hamiltonians, differences in functional forms, and the sensitivity of stabilities to more accurate potentials will be discussed in more detail in Chapter 4.

2.1.3 The Need for Ambient Temperature Structures

Up to this point in the discussion, crystal structures have been represented as a single low energy configuration. However, the depiction of organic crystals as static lattices will only rigorously apply in the limit of zero Kelvin, where (classically) no atomic motions are present and energies fully determine conformation stability. In real finite-temperature crystals, molecules vibrate around their crystalline positions and the *free energy* of the resulting ensemble determines the stability of the overall crystal structure. In previous crystal structure prediction studies based on lattice energies where the model failed to predict the experimentally known structure, the authors frequently attributed the failure to the neglect of the thermal motions and entropic contributions present in real crystals [111,140,141].

As many as 50 possible configurations can exist within 5kJ/mol of the globally stable structure in lattice energy based structure prediction studies [11]. Furthermore, roughly half of known polymorphic pairs have a lattice energy difference less than 2kJ/mol [142]. Entropy could influence the stability ranking of the polymorphs if the temperature and entropy contributions to the free energy exceed this value. A review by Gavezzoti showed that real crystals can have entropy differences between polymorphs as large as 15 J/molK which corresponds to a free energy differences of 4.5 kJ/mol at room temperature. A recent study by Nyman and Day [142] also estimated relative entropic contributions between polymorphic pairs using a harmonic approximation in 1001 systems and identified temperature-mediated transformations in the range 0 K - 300 K in 9% of the systems examined. These studies clearly demonstrate that entropy effects can lead to different crystal stability rankings at 0 K and at ambient conditions in real systems.

A common physical manifestation of the effects of entropy in the crystal landscape is the formation of disordered crystal phases which appear in nitrous oxide [143], caffeine [72], and cyclopentane [100] and many other systems. In these disordered phases, the observed crystal structure is energetically less stable than the global lattice minima but observed nevertheless because of entropic favorability.

The inability of static lattice models to faithfully represent true crystal structures has led investigators in recent years to begin generating full configuration ensembles in polymorph prediction studies using free energy simulation methods [144–154]. The free energy of a crystal structure can be estimated in molecular simulations with either a harmonic approximation, or with a fully atomistic Monte Carlo or molecular dynamics sampling of the phase space around each lattice energy minima. Both of these free energy estimation methods are described in the next sections. The differences between free energy estimates with QHA and with MD are discussed for a number of real systems in Chapter 3.

The free energies of real polymorphs can also be determined experimentally using various methods include the heat capacity method [46], the solubility method [155], and the eutectic melting method [46]. Relative entropies between of polymorphs can be computed by measuring the temperature and enthalpy change at the transition between forms with DSC [93, 98, 101, 102]. Comparison between estimated and measured entropies and free energies are presented in the systems in Chapter 3. Relative enthalpies and free energies can also be validated qualitatively experimentally by examining how the phase diagram changes with temperature and pressure.

2.2 Statistical Mechanics Formalism

The terms and formalism for the statistical mechanics approach to modeling molecular systems is briefly reviewed here before introducing both molecular dynamics and the quasi-harmonic approximation. In statistical mechanics theory, any equilibrium property of a macroscopic systems is determined by an average over the collection of all possible microstates. Let A_i be the value of some system observable in a given microstate i . The expected value of the observable $\langle A \rangle$ is computed from the standard formula for expectations from probability theory:

$$\langle A \rangle = \sum_1^N A_i p_i \quad (2.1)$$

where p_i is the probability of observing the system in microstate i , and N is the total number of possible microstates. To calculate any observable of a system, we need an expression for the probability of each microstate. The complete delineation of every possible microstate along with a corresponding probability is called a ‘thermodynamic ensemble’ [156].

The fundamental postulate at the foundation of statistical mechanics is that all microstates with the same energy have an equal probability of being explored by the system (aptly called the *equal a priori postulate* [157]). Thus for a totally isolated system with no exchange of mass, volume, or energy with the external environment, all microstates are equally likely and the system observable becomes a linear average over all possible microstates $\langle A \rangle = \frac{1}{N} \sum_1^N A_i$.

In most applications of scientific interest, the system is in equilibrium with the external surroundings at a constant temperature, T . The system can now, in principle, access any energy level, U_{sys} , by drawing energy from the surrounding environment. However, the surroundings will gain or lose an equal amount of energy to conserve the total energy of the universe. When the surroundings loses an amount of energy ΔU (and the system goes to energy $U + \Delta U$ the entropy of the surroundings will decrease, indicating that the microstate with system energy $U + \Delta U$ is less likely than one with just the energy U . As a result, the possible microstates no longer have an equal probability of occurring when they have different energy levels. So how does one represent the probability of a microstate i with energy level U ?

In order for the system to randomly fluctuate at equilibrium from an energy of U_{sys} up to an energy of $U_{sys} + \Delta U$, the surroundings must change from a microstate with energy U_{sur} to one with energy $U_{sur} - \Delta U$. By positing that each microstate of the surroundings is equally likely, the probability of the system being at an energy level U_{sys} relative to energy level $U_{sys} + \Delta U$ corresponds with the ratio of total microstates that the surroundings can be at energy levels

U_{sur} and $U_{sur} - \Delta U$.

$$\frac{p_2}{p_1} = \frac{W_2}{W_1} \quad (2.2)$$

where p_2 and p_1 are the probabilities of the system having energy U_{sys} and $U_{sys} + \Delta U$, respectively, and W_2 and W_1 are the total number of microstates of the surroundings with energy $U_{sur} - \Delta U$ and U_{sur} , respectively. The total number of microstates of the surroundings as a function of energy is directly related to the temperature and the entropy. The statistical mechanics definition of temperature is $T = \frac{dU}{dS}$. Thus, if the surroundings is at a temperature T , transferring an energy ΔU will come at the expense of $\frac{\Delta U}{T}$ less entropy in the surroundings. Boltzmann famously showed in the late 1800s that the entropy of a system is related to the number of indistinguishable microstates through the equation

$$S = k_B \ln W \quad (2.3)$$

where W is the total number of microstates, and k_B is the Boltzmann constant. Thus, the change in total microstates after losing a quantity of energy ΔU at a temperature T can be expressed as:

$$S_2 - S_1 = k_B \ln W_2 - k_B \ln W_1 = k_B \ln \frac{W_2}{W_1} = \frac{-\Delta U_{12}}{T} \quad (2.4)$$

which rearranges to

$$\frac{W_2}{W_1} = \frac{p_2}{p_1} = e^{\frac{-\Delta U}{k_B T}} \quad (2.5)$$

This remarkable result in equation 2.5 demonstrates that the probability of the system being in any microstate is proportional to the factor $e^{\frac{-U_{sys}}{k_B T}}$ at constant temperature. When the system is allowed to also exchange volume and particles with the surroundings, this factor is generalized to $e^{\frac{-U+PV+\sum_1^M \mu_i N_i}{k_B T}}$, where P is the pressure of the surroundings and μ_i is the chemical potential of species i in the surroundings.

In light of this result, the concepts of statistical mechanics in any system in any ensemble condenses to the following simple idea. The probability of any system microstate at equilibrium is determined by the amount of total microstates lost by the surroundings in order to achieve the necessary energy, volume, and number of particles. Each of these three quantities energy, volume, and particles can be increased subject to a microstate penalty to the surroundings precisely dictated by $\frac{1}{k_B T}$, $\frac{P}{k_B T}$, and $\frac{\mu_i}{T}$. The probability of any microstate i is then proportional to $p_i \propto e^{\frac{-U+PV+\sum_1^M \mu_i N_i}{k_B T}}$. Combining this powerful result with the equation 2.1 earlier gives the expected value of any system observable, typically expressed as an integral over the whole ensemble:

$$\langle A \rangle = \int_X A(x) Z^{-1} e^{\frac{-U+PV+\sum_1^M \mu_i N_i}{k_B T}} \quad (2.6)$$

where X is the full phase space, and Z is a normalization constant called the ‘Partition function’. The free energy, F , of the system at any thermodynamic state is related to the partition function through $F = -k_B T \ln Z$. In the NVE, NVT and NPT ensembles, the ‘free energy’ quantity that is minimized at equilibrium is the entropy ($-S$), the Helmholtz free energy (A), and the Gibbs free energy (G), respectively.

The goal of nearly all statistical mechanics molecular modeling efforts can be summarized as estimating the free energy difference between two or more thermodynamic states. Estimating the free energy difference is, at bottom, the process of estimating the ratio of the partition functions between the two states. A true computation of the free energy and partition function would require knowing all possible microstates for a given system. This calculation is

not feasible in practice because the number of possible microstates is intractably large in most quantum systems and essentially infinite in classical systems.

Current molecular modeling techniques instead use numerical and analytical approaches to estimate the ratio of the partition functions without computing either partition function directly. These approaches involve randomly sampling configurations in such a way that they appear with a probability $p_i \propto e^{\frac{-U + PV + \sum_1^M \mu_i N_i}{k_B T}}$. After a sufficiently large amount of sampling in the two states, the ratio of the partition functions (and thus the free energy difference) can be estimated. There are a number of different ways to implement this free energy estimation either numerically with Monte Carlo and molecular dynamics simulations or analytically with a quasi-harmonic approximation. The latter two methods will be developed in the following sections.

2.3 Estimating Entropy with a Quasi-Harmonic Approximation

2.3.1 The Harmonic Approximation

Computational models that include entropy effects typically estimate entropy differences using a harmonic approximation. [92, 142, 158, 159]. In this approximation, all motions in the crystal are assumed to be harmonic about the energy minimized structure. Any crystal configurations that require an anharmonic translation of the atoms away from the minimized structure are assumed to contribute negligibly to the overall ensemble. The simplification of harmonic motions significantly reduces the amount of computational overhead needed to compute free energy differences between polymorphs.

The classical free energy of a crystal polymorph for a given configuration with the harmonic approximation is:

$$A(T) = \sum_k \frac{1}{\beta} \ln(\beta \hbar \omega_k) \quad (2.7)$$

where $\beta = (k_B T)^{-1}$ where k_B is the Boltzmann constant and T is temperature, \hbar is the reduced Planck's constant, and ω_k is the k th angular frequency calculated from the mass weighted Hessian of the local lattice minimum at which the harmonic approximation is performed.

2.3.2 The Quasi-Harmonic Approximation

The harmonic approximation described above computes the Helmholtz free energy of a polymorph from a given configuration at a specific volume. However, at finite temperatures and pressures, crystal lattices expand to volumes larger than the 0 K minimized structure. The harmonic approximation is frequently extended to include volume changes in order to account for the thermal expansion of real crystals. These approaches are referred to as the 'quasi-harmonic approximation'. The following assumptions are used in the quasi-harmonic formalism adopted here: 1) all vibrations are treated as harmonic oscillators, 2) thermal expansion behaves isotropically, and 3) wavenumbers are lattice volume dependent only [160]. The free energy of a crystal structure with the quasi-harmonic approximation is computed by:

$$A_v(V, T) = \sum_k \frac{1}{\beta} \ln(\beta \hbar \omega_k(V)) \quad (2.8)$$

$$G(T) = U(V) + A_v(V, T) + PV \quad (2.9)$$

The Gibbs free energy from equation 2.9 is computed at the volume where the free energy is minimized at the temperature, T .

2.3.3 Limitations of the Quasi-Harmonic Approximation

The quasi-harmonic approximation improves upon the harmonic approximation by including changes in the entropy due to thermal expansion [161]. However, both HA and QHA exclude any anharmonic motions in the system of interest by the very nature of the approach. Full ambient temperature entropy differences, including anharmonic motions within a given

crystallographic cell, have not yet been directly computationally reported for organic crystals, and may not agree with the values estimated from these harmonic approximations for typical organic crystals. This approximation also suffers from the known problem that all energy minima will always be identified as also corresponding to a free energy minima. The harmonic approximation method cannot therefore determine when two energy minima correspond with configurations from the same polymorph ensemble [31].

Molecular dynamics simulations can in principle capture all entropic contributions to the ambient temperature stability of different crystal forms including those from anharmonic motion, lattice expansion, intramolecular fluctuations, and temperature-dependent vibrational frequencies. In the next section the relevant concepts of molecular dynamics are introduced along with their connection to modeling crystal polymorphs.

2.4 Introduction to Molecular Dynamics

2.4.1 Sampling the Configuration Space

Molecular dynamics (MD) has emerged in recent years as a popular method for sampling ensembles of configurations and computing free energy differences in molecular systems. These simulations will, in theory, sample the complete phase space including the anharmonic configurations neglected by the harmonic approximation. MD simulations produce configurations by numerically solving Newton’s equations of motion for a system of N particles interacting through a potential energy function $U(x_1(t), x_2(t) \dots x_N(t))$, where x_i is the spatial positions of particle i at some time t . The ‘particles’ in computational chemistry models are the atoms in the system, and the potential energy function captures the inter- and intra- molecular interactions between all the atoms in the system.

The forces on each particle i at time t are determined by the derivative of the potential energy with respect to the particles position.

$$F_i(t) = -\frac{dF_i}{dt} = -\frac{dU}{dx_i} \quad (2.10)$$

Once the forces of all particles are determined, the position and velocities of particles as a function of time are determined through a numerical integration of Newton's equation:

$$F_i(t) = -m_i \frac{d^2 x_i}{dt^2} \quad (2.11)$$

The end product of a molecular dynamics simulation is a string of configuration snapshots sampled at periodic intervals throughout the trajectory. This set of configurations provides a decent representation of the full ensemble as long as the system has reached equilibrium and a sufficiently large sample size of configurations has been collected. Molecular dynamics simulations can be viewed as essentially a Monte Carlo walk through a configuration space in the limit of a 100% acceptance rate and highly correlated adjacent samples.

Direct Monte Carlo sampling of a configuration space can also be done, where a new proposed configuration is selected at random and either accepted or rejected based on the change in energy. This Monte Carlo sampling is very efficient for systems with few correlated motions, such as an isolated vacuum phase molecule. However, in solution phase and solid phase chemistry, adjacent molecules will generally have highly correlated motions in which two molecules move together due to the close-packed / low free volume nature of the phase. In such systems, Monte Carlo methods can have an exceedingly low acceptance rate which will make the method prohibitively expensive. Molecular dynamics simulations seamlessly facilitate collective motions between molecules because the intermolecular forces of the closely packed molecules will lead to collective motion when Newton's equation is numerically integrated. Molecular dynamics therefore represents a useful method for sampling the complete configuration space of crystal polymorphs including the configurations neglected by a harmonic approximation.

A number of independent molecular dynamics simulation packages are available to facilitate sampling including GROMACS [138], TINKER [137], NAMD [162], CHARMM [163], AMBER [164], LAMMPS [165] and many more. In this work, all molecular dynamics simulations are carried out with either GROMACS or TINKER.

2.4.2 Periodic Boundaries and Long-range Interactions

Real crystals in the physical world are essentially infinite in extent from the perspective of an individual atom at the nm length-scale. However, a direct simulation of a full lab-scale crystal would require a MD simulation of $\sim 10^{23}$ atoms. This vast quantity of atoms cannot be tracked within the reasonable limits of current computer memory. Instead, molecular models of crystals make the assumption that the full system can be reasonably approximated by modeling a smaller subsystem surrounded by identical copies. A smaller crystal on the order of $\sim 10^3$ atoms is placed within a small period box. This box is surrounded on all sides by an identical copy of the box such that each atom in the original box has a corresponding atom in every surrounding box which behaves identically and moves identically to the original atom. This system is now said to have ‘periodic boundary conditions’. With this approximation, only the positions and momenta of the $\sim 10^3$ atoms in the main box need to be tracked throughout the course of the simulation.

Particles in periodic box copies will still interact with those in the main system, even though the motions are not independent. The periodic particles interact through the long-range van der Waals and Coulombic interactions. Thus, a sufficient number of box copies must be included in order to have a realistic estimate of the long-range interactions. Strictly speaking, there is no true cut-off for intermolecular interactions. All particles in the observable universe are interacting with all other particles at every moment in time. However, both types of intermolecular interactions decay in strength as the distance between particles increases (exact functional forms for intermolecular interactions in this work are presented in Chapter 4).

One approach to treating long-range interactions is to introduce an explicit interaction cut-off, wherein all interactions beyond a certain distance are ignored as negligible. This type of hard-cutoff is generally discouraged because a very large cutoff is usually needed before the interactions truly become negligible, particularly for Coulombic interactions. The van der Waals interactions tend to decrease as r^{-6} , which decreases rapidly with distance. However,

the Coulombic interactions decrease as r^{-1} and do not become negligible with any practical hard cutoff value.

In this work, the long range interactions are accounted for using Particle Mesh Ewald (PME) summation [166]. Ewald summation [167] partitions intermolecular interactions into both a ‘short-range’ and ‘long-range’ component. The short-range interactions are calculated directly by summing the interactions in each of the adjacent periodic boxes. The long-range interactions are calculated in reciprocal space using Fourier transforms. This summation in reciprocal space converges faster than the direct sum in real space and thus requires less computational overhead. For increased efficiency, Ewald summation is generally implemented using a smooth Particle Mesh Ewald (SPME) approach [168]. In SPME, the charges in the system are concentrated onto a grid of points using spline interpolation. These fixed grid points are then summed using Ewald summation. SPME is in practice substantially faster than direct Ewald summation.

Particle Mesh Ewald is used in this work to compute long-range van der Waals interactions in addition to long-range Coulombic interactions. PME is a commonly used method for long-range Coulombic interactions. However, van der Waals interactions can also be computed with the same PME approach. This works particularly well for crystalline solid systems because the bulk structure is already quite periodic and so PME converges quickly with only a few k-space points.

2.4.3 Temperature and Pressure Coupling

The ‘natural’ ensemble for a molecular dynamics simulation is the constant volume, constant energy (NVE) ensemble. MD simulations operate by placing molecules into a rigid periodic box and integrating the equations of motion as described earlier. As long as the timestep for the numerical integrator is relatively small, the system will maintain a constant energy in a constant box for the entire duration of the simulation. However, most systems of scientific

interest operate in either the NVT or NPT ensembles. In order to simulate in these ensembles, the dynamic simulations are augmented with auxiliary subroutines to mimic the effects of interacting with a surrounding environment at a temperature, T , and/or pressure, P .

The routines in MD for mimicking coupling to a constant temperature surrounding are called thermostats. The general idea behind these temperature coupling algorithms is that the velocities of particles in the system are periodically adjusted to ensure that a proper Maxwell-Boltzmann distribution of velocities is maintained throughout the simulation. This essentially substitutes for the process of particles colliding with the environment and attaining new velocities, which would occur in a real physical system. One simple approach to maintain this velocity distribution is to assign random new velocities to all particles in the system drawn from the Maxwell-Boltzmann distribution at periodic intervals throughout the simulation run-time. This technique is known as the Andersen thermostat [169]. There are also a number of non-stochastic methods to maintain the correct velocity distribution by introducing a friction term to the equations of motion or through periodically scaling the velocities of all particles [170–172]. Finally, one can ensure the correct distribution of velocities by using a stochastic integrator using the Langevin equations of motion to add a friction and random noise term when integrating the equations of motion [173].

The routines for coupling an MD simulation to a constant pressure surrounding are referred to as barostats. These pressure coupling algorithms act in a similar manner to the above temperature coupling methods except that they adjust the box volume to give a volume distribution consistent with the system's Helmholtz free energy as a function of volume, $A(V)$. Non-stochastic barostats include the Parrinello-Rahman barostat [174] and the Martyna-Tuckerman-Tobias-Klein (MTTK) barostat [175]. Briefly, the box vectors are given their own equations of motion which are integrated through time. The net effect is coupling the box to a piston, which checks the box's internal pressure (as measured by the virial) against a reference external pressure and adjusts the box dimensions accordingly. Pressure coupling can also be done stochastically by Monte Carlo sampling of the box vectors [176].

It's important to note that crystals are anisotropic, meaning that the system behaves differently depending on the direction. A visible manifestation of anisotropy appears in the thermal expansion of materials, wherein one dimension may expand more than another when a solid is heated [177,178]. For this reason, pressure coupling algorithms for solids must be anisotropic, such that the correct distribution of both the box volume and the box shape are sampled.

One final note on temperature and pressure coupling algorithms is that a truly accurate thermostat and barostat that mimics coupling to an external environment should produce both correct averages *and* correct fluctuations around the average. Some temperature and pressure coupling algorithms like the Berendsen algorithm [170] give rise to correct averaged structures (average energy/ average volume), but not the correct fluctuations around the average. These fluctuations are necessary for computing the correct free energy differences between thermodynamic states. Thus, weak coupling algorithms like Berendsen can be used for initial equilibration, but should not be used for full production simulations when free energy differences are desired. The accuracy of thermostats and barostats in producing the correct distributions around the average can be checked using the simple approach presented in ref 179.

2.5 Computing Crystal Free Energies

2.5.1 Generating Crystal Ensembles with Molecular Dynamics

Molecular dynamics simulations approximate the crystal ensemble by generating configurations with occurrences proportional to their probability in the ensemble. A molecular dynamics simulation starts from an initial crystal structure and then bring the system to a given temperature and pressure and solve newtons equations of motion as described in section 2.4. The resulting configurations in the trajectory are saved at regular intervals and used to approximate the crystal ensemble for that polymorph at that particular temperature and pressure. The difference in free energy between two crystal ensembles produced by this method can be calculated by using the energies of all configurations produced. The precise method and equations for estimating the free energy difference are presented later in this in section.

2.5.2 Defining a Free Energy Minima

When estimating a free energy for a particular polymorph, a question arises as to how to distinguish and uniquely identify a crystal polymorph ensemble. In crystal energy minimization studies, stable polymorphs appear as local minima on the energy landscape. This energy landscape represents the lattice energy of the system as a function of the spatial positions of all atoms in the system. Specifying all the positions of atoms for a given minima is therefore both necessary and sufficient to define a lattice energy minima.

Minima on the *free energy* landscape, by contrast, represent ensembles of configurations and therefore require a more complex definition than the spatial coordinates of a single structure. Attempts to define minima on an arbitrary free energy landscape generally invoke the use of order parameters [149, 180–182]. Such approaches assume that for an arbitrary free energy surface, there exists a set of order parameters such that the exact distribution of those parameters for a cluster of configurations will be both necessary and sufficient to precisely determine which minima the configurations corresponds with.

For simple systems, such as monotonic fluids, highly symmetric molecules, or colloids, the so-called Steinhardt parameters provide one effective set of order-parameters to distinguish free energy minima [180]. In peptide systems, the Ramachandran map can be used to construct the free energy surface and define free energy minima [182]. For distinguishing between solid and fluid states, one can use the mobility (diffusivity) of the molecules as a metric. In the case of crystal polymorphs, Santiso et al. have previously proposed an algorithm which distinguishes free energy minima using the intramolecular orientation of the molecules as well as their center of mass distance [180]. This approach was effective in distinguishing polymorphs of terephthalic acid.

Clearly these order parameter conditions can be satisfied some system. However it is not guaranteed that all free energy minima will have an easily identified and clearly separated order parameter distribution that can distinguish between polymorphs, particularly when high

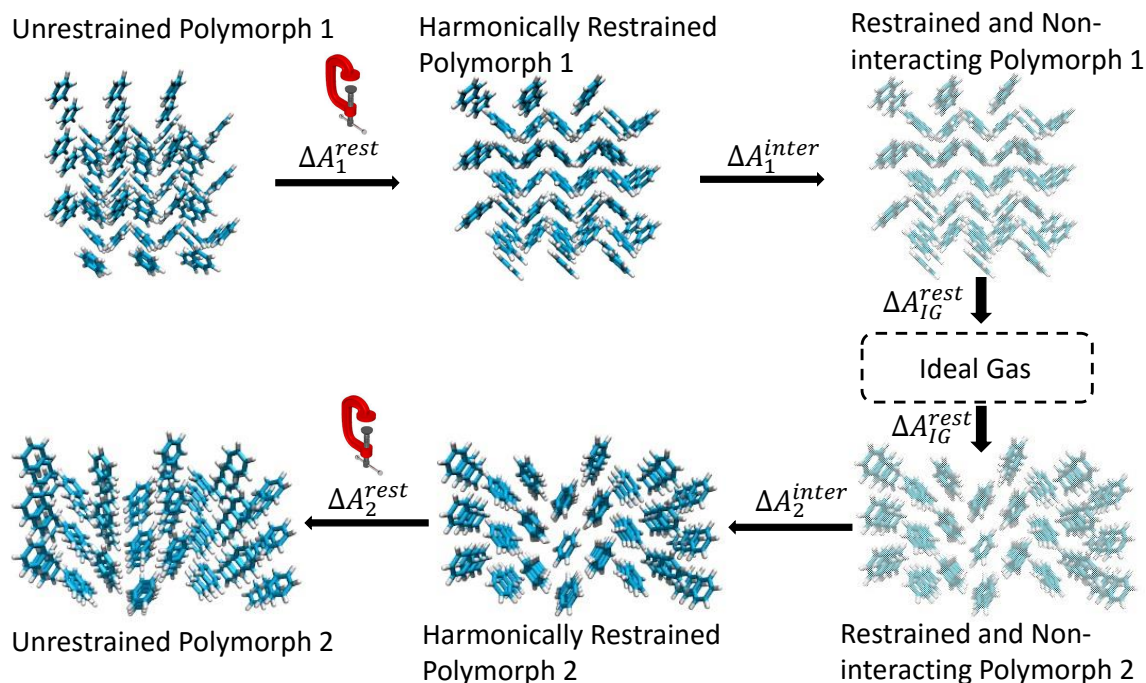


Figure 2.2: The crystal polymorphs are interconverted through a non-interacting ideal gas state using a series of simulations in which the benzene atoms are restrained to their crystalline positions and the intermolecular interactions are turned off.

temperatures cause a merging of two free energy minima. Ultimately, the rigorous and unambiguous definition of a ‘free energy minima’ remains an open research question that is outside the scope of this text. In this work, we define a free energy minima as the set of all configurations which are mutually kinetically accessible without passing through another phase. To distinguish between the two different types of ‘crystal minima’, finite temperature polymorph ensembles will henceforth be referred to as ‘free-energy minima’ and minima on the lattice energy surface will be referred to as either ‘lattice minima’ or ‘energy minima’.

2.5.3 Interconverting Crystal Polymorphs along a Pseudo-supercritical Path

Crystal polymorphs are geometrically distinct and kinetically isolated ensembles of configurations which by definition do not share any configurations. To calculate the free energy difference between polymorphs of organic crystals, each ensemble must either be driven along a thermodynamic path in which the polymorphs interconvert, or alternatively each ensemble must be driven along a path to a reference state where the free energy is known analytically. In this work, the three benzene polymorphs are interconverted through an intermediate ideal gas state using the pseudo-supercritical path (PSCP) method of Maginn and coworkers [183–186]. In this approach, each crystal structure is transformed into an analytically tractable reference state by harmonically restraining the atoms to their crystalline lattice positions followed by turning off all intermolecular interactions. The resulting state with restrained and non-interacting molecules resembles the Einstein crystal model, except that intramolecular interactions are left fully on. The harmonic restraints are then removed analytically to yield an ideal gas state. The pathway is illustrated in Figure 2.2. The total potential energy function used to transform the system from the crystal state to the restrained ideal reference state is given by

$$U_{tot}(\lambda_{rest}, \lambda_{inter}) = U_{intra} + \lambda_{inter}^2 U_{inter} + (1 - \lambda_{rest})^4 \frac{1}{2} k_x (x - x_0)^2 \quad (2.12)$$

where U_{tot} is the total potential energy of the configuration, U_{intra} and U_{inter} are the sums of all intramolecular and intermolecular interactions, k_x is the harmonic restraint constant, and x and x_0 are the current position and restraint position of all atoms in the system, respectively. The λ_{rest} and λ_{inter} terms are coupling parameters used to drive the system from the crystal to gas state. $\lambda_{rest} = 0$ and $\lambda_{rest} = 1$ correspond with the unrestrained and restrained states, respectively, and $\lambda_{inter} = 0$ and $\lambda_{inter} = 1$ correspond with the noninteracting and interacting states. The choice of the exponentials for the removal of interactions and addition of restraints is arbitrary, however we observed that using two and four, respectively, provided high overlap across the entire pathway. The total Helmholtz free energy difference to convert polymorph i into the ideal gas state is

$$\Delta A_i^{PSCP} = \Delta A_i^{rest}(\lambda_{rest} = 0 \rightarrow 1) + \Delta A_i^{inter}(\lambda_{inter} = 1 \rightarrow 0) + \Delta A_{IG}^{rest} \quad (2.13)$$

where ΔA_i^{PSCP} is the Helmholtz free energy for polymorph i , ΔA_i^{rest} is the free energy to restrain the crystalline lattice, ΔA_i^{inter} is the free energy to turn off the interactions in the restrained benzene system, and ΔA_{IG}^{rest} is the free energy to remove restraints from the non-interacting ideal gas state which is the same for all polymorphs. The PSCP is traversed in the NVT ensemble and yields the Helmholtz free energy difference to transform a polymorph from the crystal to the ideal gas state at a given volume. To calculate the constant pressure Gibbs free energy difference for polymorph i , the Helmholtz free energy is integrated over all system volumes with

$$G_i = -k_B T \ln \frac{1}{V_0} \int_V e^{-\beta(A_i(V)+PV)} dV \quad (2.14)$$

where G_j is the Gibbs free energy of polymorph i . The free energy from equation 2.14 will be accurate up to a reference volume V_0 that depends on the units of volume used in the integration. However, this arbitrary constant will cancel when examining free energy differences between polymorphs. A derivation for equation 2.14 is provided in the Supporting Information section A.1 and is similar to the result found by Chong and Ham in Ref. 187. The Helmholtz free energy, relative to the restrained ideal gas state, for polymorph i as a function of volume is:

$$A_i(V) = - \int_{V_{ref}}^V \langle P \rangle dV' + A_i^{PSCP} \quad (2.15)$$

Where A_i^{PSCP} is the free energy change to transform the polymorph into the ideal gas state at the initial reference volume V_{ref} and is calculated from the PSCP. The integral in equation 2.15 is evaluated numerically by integrating the pressure-volume curves for each polymorph fitted

to a third-order polynomial (results shown in the Supporting Information section C.4).

2.5.4 Multistate Reweighting of MD Ensembles with MBAR

The free energies of all independent MD simulations in this work are solved in a self consistent manner with multistate reweighting. There are numerous independently developed methods for implementing multistate reweighting. In this work we specifically use the Multistate Bennett Acceptance Ratio (MBAR) method. A comparison of different reweighting methods and a justification for using MBAR are provided later in Chapter 6. The equations for MBAR are presented here briefly.

The Multistate Bennett Acceptance Ratio was introduced by Chodera and Shirts [188, 189] and provides the provably optimal statistical estimator of thermodynamic quantities using samples generated from multiple thermodynamic states. Given a set of reduced energies $u_i(x_{jn})$ from samples collected in state j and evaluated in state i , the MBAR estimate of the dimensionless free energy for all states is given by

$$f_i = -\ln \sum_{j=1}^K \sum_{n=1}^{N_j} \frac{e^{-u_i(x_{jn})}}{\sum_{k=1}^K N_k e^{f_k - u_k(x_{jn})}} \quad (2.16)$$

where K is the total number of states and N_k is the total number of samples drawn from state K . The reduced energy corresponds with either $u_i(x_{jn}) = \beta U_i(x_{jn})$ at constant volume or $u_i(x_{jn}) = \beta(U_i(x_{jn}) + PV_{kn})$ at constant pressure. The dimensionalized free energy can be recovered from equation 2.16 using the relation $F_k = \beta f_k$ where F is either the Helmholtz free energy or the Gibbs free energy depending of whether the simulations were generated in NVT or in NPT. In the limiting case of two sampled states, equation 2.16 collapses to the Bennett Acceptance Ratio (BAR) [190] and in the case of one sampled and one unsampled state, equation 2.16 is equivalent to the Zwanzig equation [191]. All free energy differences in this work, including the free energies along the PSCP and the free energies to reweight between potentials, are calculated using equation 2.16.

In many instances, the free energy difference across a range of temperatures is desired in order to estimate the thermal stability of different solid forms as well as to determine relative entropies. The exact free energy of the polymorphs as a function of temperature ($G(T)$) cannot be determined with MBAR. However, the relative free energy *difference* of the polymorphs as a function of temperature ($\Delta G(T)$) can be found from the reduced free energies of equation 2.16 using:

$$\Delta G_{ij}(T) = k_B T \left(\Delta f_{ij}(T) - \Delta f_{ij}(T_{ref}) \right) + \frac{T}{T_{ref}} \Delta G_{ij}(T_{ref}) \quad (2.17)$$

where $G_{ij}(T)$ is the free energy difference between polymorph i and j at arbitrary temperature T , $\Delta f_{ij}(T) = \frac{f_i}{k_B T} - \frac{f_j}{k_B T}$ is the reduced free energy difference computed here using MBAR, and $\Delta G_{ij}(T_{ref})$ is the free energy difference at reference temperature T_{ref} supplied by the PSCP. A derivation for equation 2.17 is provided in Appendix A A.1.

2.5.5 The Limits of Thermodynamics-based Crystal Structure Prediction

It is important at this stage to stress the difference between kinetic and thermodynamic stability. The methods above allow one to compute free energy differences between solid phase ensembles. Free energy describes the *thermodynamic* stability of different solid forms. However, showing that a new solid form is more thermodynamically stable than a known structure does not guarantee a conversion under practical conditions.

De Beers[®] has billions of dollars riding on the fact that their metastable diamonds will not convert into the thermodynamically stable graphite under working conditions. Millions of happy couples have planned engagements on this same assumption, despite clear evidence that diamonds are unstable. This example illustrates the important distinction between kinetic and thermodynamic stability. Having a more thermodynamically stable phase with a lower free energy does not mandate that a conversion will take place because high kinetic barriers may separate the two forms.

The structure with the lowest free energy is also not obligated to appear in the first crystallization attempt. The case studies presented in section 1.2 for rotigotine and ritonavir illustrate where the most stable form was not found in any lab until years after the initial crystallization. Strictly speaking, the formation of a particular solid form is controlled by the nucleation barriers for each possible structure during the crystallization process. The structure with the lowest barrier for nucleation and kinetic stability on everyday timescales is the one that will appear first, regardless of the thermodynamical stability. With this in mind, it is tempting to perceive that modeling crystal *kinetics* is the only true way to predict the crystal structures that will appear in the lab.

Predicting crystal structures through direct simulation of nucleation from the melt is not currently feasible for numerous reasons, both practically and philosophically. From a practical standpoint, nucleation of a crystal from the melt can occur on the timescales of seconds to days. Standard molecular dynamics simulations occur on the timescales of nanoseconds to microseconds and can in some extreme cases be pushed to milliseconds. These simulations are simply not long enough to capture spontaneous nucleation without a directive influence. A number of research efforts to alleviate this problem and directly simulate nucleation are ongoing [43,192–195]. However, this technology is still many breakthroughs away from being a mature method.

Predicting structures through directly computing nucleation barriers also has a number of philosophical challenges. Asking ‘which structure has the lowest kinetic barrier to form?’ is a far more ambiguous question than ‘which structure is globally stable?’. The exact kinetic barrier to reach a particular solid form depends on both the precise initial state and the transition state. Is the structure undergoing solid-solid transformation, or is it passing through a melted state? Is the crystallization being driven by slow cooling, or by solvent evaporation? Which solvent is being used? How forcefully is the system being agitated? The exact kinetic barrier to form a particular structure is intimately related to the answer to each of these questions, and therefore ‘which structure has the lowest kinetic barrier to form’ is quite situation dependent.

The molecule 5-fluorouracil has lower barriers to form a metastable structure when crystallized through solvent evaporation in aqueous solution [53, 196, 197]. In this particular crystallization process, the thermodynamically stable form has prohibitively high barriers to nucleate [54]. However, the globally stable form can be easily formed when crystallizing in nonpolar nitromethane. Therefore, nucleation models can predict crystal structures only insofar as the correct questions are being asked of the model. In contrast, the answer to ‘which structure is globally stable’ will always be the same regardless of the crystallization process.

Ultimately, estimating the thermodynamic stability of different crystal forms can provide useful insights into solid form design without any kinetic detail. Once the thermodynamically stable crystal has been established in a manufacturing process, it is unlikely that the material will revert into a new form without significant changes in the external conditions. Therefore, studies which definitively show that the globally stable structure has been found can ensure longevity of the commercial product regardless of the barriers to any other metastable form. Conversely, a clear identification of a new thermodynamically more stable form can warn of an impending transformation and allow preventative measures to be taken before a product hits the market. Therefore in this work, we focus on methods to address the well-defined problem of thermodynamic stability between solid forms. Methods to simulation nucleation and identify favorable crystallization conditions are left for future research efforts.

2.6 Eliminating Unobservable Crystals with Molecular Dynamics

In Chapter 1, we noted that crystal structure predictions often produce far more lattice energy minima than observed polymorphs. Furthermore, there are three prominent explanations put forward for why lattice-energy based computer models produce many minima that are not observed experimentally. In this chapter, we developed the methods for modeling crystal structures with atomistic simulations to produce configuration ensembles around each energy-minimized structure. Eliminating unobserved crystal forms can be accomplished by improving the physical accuracy of the model used to represent each crystal form. The two

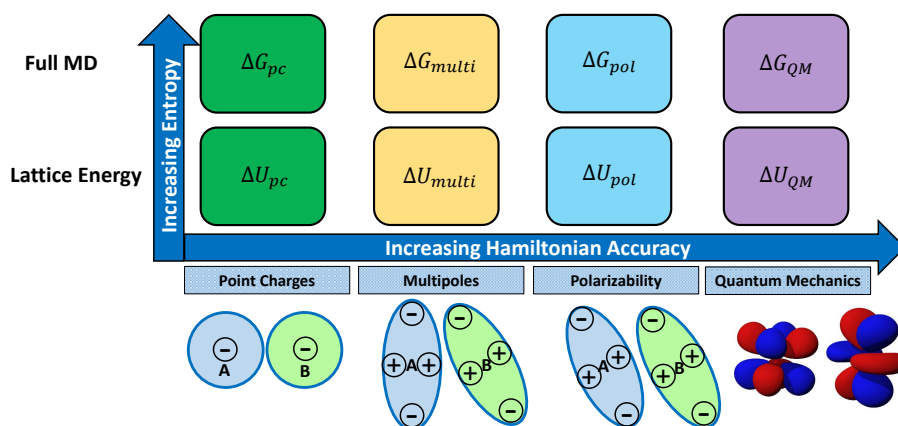


Figure 2.3: The space of all possible ways to model crystals has two dimensions. The first dimension is the ability to capture entropy differences between crystals. The second dimension is the accuracy of the Hamiltonian in capturing the energetic interactions between molecules in the crystals. Each system may require a different level of theory within this space in order to correctly rank crystal stabilities.

dimensions of improvements to the crystal model are shown in Figure 2.3. Below, we describe how molecular simulation approaches can be used to eliminate dead-end structures for each of the three proposals in Chapter 1.

1. Temperature-mediated Stability Reranking

Crystal structures that are stable with low energies at 0 K may become unstable at higher temperature due to a low entropy. At finite temperatures, the correct measure of thermodynamic stability is the free energy. The free energy stability rankings will only match the lattice energy stability rankings if entropies are essentially polymorph independent and therefore cancel.

With molecular simulations, the relative entropy and free energy between crystal forms can be directly computed as a function of temperature from 0 K up to ambient conditions. Equations 2.12–2.17 detail how relative polymorph free energies can be estimated. If a minima is energetically favorable but entropically unfavorable, the free energy will

become less stable as temperature increases. Thus, dead-end structures that rerank and become metastable at ambient conditions will be detected with this methodology. Specific examples where temperature-mediated reranking event are identified with this approach are presented in Chapter 3.

2. Sensitivity to the Hamiltonian

Polymorphs that are favorable in a cheap Hamiltonian may become far less stable when additional terms are added to the electrostatic description. Anisotropic effects like multipoles and induced polarization will be polymorph dependent, indicating that the effects can feasibly reorder the relative stabilities of candidate structures. Molecular simulations can determine this change in stability by computing the relative free energies of all candidate structures in both cheaper and more expensive levels of theory. These effects of more accurate potentials should be probed at finite temperatures rather than at 0 K in order to be decoupled from the hypothesis above.

The finite temperature free energies in more expensive potentials can in principle be computed through direct simulation using equations 2.12 - 2.17. However, the free energies can also be computed efficiently by taking cheaper simulations and adding on a post simulation correction factor to account for the differences to the more expensive potential. This endstate correction factor approach using cheaper simulations is developed in Chapter 4. The effects of going from a point-charge potential to a more expensive polarizable potential in three crystal polymorphs of benzene are also presented in Chapter 4.

3. Lattice Minima Interconversion at Ambient Temperatures

Multiple lattice minima that have low barriers to a more stable structures will appear as a single solid form at ambient conditions. This represents a clear explanation for why a 0 K crystal prediction study would produce many structures that do not appear as unique experimentally observable polymorphs. The interconversion of lattice minima can be detected with molecular dynamics simulations by heating a collection of structures up

to ambient conditions and comparing the configuration ensembles that are produced from each simulation. Lattice minima that have interconverted into the same ensemble will have identical distributions for any of the possible order parameters describing the ensemble. In this work, we use the average energy, cell volume, and free energy as the distinguishing order parameters. In other words, two clusters of configurations are said to belong to the same ensemble if they have an indistinguishable average energy, free energy, and cell volume. More complex order parameters involving the positions and rotations of the molecules may be necessary in certain systems. However, the derivation and implementation of these order parameters is left to a future study.

In chapter 5, the conversion of metastable minima into more stable structures at ambient conditions is presented for multiple rigid and flexible molecules. The structures that interconvert are demonstrated to have indistinguishable energies, free energies, and box volumes as described above. The above procedure is planned to be utilized on a real CSP of resorcinol [198] to determine the number of candidate structures that interconvert at ambient conditions.

Chapter 3

Predicting Temperature-mediated Polymorphic Transformations

3.1 Introduction

Many solid materials do not have a single unique crystal structure and can instead be observed in one of many stable solid forms referred to as polymorphs. Thermodynamic properties such as solubility [14,48,50,96,102,199], hardness [7,34,56], and conductivity [26–30,37,38] can vary substantially between different solid forms. Therefore, commercial interest in a material can be as much dependent on the crystal structure as on the chemical constituents. Commercial manufacturers will often send a candidate material through extensive polymorphic screening experiments to identify all possible solid forms, either to improve the material's performance [27,33,35] or to avoid unexpected transformation in the future [14,48,50].

A complicating factor in the search for commercially viable solid phases is that the stable forms of a material can change based on external factors such as temperature [18,46,92–103], pressure [76–81] and humidity [47,82,83,200]. A complete experimental mapping of all possible solid forms for a given material would therefore require screening at every condition the material will encounter.

Computational crystal structure prediction methods represent an alternative approach to traditional experimental screenings for finding the most stable crystal forms. Computer-based

models can rapidly estimate the properties of a material at hundreds or thousands of thermodynamic states simultaneously with appropriate parallelization. Ideally, these computational methods could be completed at a significantly reduced cost compared with sequential experimental polymorph screening techniques. However, computer models must sufficiently capture all of the relevant molecular details of a physical system in order to capture polymorphic transformations in response to external conditions.

Temperature-mediated polymorphic transformations involve a crystal structure that is globally stable at low temperatures becoming metastable at higher temperatures and subsequently restructuring to a more stable form. Polymorphs that reversibly change stability order as a function of temperature in this manner are referred to as enantiotropically related forms and have been observed experimentally for a number of real organic molecules [18, 46, 92–103]. These temperature-mediated transformations arise from a difference in entropy and heat capacity between the low-temperature and high-temperature crystal forms. The challenge of predicting temperature-mediated transformations therefore reduces to accurately estimating entropy and enthalpy differences between solid forms.

The most common computational crystal structure prediction methodologies compare candidate structures based on their minimized lattice energies [11, 100, 201, 202]. A tacit assumption in these models is that entropy differences between candidate structures contribute negligibly to the relative stability. Excluding entropy significantly reduces the expense of computing the crystal stability. However, the exact contribution of the ambient temperature entropy to the relative crystal stability is an essentially unknown quantity. In many cases, static lattice models without entropy are still sufficiently accurate to identify the experimental structure as the most stable form [11, 100, 112, 113, 201, 203–205]. However, neglecting entropy precludes these models from being able to predict temperature-mediated transformations.

Computational models that include entropy effects typically estimate entropy differences using methods assuming harmonic behavior [92, 142, 158, 159]. In these approximations, all motions in the crystal are assumed to be harmonic about the energy minimized structure. A recent study by Nyman and Day [142] estimated relative entropic contributions between

polymorphic pairs of rigid molecules using a harmonic approximation in over 500 systems and identified temperature-mediated transformations in the range 0 K–300 K in 9% of the systems examined. Harmonic approximations, therefore, represent a possible method for predicting thermal solid-solid transformations. However, the harmonic analysis excludes any anharmonic motions in the system of interest by the very nature of the approach. The quasi-harmonic approximation adds some anharmonic contribution by including changes in the free energy due to thermal expansion [161]. However, full ambient temperature entropy differences, including contributions from anharmonic molecular motions, have not yet been directly computationally reported, and may not agree with the values estimated from these harmonic approximations for typical organic crystals. Additionally, the Nyman et al. study examined only rigid molecules using a rigid harmonic approximation, meaning that this estimate is a lower bound of the presence of solid-solid transformations in organic crystals in general.

Molecular dynamics (MD) simulations can in principle capture all entropic contributions to the ambient temperature stability of different crystal forms including those from anharmonic motion, lattice expansion, intramolecular fluctuations, and temperature-dependent vibrational frequencies. When molecular dynamics models are applied to systems with temperature-mediated transformations, the models should reveal that the high temperature crystal form has a higher entropy than the low temperature form as long as the energy function is accurate enough to describe the local energy landscape around each crystal ensemble.

In this study we examine whether including the entropies from classical point-charge MD simulations in the relative crystal stability estimates is sufficient to predict temperature-mediated crystal transformations in small molecule organic crystals. Using molecular dynamics to compute the relative entropies could overcome a prominent shortcoming in lattice-energy based models which neglect harmonic and/or anharmonic entropy contributions. In this study, we simulate 12 polymorphic systems that are known experimentally to have a temperature-mediated transformation. These systems were chosen to cover a range of sizes, flexibilities, and crystalline order. In each system we examine whether the molecular dynamics model correctly predicts a significantly higher entropy for the high temperature form relative to the

most stable 0 K form. The free energies computed from MD are directly compared against those with a cheaper quasi-harmonic approximation to capture the importance of anharmonic motions in the thermodynamic properties estimates of crystal ensembles at ambient temperature. We also compute the enthalpy difference as a function of temperature to determine how much high temperature enthalpy differences deviate from the 0 K lattice energy difference and potentially facilitate reranking.

Finally, we examine the frequency at which OPLS-AA [206–208] correctly predicts a larger entropy for the high temperature polymorph relative to the 0 K structure. We compare this to the frequency that OPLS-AA correctly identifies the low temperature form as having a lower minimized lattice energy in order to determine whether the relatively cheap point-charge OPLS-AA potential performs well at estimating relative entropies even in systems where the relative lattice energies are not well represented. As a final measure of accuracy, we compare the estimates of entropy and enthalpy from OPLS-AA to experimental values in the literature.

3.2 Methods

3.2.1 Initial System Structures

A total of 12 molecules were chosen for this study, shown in figure 3.1. All of these systems are known experimentally to undergo a temperature-mediated polymorphic transformation. Therefore, they represent ideal candidates to test whether molecular dynamics simulations with OPLS-AA can effectively predict a large positive entropy difference between high and low temperature crystal forms. Additionally, these systems had well-characterized structures for both the high and low temperature form in the Cambridge Structural Database (CSD). The experimentally determined polymorphic transition temperatures and CSD refcodes for each crystal are shown in Table B.1 of Appendix B for reference. The molecule set includes both small and large molecules as well as rigid and flexible molecules. Additionally, two systems (cyclopentane and succinonitrile) with dynamically disordered solid forms are also included.

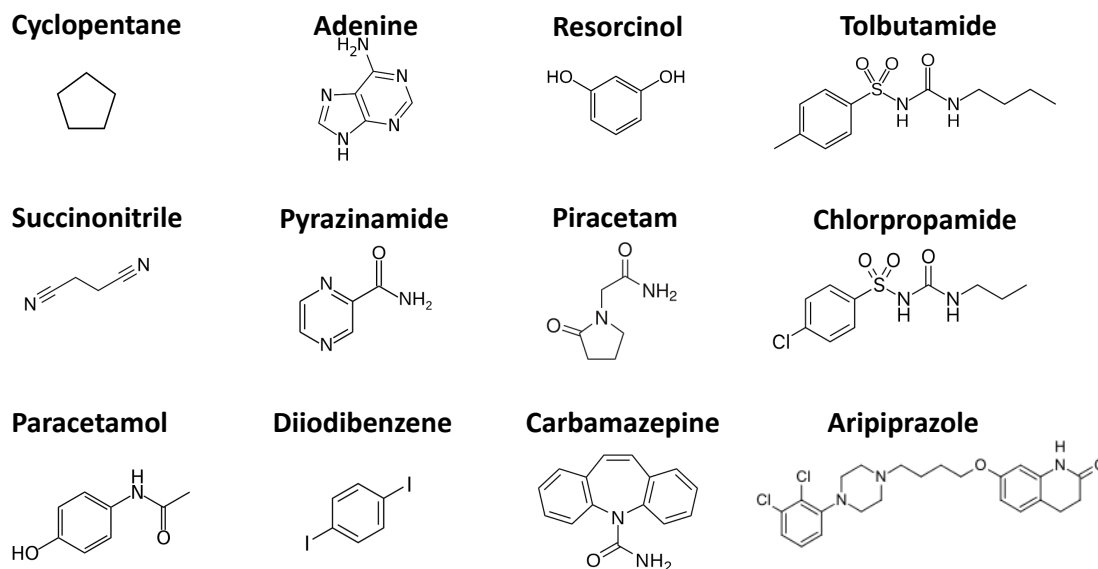


Figure 3.1: Molecules examined in this study.

The starting crystal structures were taken directly from the Cambridge Structural Database (CSD). The unit cells were replicated into larger crystal supercells such that all supercell dimensions were > 1.6 nm. The choice of 1.6 nm satisfies the condition in GROMACS of supercell dimensions longer than twice the long-range interaction cutoff of 0.8 nm. Interaction cutoffs larger than 0.8 were previously shown to have a negligible effect on the average energy difference between crystal polymorphs with proper treatment of longer range periodic interactions [119]. Additionally, the large supercells accommodate between 32–96 independently moving molecules in each system which approximates real crystals better than single unit cells. The resulting supercells from this procedure were then minimized to find the closest lattice energy minima of each structure in the OPLS-AA potential.

We observed during the course of this work that some of the minimized crystals from the

CSD restructured into a new more stable form upon heating. These restructuring events represent irreversible kinetic escapes from metastable lattice minima and are distinct from the reversible transformations of enantiotropic polymorph pairs at a coexistence temperature. The restructured configurations from these events were subsequently crystal minimized according to the above process and used for all further ambient temperature simulations. A more detailed discussion of this restructuring process and the implications on computed CSP are given in Chapter 5.

The dynamically disordered form I of cyclopentane and form II succinonitrile did not have a well-defined crystal structure in the Cambridge Structural Database (CSD). To obtain the putative lattice minima for form I cyclopentane, the form III structure was heated to 140 K. At this temperature the structure spontaneously changed to a new lattice minima which had box vectors within 2% of those of the experimental form I at 140 K [100]. This spontaneous restructuring in a short MD simulation was also seen in other MD simulation studies of cyclopentane [100]. Form II succinonitrile was obtained by changing the box vectors of the form I lattice minima to match the experimental form II and shifting the molecules' center of masses proportionally. This structure was then heated to 200 K and then crystal minimized. There is no rigorous method to prove that the resulting structure is the experimental form II succinonitrile. However, the similar box vectors and space group to experiment [209] as well as the observed dynamic disorder and higher entropy compared with form I all strongly suggest that this structure corresponds closely with experimental form II succinonitrile.

3.2.2 Free Energy Estimation with MBAR

The relative stability of different polymorphs were determined from 0 K through 300 K by computing relative Gibbs free energy curves as a function of temperature. The precise method for generating the temperature-dependent Gibbs free energy curves was similar to the process we developed previously [119]. The method is briefly summarized here.

First, the relative free energies of all polymorphs from molecular dynamics are determined at a single reference temperature and pressure by driving each crystal structure to an ideal

reference state using a variant of the pseudo-supercritical Path (PSCP). The details of this approach are described in this previous study [119]. Briefly, the molecules in the crystal structure are atomically restrained to their minimized lattice positions. The intermolecular interactions are then alchemically removed from the simulation. Finally, the effects of the harmonic restraints are removed from the system analytically, yielding the ideal gas state. In most rigid small molecules, the intramolecular energies in the minimized structure are polymorph independent and therefore contribute equally to the free energy [119]. However, for the more flexible molecules aripiprazole, tolbutamide, chlorpropamide, and the disordered succinonitrile, the torsion and 1-4 Coulomb and van der Waals interactions were alchemically removed along with the intermolecular interactions. This last step is necessary to account for differences in the 1-4 interactions and equilibrium torsion energies between different polymorphs. We note that at ambient conditions the PV contribution to the free energy is negligible, so that $\Delta A \approx \Delta G$ as observed in other systems [119,185].

The temperature-dependence of the relative free energies were determined by simulating each polymorph at a range of temperatures between 10 K through 300 K (or the melting point) at ambient pressure in intervals of 10 K. We simulated with smaller 1 K intervals when the overlap (following the definition in Ref 119) between adjacent temperature states dropped below 0.0001. Simulating below 10 K is challenging due to kinetic trapping within the coarse energy landscape. However, the free energy is estimated to make a linear transition from 10 K to 0 K in all systems based on the free energy intersecting the 0 K minimized lattice energy within 0.005 kcal·mol⁻¹ when the free energy is extrapolated from both 20 K and 10 K (see Appendix B section B.2). We note that in many systems the free energy follows a linear trend well above 10 K and would therefore require significantly less MD simulations than we used here to create the full temperature-dependent free energy. However, for completeness we directly simulated all systems down to 10 K regardless of when the system free energy converged to a linear descent to 0 K.

The free energy to alchemically transform the crystals along the PSCP and the free energy to heat the crystals to high temperature were computed using the Multistate Bennett Acceptance

Ratio (MBAR) method. Given a set of reduced energies $u_i(x_{jn})$ for samples collected from state j and evaluated in state i , the MBAR estimate of the dimensionless free energy for all states is given by:

$$f_i = -\ln \sum_{n=1}^N \frac{e^{-u_i(x_n)}}{\sum_{k=1}^K N_k e^{f_k - u_k(x_n)}} \quad (3.1)$$

where K is the total number of states and N_k samples are drawn from each of the k states, with $N = \sum_k N_k$ total samples drawn from all states. The reduced energy corresponds with either $u_i(x_n) = \beta U_i(x_n)$ at constant volume or $u_i(x_n) = \beta(U_i(x_n) + PV_n)$ at constant pressure. The dimensionalized free energy can be recovered from equation 3.1 using the relation $F_k = \beta f_k$ where F is either the Helmholtz free energy or the Gibbs free energy depending of whether the simulations were generated in NVT or in NPT. Expectation averages of u_i can be computed by

$$\langle u \rangle_i = \sum_{n=1}^N u_i(x_n) W_i(x_n) \quad (3.2)$$

$$W_i(x_n) = \frac{e^{f_i - u_i(x_n)}}{\sum_{k=1}^K N_k e^{f_k - u_k(x_n)}} \quad (3.3)$$

Entropy differences between polymorphs are computed from the temperature-dependent free energy estimates using the relationship $\Delta S(T) = \frac{\Delta U(T) - \Delta G(T)}{T}$, calculated with MBAR, which can take into account statistical correlations between the expectation of $\Delta U(T)$ and the log partition function $\Delta G(T)$ take into account [210].

3.2.3 Exploring the Magnitude of Entropic Contributions

Lattice energy-based CSP studies use the minimized lattice energy to approximate the stability difference at ambient conditions. The true stability difference between solid forms at any temperature, T , is given by $\Delta G(T) = \Delta U(T) - T\Delta S(T)$. The difference in minimized lattice

energy will be an accurate approximation of $\Delta G(T)$ if the energy difference is roughly independent of temperature such that $(\Delta U(T) - \Delta U(0K)) \ll \Delta U(0K)$ and if the entropy difference is small such that $T\Delta S(T) \ll \Delta U(0K)$.

We critically examine how well the minimized lattice energy approximates the free energy for each of the 12 systems studied in this study by directly computing the magnitudes of $T\Delta S(T)$ and $(\Delta U(T) - \Delta U(0K))$ as a function of temperature. The absolute magnitude of $T\Delta S(T)$ will typically increase with temperature, and it is therefore necessary to specify a temperature in order to assess the significance of the entropy in relative crystal stabilities. In this study, we choose to compute the magnitudes of $T\Delta S(T)$ and $(\Delta U(T) - \Delta U(0K))$ at ambient temperature (300 K) for all systems regardless of the transition temperature or melting point.

A completely accurate prediction of the transition temperature requires that both the lattice energy *and* the entropy be correctly estimated in the sampled potential. However, a large entropy difference is indicative of a temperature-mediated transition even in systems where the lattice energy is poorly estimated with OPLS-AA. Therefore, we assess the prediction of a polymorphic transformation based on a sufficiently large entropy difference between two forms at ambient conditions. This allows us to evaluate the model's ability to compute entropy differences irrespective of the accuracy of the minimized lattice energy difference.

The precise value wherein an entropy difference becomes 'significant' and predicts a re-ranking is inherently subjective and case specific. All polymorphic pairs will have some non-zero entropy differences at all temperatures. However, entropy differences that are sufficiently small are unlikely to change the polymorph stability rankings before the material melts. Therefore it is necessary to define a magnitude of $T\Delta S$ at ambient temperature which makes a solid-solid transformation tenable. In the recent study by Nyman and Day [142], roughly 10% of the 508 polymorphic pairs had lattice energy differences less than $0.5 \text{ kJ}\cdot\text{mol}^{-1}$ (or $0.12 \text{ kcal}\cdot\text{mol}^{-1}$). Therefore, to be 90% confident that a transformation will not occur, the contribution of $T\Delta S(T)$ to the free energy should be smaller than $0.12 \text{ kcal}\cdot\text{mol}^{-1}$ at ambient conditions. Conversely, a $T\Delta S(T)$ contribution greater than $0.12 \text{ kcal}\cdot\text{mol}^{-1}$ signifies that a temperature-mediated transformation cannot be immediately ruled out.

Decoupling the assessment of accurate minimized lattice energies from accurate entropy contributions allows us to compare the reliability of OPLS-AA in generating the enthalpy and entropy individually. A previous free energy study of benzene in both a cheap point-charge potential and a more expensive polarizable Hamiltonian revealed that the choice of potential had a stronger effect on the lattice energy at 0 K than on the free energy at 250 K [119]. This suggests that free energy differences (and potentially entropy differences) may be faithfully represented in cheap potentials even if the minimized lattice energies are incorrectly ranked. Here we directly compare the accuracy of the OPLS-AA potential in producing correct relative lattice energy rankings as well as correct predictions of $T\Delta S(T) > 0.12 \text{ kcal}\cdot\text{mol}^{-1}$.

As a final assessment of OPLS-AA accuracy, we compare the computed minimized lattice energies and entropies to estimates from experimental measurements. The minimized lattice energy corresponds physically with temperatures approaching 0 K and is used here to probe the accuracy of OPLS-AA in models without temperature effects. The entropy at 300 K is used to test the accuracy of OPLS-AA in capturing the effects of thermal motion at ambient conditions independent of the minimized lattice energy. However, experimental measurements of relative entropy and enthalpy are generally performed at a single transition temperature, which is different for each system and may not correspond with either 0 K or 300 K. Therefore in our comparison we make the modest assumption that relative entropies and enthalpies do not change significantly with temperature, which we show to be true in section 3.3.1.

3.2.4 Quasi-harmonic Approximation

The equations and methodology for computing free energy differences with a quasi-harmonic approximation are provided in Chapter 2 section 2.3.2. The equations for computing the Gibbs free energy are reproduced below:

$$A_v(V, T) = \sum_k \frac{1}{\beta} \ln(\beta \hbar \omega_k(V)) \quad (3.4)$$

$$G(T) = U(V) + A_v(V, T) + PV \quad (3.5)$$

where $A_v(V, T)$ is the Helmholtz free energy at a given volume V and temperature T and $G(T)$ is the Gibbs free energy. The Gibbs free energy from equation 3.5 is computed at the volume where the free energy is minimized at the temperature, T .

3.2.5 Simulation Details

All molecular dynamics simulations were carried out with the GROMACS 5.0.4 [138] simulation package. The simulations were temperature controlled with stochastic dynamics [173] using a characteristic time constant of 1 ps. Pressure was maintained at 1 bar using the anisotropic Parrinello-Rahman barostat [174] with a time constant of 10 ps. All GROMACS simulations were run for 5 ns with the first 0.5 ns omitted for equilibration. All molecules were modeled using the point-charge OPLS-AA potential. Long-range Lennard-Jones and electrostatic interactions were treated with Particle Mesh Ewald summation at a cutoff of 0.8 nm with a Fourier grid spacing of 0.13 nm. The energy difference between organic crystal polymorphs is insensitive to the cutoff distance with these settings near 0.8 nm as shown previously [119].

The initial crystal structures were taken directly from the Cambridge Structural Database (CSD). The unit cells were replicated into larger crystal supercells such that all box dimensions were larger than twice the long-range cutoff of 0.8 nm (Table B.1). The resulting supercells were then minimized to a RMS gradient of $0.01 \text{ kcal} \cdot \text{mol}^{-1} \cdot \text{\AA}^{-1}$ using the XTALMIN subroutine within TINKER. The structures' box vectors and atom positions were then further minimized in GROMACS using an energy minimization to a maximum force of $< 2 \text{ kJ} \cdot \text{mol}^{-1} \cdot \text{nm}^{-1}$ alternating self-consistently with a Nelder-Mead minimization algorithm (as implemented in numpy version 1.8.2) with the box vectors as optimization variables until a final minimum energy crystal structure is found.

The form III cyclopentane crystal was found to be unstable with the anisotropic Parrinello-Rahman barostat, leading to unphysical box vector fluctuations. However, these unphysical fluctuations were not observed with either the anisotropic Berendsen barostat or the isotropic Parrinello-Rahman barostat. We therefore chose to equilibrate the crystal at each temperature using the anisotropic Berendsen barostat and run the resulting structure with the isotropic Parrinello-Rahman barostat for the production simulation.

The alchemical pseudo-supercritical path calculations were run with the same simulation settings as our previous work [119]. 20 quartically spaced intermediate states were used to add harmonic restraints, and 10 quadratically spaced intermediate states were run to remove intermolecular interactions. The positions of all harmonic restraints were generated by energy minimizing each polymorph without intermolecular interactions to find the relaxed vacuum geometry for all molecules.

The convergence of the free energies along the thermodynamic path were verified by comparing the free energy difference estimate at 100 K through a direct PSCP calculation and through a PSCP calculation at 200 K combined with the free energy difference between 200 K and 100 K using equation 3.1. The direct PSCP at 100 K gave the same result within uncertainty to the estimate produced by starting at 200 K and decreasing the temperature to 100 K. For cyclopentane, which melts below 200 K the PSCP was run at 50 K and closed at 30 K. For the more flexible aripiprazole molecule, the closure was verified at 300 K due to better sampling along the PSCP at high temperatures. The thermodynamic cycle closure for all systems is shown in Appendix B for reference section B.3.

Uncertainties in the free energy estimates were determined from the variance of 200 independent bootstrap samples of the uncorrelated data. The trajectories were uncorrelated using the method used first by Muller-Krumbhaar and later by Chodera [211–214]. The uncorrelated trajectories for both the PSCP and the simulations over the temperature range were subsampled 200 times and used to generate 200 independent estimates of the temperature-dependent free energy. The variance at each temperature was used as the uncertainty in the free energy at that point.

Vibrational analysis, structure minimizations, and lattice energies for quasi-harmonic analysis were carried out with the Tinker 7.1.3 molecular modeling package using the same OPLS-AA potential as used for molecular dynamics simulations. We note that long-range interactions in Tinker were treated differently than in GROMACS because Tinker does not have an implementation of PME for van der Waals interactions and Tinker uses a grid-size parameter for PME as opposed to the grid-spacing used in GROMACS. However, the differences in minimized energy difference between Tinker and GROMACS were below $0.005 \text{ kcal}\cdot\text{mol}^{-1}$ in all cases.

Structures discussed in section 3.1 were isotropically compressed and expanded to 100 reference states of different volumes. First, the structures were compressed or expanded from their initial structure by multiplying each lattice vector by a constant; similarly, the molecules are adjusted to keep the center of mass at the same fractional position with respect to the the lattice vectors. Next, the structure was minimized using a BFGS nonlinear optimization implemented in Tinker's `minimize` executable. This minimization routine keeps the lattice parameters constant. Lastly, the minimized structure is used to determine a denser (or less dense) structure, repeating the process. Compression was performed for lattice vector fractions from 1.0 to 0.97, at a step size of 0.001. Expansion was performed for lattice vector fractions from 1.0 to 1.06, using the same step size. In practice, the Gibbs free energy of all systems used lattice structures with lattice vector fraction between 1.0 and 1.04.

The lattice minimum energy, for each compressed and expanded structure, was found using Tinker's `analyze` executable. The Gruneisen parameter for each polymorph was found to determine how the vibrational spectra changes with expansion. The vibrational frequencies of the two references states used in the Gruneisen parameter were found using the `testhess` executable, which calculates the frequencies through diagonalization of the structure's numerical mass-weighted Hessian. For each polymorph, the vibrational spectra is calculated for the global lattice minimum structure and a structure expanded by a lattice vector fraction of 1.001. With these two reference states and their vibrational spectra a Gruneisen parameter for each k th wavenumber can be calculated. In equation 3.6, γ_k is the Gruneisen parameter, ω_k is the

wavenumbers, $\omega_k(V)$ is vibrational frequency, and V is the volume.

$$\gamma_k = -\frac{d \ln \omega_k(V)}{d \ln V} \quad (3.6)$$

Equation 3.6 can be solved numerically and γ_k can be used in equation 3.7 to determine the vibrational frequencies at any given isotropic volume. In equation 3.7 $\omega_k(V)$ is the vibrational frequency at a new specific volume, $\omega_k(V_0)$ is the vibrational frequency at the global lattice minimum structure, V is the volume of the new structure, v_0 is the volume of the original structure, and γ_k is the Gruneisen parameter for that specific frequency.

$$\omega_k(V) = \omega_k(V_0) \left(\frac{V}{V_0} \right)^{-\gamma_k} \quad (3.7)$$

Using equations 2.8 and 2.9, Gibbs free energy curves were created for a temperatures of 1, 5, 10 K and steps of 10 K until the temperature cutoff for used in MD simulations is reached.

3.3 Results and Discussion

3.3.1 Entropy and Enthalpy at Finite Temperature

Lattice-energy based crystal prediction models assume that entropic contributions ($T\Delta S$) between crystal forms are much smaller than the difference in minimized lattice energy. The known experimental temperature transformations in the systems examined here indicate that this condition should break down, with ($T\Delta S$) being relatively large around 0.12 kcal·mol⁻¹ or higher. The entropic contribution to the free energy difference at 300 K between the low and high temperature forms of all systems are shown in Figure 3.2. The contribution of $T\Delta S > 0.12$ kcal·mol⁻¹ in 9 of the 12 systems confirming that entropy effects cannot be ignored. Therefore, the molecular dynamics model in the OPLS-AA potential clearly identifies that entropic contributions to the free energy are not negligible in these system, as hypothesized.

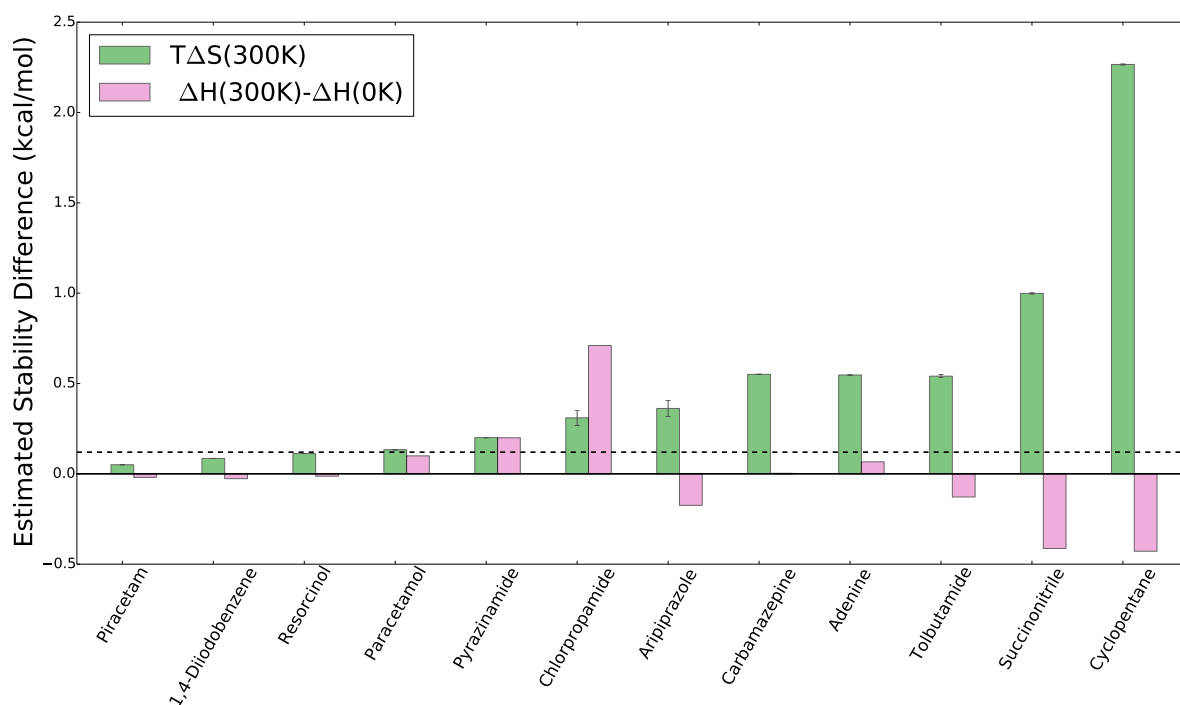


Figure 3.2: Temperature and entropy contributions to the relative stability of the high temperature form at 300 K are significant in 9 of the 12 systems studied. The dashed line indicates $0.12 \text{ kcal}\cdot\text{mol}^{-1}$ which in this study indicates a significant entropy contribution to the relative stability. The change in enthalpy difference from 0 K to 300 K are generally smaller than the temperature/entropy effects at 300 K. However, the change in enthalpy difference is larger in one of the more flexible molecules. Note that the entropy difference for cyclopentane was used at 110 K because the system melts well before 300 K.

The decrease in relative free energy of the high temperature form at hotter conditions could in theory result from a change in the enthalpy difference rather than due to a difference in entropy. However, the change in the enthalpy difference between 0 K and 300 K is generally smaller than $T\Delta S$ in a majority of the systems examined in this study (Figure 3.2). The change in energy difference is only larger than $T\Delta S$ in the flexible chlorpropamide system. The large contributions of $T\Delta S$ coupled with the relatively small contributions of the enthalpy change confirm that the experimental re-ranking of the polymorph stabilities in these systems likely result from a difference in entropy between the low temperature and high temperature crystal forms rather than a change in enthalpy difference.

The large entropy differences and temperature-mediated transformations presented in this study illustrate the importance of including entropy in crystal structure prediction studies. As one specific example of the importance of entropy, the estimated free energy difference between carbamazepine form I and form III was extended to 500 K and shown in Figure 3.3. The minimized lattice energy difference in OPLS-AA indicates that carbamazepine form I is less stable than form III by more than $0.5 \text{ kcal}\cdot\text{mol}^{-1}$. Experiments performed from room temperature up to the melting point of 463 K reveal that both form I and form III are observable, with form I being more stable above 435 K [93]. The large difference in minimized lattice energy, produced by a model that does not include entropy, may initially be interpreted as an inconsistency with experiments and may falsely be attributed to an insufficient accuracy of the OPLS-AA potential. However, the Gibbs free energy differences in OPLS-AA at 300 K shows that both form I and form III have similar stabilities, and above 435 K form I is correctly identified as being the more stable form. Therefore, carbamazepine represents a clear example where entropy needs to be included for an accurate comparison between theory and experiments performed near ambient conditions. Indeed, many previous failures of lattice-energy based CSP studies to identify the most stable experimental structure at high temperature may reflect a failure of the model to account for entropy rather than a failure of the Hamiltonian to describe the potential energy surface.

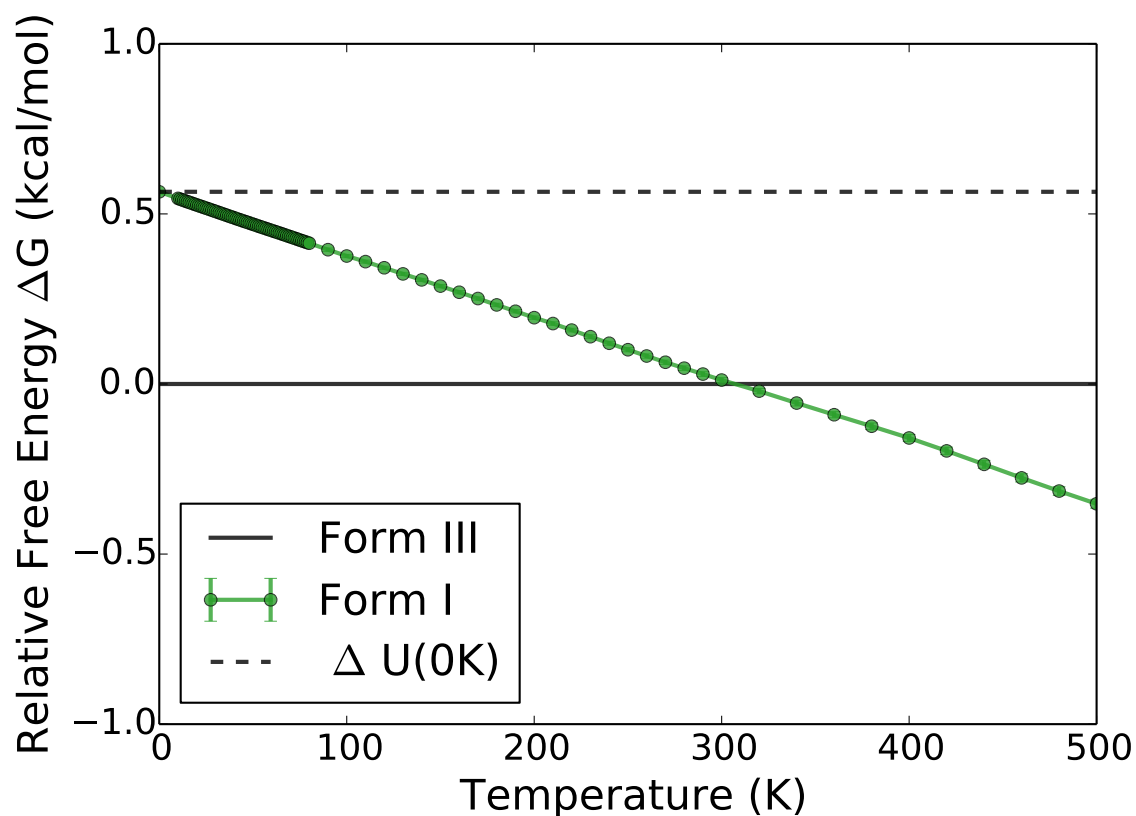


Figure 3.3: The carbamazepine stability difference without entropy suggests that form I is significantly less stable than form III. The relative free energy including temperature and entropy reveals that form I and form III have similar stabilities at ambient conditions, consistent with both forms being experimentally observable.

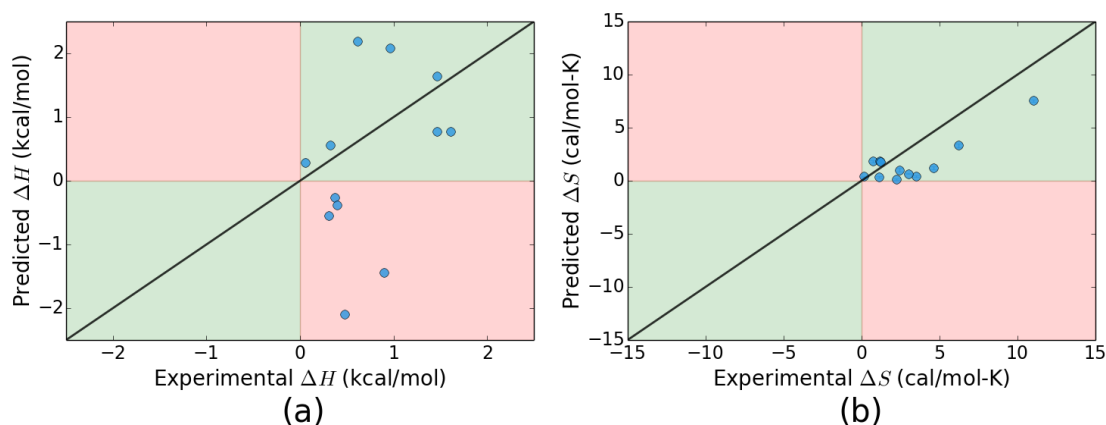


Figure 3.4: OPLS-AA estimates the relative entropy of the high and low temperature form more accurately than it estimates the enthalpy difference when compared to experiments. The green quadrants indicate where the model and experiments agree on the relative ranking of the two structures. Red quadrants indicate where the model and experiments disagree on the ranking. Only 58% of the low temperature forms using the OPLS-AA force field had the lower minimized lattice energy, but all high temperature forms using OPLS-AA have a higher entropy than the low temperature form.

The OPLS-AA potential generally provides a more accurate estimate of the entropy difference than the enthalpy difference in the systems studied here. Both the computed lattice energy difference and entropy difference at 300 K are compared against experimental estimates of the entropy and enthalpy of transition in Figure 3.4. The green colored quadrants in Figure 3.4 indicate where the OPLS-AA potential and experiments agree on the relative ranking of the two crystal forms. The model correctly identifies the high temperature form as having a higher entropy than the low temperature form in all 12 systems examined. Conversely, OPLS-AA identifies the low temperature form as having a lower minimized lattice energy in 7 of the 12 systems (58%). This result suggests that cheaper potentials are better able to capture entropic differences between polymorphs than they are in describing the minimized lattice energy difference at 0 K, at least in these systems.

3.3.2 Comparison of Molecular Dynamics and QHA

The quasi-harmonic approximation effectively produces free energy estimates in the small rigid molecules in this work. Figure 3.5 shows the relative free energy as a function of temperature for six small rigid molecules in this work estimated with both molecular dynamics and the quasi-harmonic approximation. In these systems, the relative free energy is essentially indistinguishable between the two estimation methods over the entire temperature range. In the case of pyrazinamide, the deviation in the free energy estimate with MD and QHA at 300 K is $0.005 \pm 0.002 \text{ kcal}\cdot\text{mol}^{-1}$ which is effectively zero. The largest deviation at 300 K in these six systems occurs for carbamazepine, with a difference in estimated free energy of $0.07 \text{ kcal}\cdot\text{mol}^{-1}$.

The close agreement between the QHA-derived and MD-derived free energies at low temperatures is not surprising. At 0 K the two methods give the same free energy difference by definition. At temperatures near 0 K, all motions in the crystal should be roughly harmonic around the energy minimized structure. However, one would expect small anharmonic motions to appear in the crystals at higher temperatures where the molecules are free to move well beyond the single most favorable lattice position. The lack of divergence between QHA

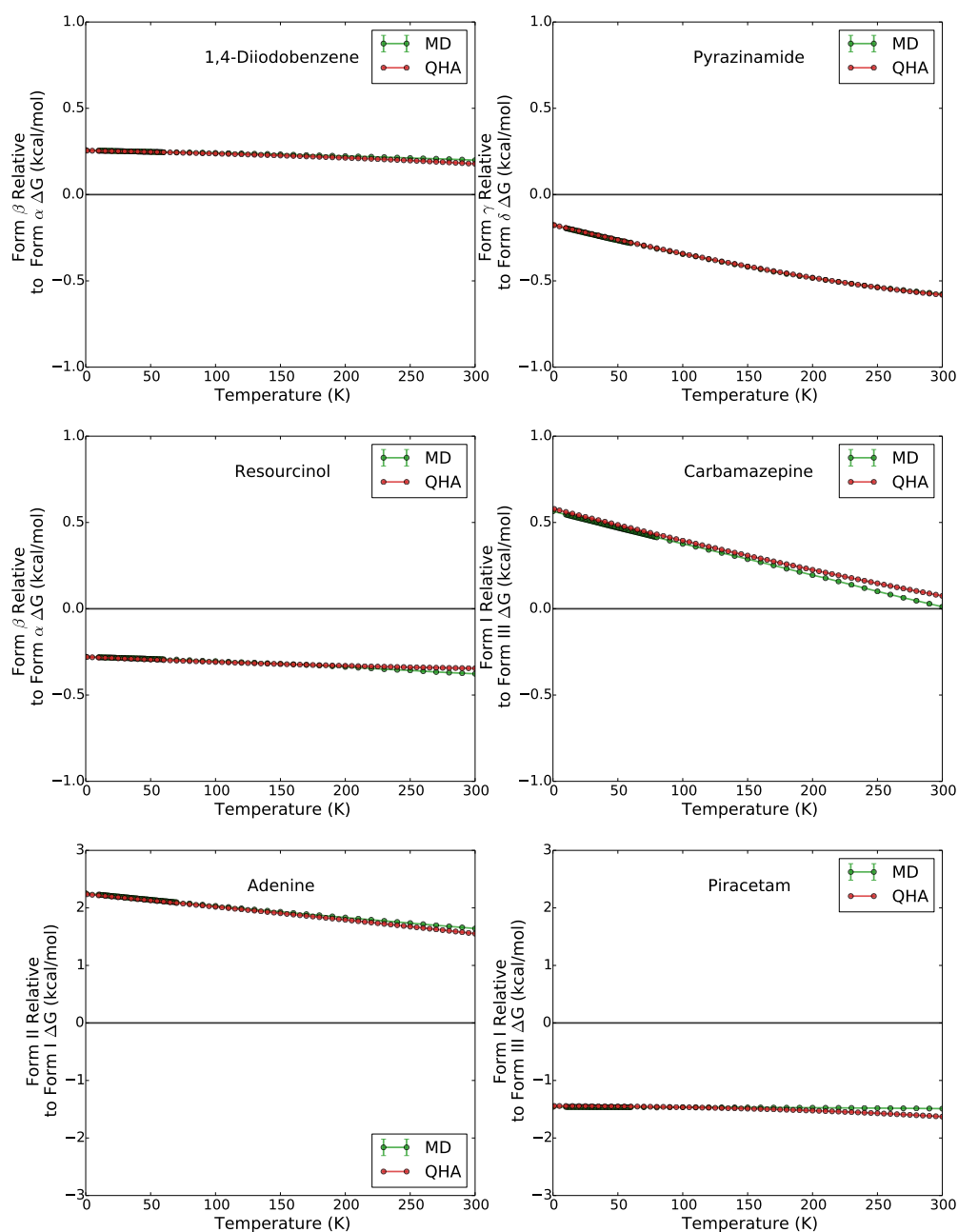


Figure 3.5: The relative polymorph free energies are well approximated with the quasi-harmonic approximation in the six small rigid molecules in this study. The free energies are essentially indistinguishable at all temperatures from 0 K to 300 K.

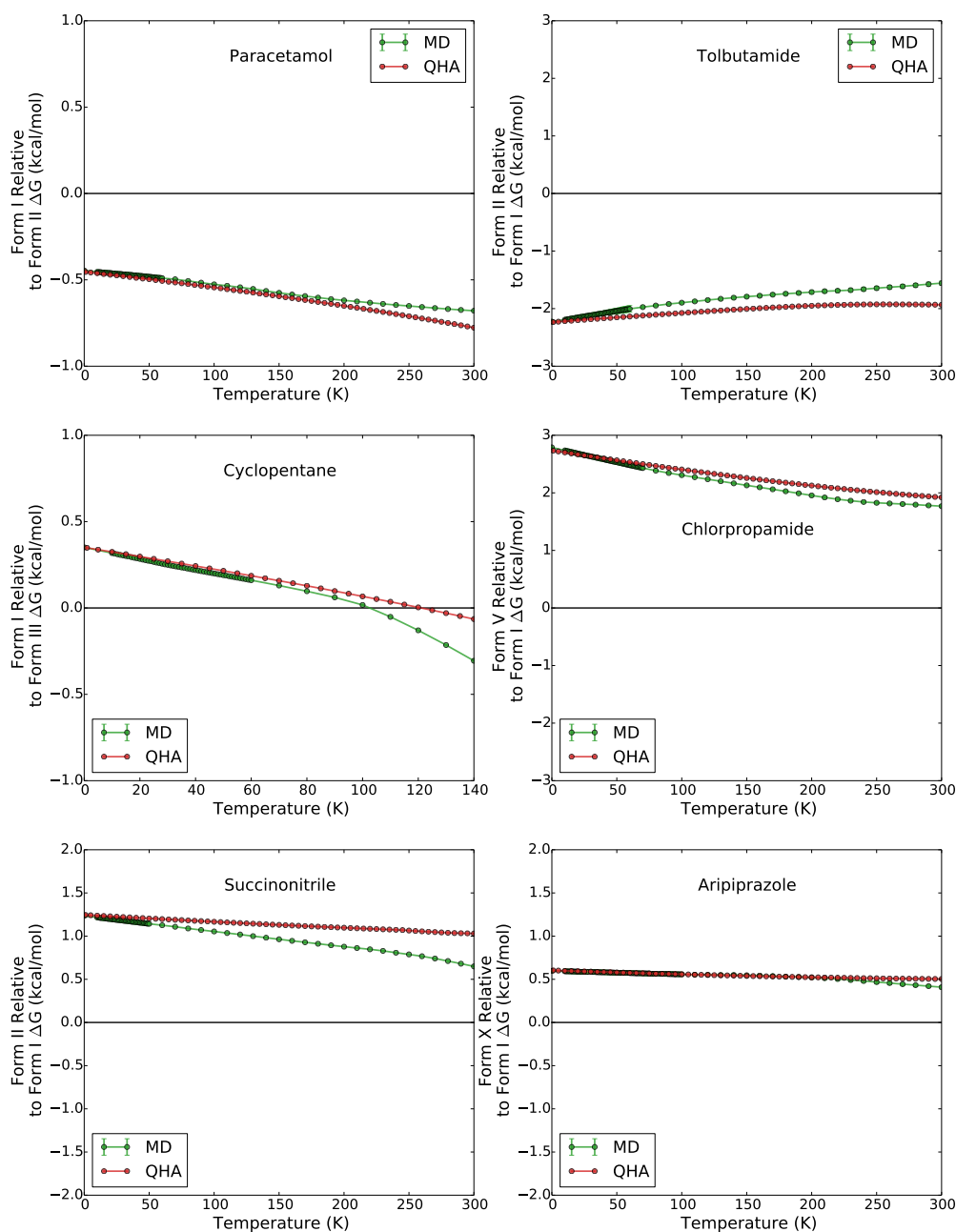


Figure 3.6: The relative polymorph free energies from molecular dynamics diverge from those with the quasi-harmonic approximation in the flexible and disordered crystals in this study. The difference in free energy between the two methods leads to different estimated polymorphic transformation temperatures.

and MD in Figure 3.5 even at higher temperatures near ambient conditions suggests that the difference in entropy between the polymorphs still comes predominately from harmonic motions in the crystals rather than anharmonic motions in these systems.

Deviations between the free energy estimates with the quasi-harmonic approximation and molecular dynamics do appear at ambient temperatures in the larger and more flexible molecules in this study. The temperature-dependent free energy estimates for the six molecules with the largest deviation between the two methods are shown in Figure 3.6. QHA and MD have similar predictions of the free energy at low temperatures. However, the free energy estimates of the two methods tend to diverge at higher temperatures near ambient conditions. This deviation is most pronounced in succinonitrile where the two methods differ on the 300 K free energy by $0.37 \text{ kcal}\cdot\text{mol}^{-1}$. The divergence suggests that anharmonic motions in these systems have a significant contribution to the relative entropy of different crystal structures that do not cancel between polymorphs.

The results from Figure 3.5 and Figure 3.6 suggest that the presence of flexibility and/or disorder dictates whether a quasi-harmonic approximation will provide an accurate estimate of the relative free energy. All six molecules examined in Figure 3.5 have no flexible or rotatable degrees of freedom. By contrast, the six molecules presented in Figure 3.6 all have at least one rotational degree of freedom or have one dynamically disordered crystal form. The divergence between QHA and MD in these six molecules, coupled with the lack of significant divergence in the small rigid molecule systems in Figure 3.5, indicates that flexibility and/or disorder are key factors in predicting the efficacy of a quasi-harmonic approach.

Differences in the estimated free energy from QHA and MD ultimately manifest in different predicted temperatures of polymorphic transformation. The temperature-dependent free energy of cyclopentane in Figure 3.6 illustrates this effect. Both MD and QHA predict that the dynamically disordered form I cyclopentane has a higher entropy than form III and therefore becomes more stable at higher temperatures. However, the order-disorder transition in form I around 90 K results in a sudden increase in entropy due to the new anharmonic full-body rotations of the cyclopentane molecules. This rapid entropy increase is captured by the

molecular dynamics approach. However, the transition is not captured by the quasi-harmonic approach because the new modes of motion are not a harmonic vibration around the minimum energy structure. The predicted polymorphic transformation temperature is the temperature at which the free energy difference between the crystal structures is zero. The different free energy estimates from QHA and MD in cyclopentane lead to a different predicted transformation temperature by 20 K as a result of QHA not accounting for the anharmonic entropy contributions in the disordered form I crystal.

The ‘large’ and ‘flexible’ systems examined in this work are noticeably smaller and less flexible than many common pharmaceutically relevant compounds. Aripiprazole is the largest molecule in the study with 57 atoms with 7 flexible degrees of freedom. The polymorphic pharmaceutical molecule ritonavir have 98 atoms and 28 flexible degrees of freedom. It is reasonable to expect that the magnitude of the anharmonic contribution to the free energy and the magnitude of the divergence between QHA and MD will be larger in ritonavir than in any of the molecules studied here.

3.4 Conclusions

The relative entropies of twelve polymorphic organic systems with known temperature-mediated transformations were computed with classical point-charge molecular dynamics simulations. The temperature and entropy contributions to the free energy were larger for the high temperature form than the low temperature form in all systems examined. This result is consistent with experimental observations of temperature-mediated transformation of the low temperature form at high temperatures and indicate that the point-charge simulations are proficient at ranking the entropies of different crystal polymorphs.

The OPLS-AA potential incorrectly assigned a lower minimized lattice energy to the high temperature form in 5 of the 12 systems examined. This incorrect ranking combined with the qualitatively correct relative entropies leads to these systems never having a coexistence point between the low and high temperature forms in the model. We therefore conclude that the

point-charge OPLS-AA potential is not sufficiently accurate for quantitative predictions of the transition temperature between enantiotropic solid forms. However, the OPLS-AA potential identified a $T\Delta S > 0.12 \text{ kcal} \cdot \text{mol}^{-1}$ at ambient conditions in 9 of the 12 systems and identified the high temperature form as having a higher entropy than the low temperature form in all systems. The higher frequency of OPLS-AA in correctly ranking the relative entropy compared to the minimized lattice energy suggests that cheaper point-charge potentials may still be useful in estimating entropy differences even in systems where a more expensive potential is needed to characterize the lattice minima. Indeed, the estimated entropy differences in OPLS-AA are closer to experimental measurements than the enthalpy estimates in the systems examined.

In many of the systems, particularly the small rigid molecules, the relative entropies remain constant over the temperature range from 0 K up to 300 K and have linear free energy differences. In these small rigid systems the free energy differences are well approximated by a simpler harmonic approximation and deviate from the MD estimates by less than $0.07 \text{ kcal} \cdot \text{mol}^{-1}$ at all temperatures up to 300 K. However, systems with dynamic disorder and systems with multiple degrees of flexible motion had significant differences in entropy between MD and QHA as large as $0.37 \text{ kcal} \cdot \text{mol}^{-1}$ at 300 K. This deviation increases the predicted polymorph transformation temperature in cyclopentane by 20 K. We therefore recommend using a full anharmonic approach, such as molecular dynamics, when computing ambient temperature stability differences between polymorphs of molecules with multiple flexible degrees of freedom or disordered crystal forms.

The thermal motions at ambient temperature in some systems also led to multiple lattice minima interconverting into an indistinguishable configuration ensemble at ambient temperatures. This interconversion results from small transition barriers between similar crystal structures that are easily crossed at 300 K. Therefore, fully atomistic molecular dynamics simulations provide an additional benefit beyond anharmonic entropy contributions by identifying degenerate minima on the lattice energy surface that correspond with a known solid form rather than a new observable polymorph.

Chapter 4

Exploring the Effects of Polarization Through Hamiltonian Reweighting

4.1 Introduction

A long sought-after goal in the field of computational chemistry is the ability to predict the crystal structure of small-molecule organic solids. The presence of multiple experimentally observable crystal structures, or polymorphs, for certain chemical compounds makes this problem particularly challenging. Predicting all of the potentially stable crystal structures is of particular interest to the pharmaceutical industry, where the presence of a more stable polymorph can render commercial drugs ineffective and lead to costly market recalls [14, 33, 48, 50–54]. The influence of polymorphism in commercial materials is not limited to pharmaceutical applications and has been observed to affect the hydraulic properties of cement [44], the reactivity of energetic materials [20, 22–24], the weather resistance of dyes [25], the transparency of optical photonics materials [42, 43] and the charge carrier mobility in semi-conductor materials [26–30].

Computer-based models can, in principle, rapidly and inexpensively take a particular

molecule as input and search through possible crystal structures, identifying the set of stable or metastable polymorphs which are likely to be observed experimentally. Previous computational methods for polymorph prediction have largely focused on ranking a set of candidate crystal structures based on lattice energies. These prediction methods then identify all structures within a few $\text{kJ}\cdot\text{mol}^{-1}$ of the global minima as potentially observable polymorphs [10, 11, 31, 33, 111–113, 116, 140, 202, 215–221].

Fully atomistic simulation approaches such as Monte Carlo or molecular dynamics simulations are better able to represent crystals at finite temperature relative to static lattice models because crystal structures are treated as ensembles of configurations rather than a single snapshot. Because MD acts on ensembles, complex effects can be observed which influence the stability and would be missed in lattice energy-based studies [144–154]. As one specific example, Yu and Tuckerman modeled six polymorphs of benzene with a free energy sampling method and proposed that the experimentally observed high-pressure benzene crystal may actually be a mixture of the benzene II and benzene III structures typically identified as distinct polymorphs in energy minimization studies [147]. This type of mixed stacking structure highlights the need for finite-temperature free energy approaches.

Although finite temperature models better approximate real crystals, the additional expense of molecular dynamics simulations generally limits the complexity of the sampling potential to simple point-charge models. Potentials with polarizability are better able to account for anisotropic or polar environments and have been shown to be critical when calculating free energies of binding and solvation in solutions [222–226] as well as ion diffusion [227], and protein folding [224, 228]. The inclusion of polarization is likely to be even more important in periodic solids due to long-range electrostatic correlations, and numerous crystallographic studies based on lattice energies have emphasized the need for distributed multipoles or full polarization for proper polymorph prediction [33, 111–113, 116, 216, 229–235]. For example, Williams and Weller modeled various azabenzene structures and concluded that no point-charge model will be able to reproduce the correct structures unless additional lone pair sites are included [229]. Reilly and Tkatchenko also showed that the stability of the observed form

I of acetylsalicylic acid (aspirin) results from a coupling of the electronic fluctuations and lattice vibrations in higher level energy models rather than a mechanical instability of form II as previously proposed [234]. Multiple other studies, including blind structure predictions, have demonstrated that lattice energy methods with more realistic Hamiltonians have better agreement with experimental measurements [109–113, 115, 115, 215, 235–237].

The improvement in lattice energy-based studies suggests that including polarization may also improve the agreement between fully atomistic molecular dynamics simulations and experimental measurements. However, the effect of these realistic Hamiltonians on free energy differences remains an open research question because of the higher computational demand of molecular dynamics simulations. We must therefore be able to calculate converged free energy differences between finite-temperature polymorphs for a range of approximations of the electrostatic energy function in order to determine the effects of these approximations on the polymorph free energy landscape. In large-scale testing, we will need to carry out these free energy methods without relying on direct brute-force simulations run with the most expensive potential.

Warshel et al. [238] introduced a method based on the pioneering work of Kirkwood [239] and Zwanzig [191] to calculate free energies in complex Hamiltonians indirectly by simulating in a cheap potential and calculating ensemble averages in the expensive potential using Boltzmann reweighting, thus avoiding direct simulation in the expensive potential [240–258]. The formula for estimating the ensemble average in a target state given samples in another state is [191]:

$$\Delta G_{12} = -\beta^{-1} \ln \left\langle e^{-\beta \Delta U_{12}(x)} \right\rangle_1 \quad (4.1)$$

$$\langle X \rangle_2 = \left\langle X(x) e^{\beta \Delta G_{12} - \beta \Delta U_{12}(x)} \right\rangle_1 \quad (4.2)$$

Where $U_{12}(x)$ is the energy difference between energy function 1 and energy function 2, ΔG_{12} is the free energy difference between the two potentials, and $X(x)$ is any observable of the system.

This reweighting scheme of calculating a free energy difference between two unsampled states indirectly through two sampled states with equation 4.1 is illustrated in Figure 4.1. The most commonly used reweighting method to calculate free energies to a target Hamiltonian is the Zwanzig relationship (also known as the EXP method) [191] which is simply equation 4.1. This method produces the correct free energy estimates in the limit of large sampling. However, the Zwanzig reweighting method requires a prohibitively large amount of sampling to converge estimates in systems with low configuration space overlap between the cheap and expensive potential [259].

The overlap O between two states 1 and 2 can be quantitatively measured with:

$$O_{12} = 2 \int \frac{P_1(x)P_2(x)}{P_1(x) + P_2(x)} dx \quad (4.3)$$

where $P_i(x)$ is the probability of configuration x in state occurring in state i . In previous studies employing Boltzmann reweighting, poor overlap was attributed to differences in the intramolecular interactions between the cheap and expensive potential [247,253]. This poor overlap has led researchers to using alternate reweighting strategies such as replacing the total energy difference in equation 4.1 with just the difference in intermolecular interactions, while holding the intramolecular interactions constant [247–250]. The target potential with this approach has intermolecular interactions that correspond with the expensive potential and intramolecular interactions at the cheaper level of theory. Thus the resulting free energy produced by equation 4.1 with this substitution reflects the difference in free energy between the full cheap potential and an unsampled expensive potential with cheap intramolecular interactions and expensive intermolecular interactions. The additional free energy cost to convert the intramolecular interactions from the cheap to expensive Hamiltonian is assumed to be negligible. The bias introduced by this approximation is in many cases small, [247,249,250] but there is presently no rigorous bounds or systematic prediction of the error introduced by this approximation. Therefore we will not investigate this approach here.

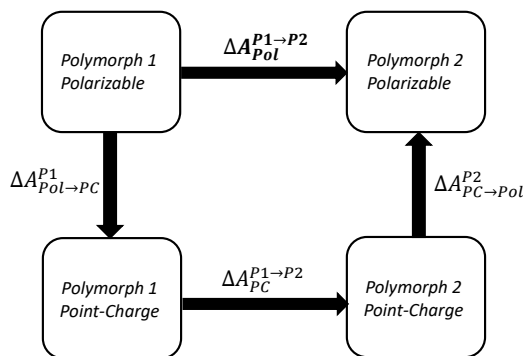


Figure 4.1: The thermodynamic cycle for interconverting two crystal polymorphs in a polarizable potential using the free energy difference in a point-charge potential. The desired free energy difference in the polarizable potential is indicated in bold.

Other studies have proposed augmenting the cheap simulations with a multistage Hybrid Monte Carlo algorithm to reject unfavorable configurations in the expensive potential [251,252]. With this approach, the original cheap potential can still be used for sampling despite poor overlap to the expensive potential. However a large number of sampled configurations will be rejected with this approach if the overlap is very low, and thus the sampling efficiency in the cheap potential will be poor. Both the modification to the Zwanzig equation and this Hybrid Monte Carlo scheme represent changes to the post-simulation analysis method to overcome poor overlap between the sampled and target potential.

One can also directly improve the overlap between the sampled and target Hamiltonian by designing customized potentials which improve the frequency that a high-probability configuration in the target Hamiltonian will be visited in the cheaper sampled state [245,246,253,260–262]. Wang et al. created custom MM potentials with forces that match those in a quantum-based Hamiltonian and this process, known as adaptive force matching (AFM), improved the

convergence of Boltzmann reweighting from a classical to DFT-based potential by over 2 orders of magnitude when estimating the dielectric constant of ice. [253,260]. Warshel et al. also conducted Hamiltonian reweighting using empirical valence bond (EVB) potentials specifically designed to match the QM states and concluded that the designed potentials enabled the free energy barrier of enzymatic reactions to be determined in ab initio potentials. [245,246].

In this chapter, we present a methodology to rapidly determine the effect of including polarizability in Hamiltonians used for crystal polymorph free energy evaluations. The methodology avoids sampling in the expensive polarizable potential by creating point-charge potentials with intramolecular parameters that match the corresponding terms in the target Hamiltonian. This process greatly increases the configurational overlap between the point-charge and polarizable Hamiltonian and enables Boltzmann reweighting with full configuration energies rather than just the intermolecular interactions. Designing these point-charge potentials requires no parameterization or optimization and only requires inputting the intramolecular parameters from the polarizable model into the point-charge potential. Using these designed point-charge potentials can in some systems increase the overlap enough to enable reweighting with sampling done only in the cheap point-charge model without any direct sampling in the polarizable potential. This reweighting without direct sampling is important in cases such as ab initio Hamiltonians, where sampling in the target potential is prohibitively expensive.

To our knowledge, this is the first use of Hamiltonian reweighting from lower-level to higher-level energy functions applied to organic crystal polymorphs. Additionally, this is the first direct comparison of the differences between point-charge and polarizable Hamiltonians in calculating the relative Gibbs free energies of sets of organic crystal polymorphs. In both the point-charge and polarizable potentials, we compare the free energy differences with the differences in minimized lattice energy to determine the relative significance of including entropy in polymorph prediction studies.

We demonstrate this reweighting technique on three polymorphs of solid crystalline benzene depicted in Figure 4.2. Benzene is known experimentally to have multiple observable polymorphs that are each globally stable at various temperatures and pressures [32,76], and

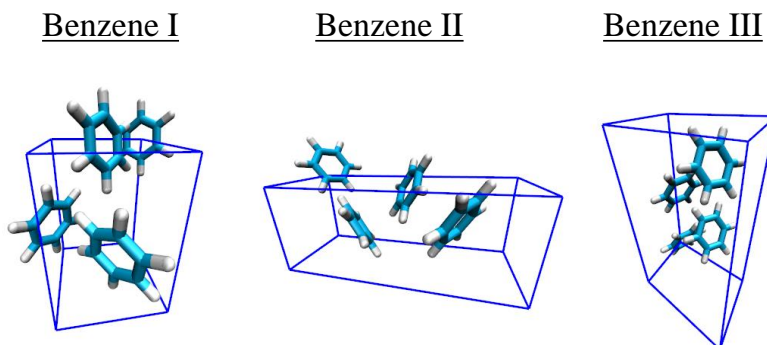


Figure 4.2: The three benzene crystal structures studied in this work.

the simplicity of the molecule makes it an ideal candidate to demonstrate the implementation of the customized reweighting methodology presented in this work. For the three polymorphs of benzene, we calculate the free energy difference between the crystal structures at finite temperature in the point-charge OPLS all-atom potential [206–208] and the point-charge GROMOS54A7 [263] potential as well as the polarizable AMOEBA09 potential [264–266]. The polymorph free energy differences in the AMOEBA09 potential are calculated without direct AMOEBA09 sampling by computing the relative polymorph free energies in a designed point-charge potential and adding on the free energy necessary to transform the Hamiltonian from the point charge to the polarizable potential as shown in Figure 4.1.

To validate the overall process of calculating a free energy difference indirectly with sampling in a designed point-charge potential, we first calculate the relative free energies of the three benzene polymorphs in the GROMOS54A7 potential both with direct brute force sampling and with the indirect pathway using sampling only in a designed point-charge potential with matching intramolecular parameters. We demonstrate here that the free energies with the direct and indirect pathways converge to the same values. We then calculate the free energies of the three benzene polymorphs in the polarizable AMOEBA potential [264–266] using the

same indirect approach without sampling in AMOEBA and show that the uncertainties from this method are low enough to distinguish the most stable polymorphs at these conditions. Finally, we compare the differences in relative free energies of the benzene polymorphs in the OPLS-AA, GROMOS54A7 and AMOEBA09 potentials over a range of temperatures and estimate the added value of using a polarizable description for the electrostatic interactions as well as the added-value of using finite temperature free energy methods over 0K lattice-energy-based methods in polymorph prediction studies.

4.2 Methods

4.2.1 Computing free energies between point charge potential polymorphs using a pseudo-supercritical path

The task of computing temperature-dependent free energy curves for one system involves computing the relative free energies of all polymorphs at a single temperature and pressure followed by simulating each polymorph over the entire desired temperature range. The polymorphs are interconverted here using the Pseudo-supercritical Path method described in Chapter 2. The relative free energies at each temperature are solved self-consistently using the MBAR multistate reweighting estimator. The relevant equation for computing free energies with MBAR is also discussed in detail in Chapter 2. The main equations for the PSCP and MBAR are reproduced below.

The total Helmholtz free energy difference to convert polymorph i into the ideal gas state is:

$$\Delta A_i^{PSCP} = \Delta A_i^{rest}(\lambda_{rest} = 0 \rightarrow 1) + \Delta A_i^{inter}(\lambda_{inter} = 1 \rightarrow 0) + \Delta A_{IG}^{rest} \quad (4.4)$$

where ΔA_i^{PSCP} is the Helmholtz free energy for polymorph i , ΔA_i^{rest} is the free energy to

restrain the crystalline lattice, ΔA_i^{inter} is the free energy to turn off the interactions in the restrained benzene system, and ΔA_{IG}^{rest} is the free energy to remove restraints from the non-interacting ideal gas state which is the same for all polymorphs. The Helmholtz free energy in the average crystal cell is essentially identical to the Gibbs free energy at ambient pressure. However, the Helmholtz free energy from equation 4.4 can be used to compute the Gibbs free energy by integrating the Helmholtz free energy over all box volumes. See chapter 2 for more details of this process.

The Multistate Bennett Acceptance Ratio [188, 189] provides the provably optimal statistical estimator of thermodynamic quantities using samples generated from multiple thermodynamic states. Given a set of reduced energies $u_i(x_{jn})$ from samples collected in state j and evaluated in state i , the MBAR estimate of the dimensionless free energy for all states is given by

$$f_i = -\ln \sum_{j=1}^K \sum_{n=1}^{N_j} \frac{e^{-u_i(x_{jn})}}{\sum_{k=1}^K N_k e^{f_k - u_k(x_{jn})}} \quad (4.5)$$

where K is the total number of states and N_k is the total number of samples drawn from state K . The reduced energy corresponds with either $u_i(x_{jn}) = \beta U_i(x_{jn})$ at constant volume or $u_i(x_{jn}) = \beta(U_i(x_{jn}) + PV_{kn})$ at constant pressure. The dimensionalized free energy can be recovered from equation 4.5 using the relation $F_k = \beta f_k$ where F is either the Helmholtz free energy or the Gibbs free energy depending of whether the simulations were generated in NVT or in NPT. All free energy differences in this work, including the free energies along the PSCP and the free energies to reweight between potentials, are calculated using equation 4.5.

Although the exact free energy change of the polymorphs as a function of temperature cannot be determined, the relative free energy *difference* of the polymorphs as a function of temperature can be found from the reduced free energies of equation 4.5 using:

$$\Delta G_{ij}(T) = k_B T \left(\Delta f_{ij}(T) - \Delta f_{ij}(T_{ref}) \right) + \frac{T}{T_{ref}} \Delta G_{ij}(T_{ref}) \quad (4.6)$$

where $G_{ij}(T)$ is the free energy difference between polymorph i and j at arbitrary temperature T , $\Delta f_{ij}(T) = \frac{f_i}{k_b T} - \frac{f_j}{k_b T}$ is the reduced free energy difference computed here using MBAR, and $\Delta G_{ij}(T_{ref})$ is the free energy difference at reference temperature T_{ref} supplied by the PSCP. A derivation for equation 4.6 is provided in Appendix A section A.1.

4.2.2 Point-charge and Polarizable Potentials

The OPLS-AA and GROMOS54A7 potentials both follow the traditional point-charge potential form in which the total potential energy is the sum of intramolecular and intermolecular interactions according to:

$$U_{tot} = U_{bond} + U_{angle} + U_{dihedral} + U_{vdw} + U_{elec} \quad (4.7)$$

where U_{tot} is the total potential energy, U_{bond} is the bond stretching term, U_{angle} represents the angle stretching interactions, $U_{dihedral}$ represents the proper and improper dihedral torsions, U_{vdw} contains the intermolecular van der Waals forces, and U_{elec} contains the Coulombic interactions between all point charges in the system. Both OPLS-AA and GROMOS54A7 have this same decomposition of total energy. However, the precise functional forms of these individual terms are slightly different between the two potentials. Even in cases where the functional forms are the same, the two potentials have slightly different parameters, such as a difference in equilibrium benzene bond lengths. A full description of the functional forms and parameters for all force fields is presented in Appendix C section C.3.

The AMOEBA force field incorporates the effects of distributed multipoles and induced polarizability in addition to the point-charge electrostatics and intramolecular interactions present in the GROMOS54A7 and OPLS-AA potentials. The total potential energy in AMOEBA is described by the equation:

$$U_{tot} = U_{bond} + U_{angle} + U_{str-bnd} + U_{torsion} + U_{oop} + U_{vdw} + U_{elec}^{perm} + U_{elec}^{ind} \quad (4.8)$$

where $U_{str-bnd}$ characterizes the stretch-bend interactions, $U_{torsion}$ is the dihedral torsion interactions, U_{oop} is the out-of-plane deformation energy, and U_{elec}^{perm} and U_{elec}^{ind} describe the permanent and induced electrostatic interactions, respectively. The additional energetic terms in the AMOEBA09 Hamiltonian relative to GROMOS54A7 and OPLS-AA come in the U_{elec}^{perm} and U_{elec}^{ind} terms of the total energy equation. A full description of the functional forms for the intramolecular interactions as well as the multipoles and induced polarization are presented in Appendix C section C.3.

The designed point-charge potentials in the present work, Designed-GROMOS and Designed-AMOEBA, were created by starting with the OPLS-AA potential form and modifying the intramolecular interaction terms to match those in the GROMOS54a7 and AMOEBA09 potential, respectively. The parameters in the designed potential can be set to any arbitrary value in order to increase the overlap to the AMOEBA09 potential. The accuracy of the designed point-charge potential in representing the true physical system is not relevant because only the polymorph end states in the expensive potential need to have physical significance.

The total potential energy for each designed potential follows the functional form of equation 4.7. The functional forms and parameters for the bonded interactions were chosen to match the corresponding terms in the target potential. The van der Waals and electrostatic interactions were taken without modification from the OPLS-AA potential. Full details of the designed point-charge potential functions and the precise values of the interaction parameters are included in Appendix C section C.3.

The unit cell shape of benzene I in AMOEBA09 is markedly different than the cell shape in the point-charge Designed-AMOEBA potential as shown in Figure 4.3. In the AMOEBA09 potential, Benzene I has a cubic $Ra3$ crystal structure, while in all examined point charge potentials Benzene I has a orthorhombic $Pbca$ structure both at 0K and at finite temperature. This difference in box vectors leads to virtually no overlap between the benzene I ensembles in the polarizable and point-charge potentials despite having matching intramolecular parameters and overlapping volume distributions. To create sufficient overlap for the benzene I crystal, the box shape in the Designed-AMOEBA point-charge potential is artificially constrained to

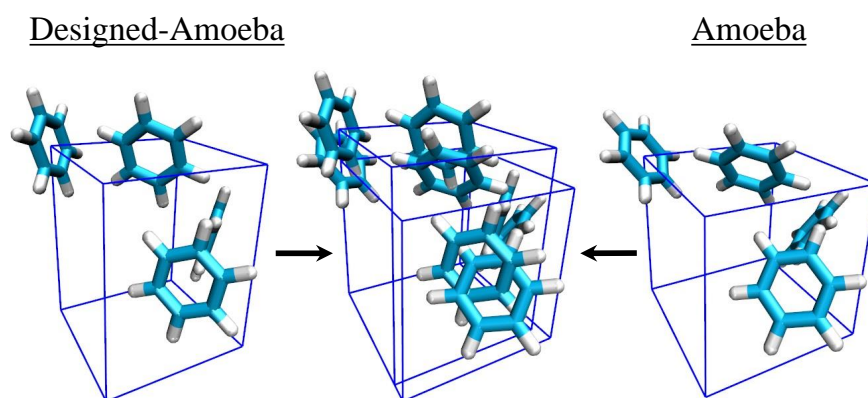


Figure 4.3: The benzene I crystal structure has a different unit cell shape and space group in the polarizable AMOEBA09 potential and the point-charge potentials. In the AMOEBA09 potential, the unit cell is perfectly cubical in the $Ra3$ spacegroup with lattice parameters 0.77, 0.77, 0.77 nm. In the point-charge potentials benzene I has the $Pbca$ spacegroup and has lattice parameters 0.73, 0.93, 0.69 nm in the Designed-AMOEBA potential. The total volume of both unit cells is 0.46 nm^3 . The unit cells depicted here were generated by minimizing the lattice energy in position space and box vector space in TINKER.

the equilibrium aspect ratio in the AMOEBA09 potential and simulated isotropically in the NPT ensemble. The configurations from the Designed-AMOEBA sampling therefore may not correspond with the true free-energy minima in the Designed-AMOEBA Hamiltonian at $P = 1$ atm and $T = 200K$. However only the true minima of the *target* Hamiltonian needs to be sampled in the indirect reweighted pathway. The Designed-AMOEBA ensemble is ultimately connected to the same restrained ideal gas state and thus the polymorph used in the PSCP can be artificially manipulated in any fashion to improve the overlap.

The aspect ratio of benzene I in the AMOEBA potential is determined without dynamic sampling by minimizing the crystal lattice energy of the QM-derived benzene I crystal structure in the AMOEBA Hamiltonian. We note that constraining the box vectors to the AMOEBA09 aspect ratio is only necessary for benzene I as the unit cell shapes for benzene II and benzene III are similar enough in the polarizable and point-charge potentials to achieve high overlap.

4.3 Simulation Details

Molecular dynamics simulations and energy calculations were carried out using the GROMACS 5.0.4 [138] and TINKER 7.0 [137] simulation packages. Initial polymorph structures were obtained by ab initio QM energy minimization at the MP2/aug-cc-pVDZ level using a recently developed HMBI fragment approach. [267].

The benzene crystal was modeled as a single unit cell with 4 independently moving benzenes in each polymorph for direct comparison with standard unit cell lattice energy minimization approaches. Because the GROMACS simulation package does not allow cutoffs longer than the system size, this small system was simulated using a 3x3x2 supercell but symmetrizing the forces so that all unit cells move identically. Because of the small size of the system, configuration space overlap between simulations with different potentials is much larger compared to the overlap obtained with larger super cells. However, the free energy differences between potentials for each polymorph were the same in the 4- and 72- benzene crystals within

0.01 kcal·mol⁻¹ suggesting that differences due to supercell size canceled between potentials. For more details see Appendix C section C.5.

The correlated motions present in smaller supercells will also artificially decrease the entropy of each polymorph compared to a crystal of infinite extent. The relative OPLS-AA free energy of benzene II and benzene III changed by 0.03 kcal·mol⁻¹ and 0.10 kcal·mol⁻¹, respectively, when using the 72- benzene supercell rather than the 4-benzene supercell. This difference is not enough to change the predicted globally stable benzene polymorph. However, it suggests that larger supercells may be necessary for other polymorphic systems. For more details see Appendix C section C.5.

The smaller 4-benzene systems were able to access conformations not accessible to larger supercells which are related to the original polymorph by a mirror plane symmetry operation. These transformations between symmetric conformations are an artifact of the correlated motions in the small benzene system and do not affect the thermodynamics of the system because the original and flipped configurations are thermodynamically equivalent. However, this symmetry is broken when harmonic restraints are added to the system, and significant statistical noise appears in the restraint process of the PSCP as a result of the transformations between symmetric configurations. To remove this noise, the symmetrically flipped configurations (determined by the angle between the planes of the benzenes relative to the initial configuration) were removed from the set of sampled configurations before the PSCP free energy calculation using equation 4.5.

Long-range Lennard-Jones and electrostatic interactions were calculated using Particle Mesh Ewald summation with a cutoff of 0.7 nm and a Fourier spacing of 0.13 nm. The change in the system potential energy due to a change in cutoff radius at this cutoff is essentially the same for all polymorphs and the free energy difference between the polymorphs is therefore insensitive to the cutoff distance in this range. For more details see Appendix C section C.6.

For the PSCP calculations, each of the 3 benzene polymorphs were initially equilibrated in the NPT ensemble at 1 bar and 200K in each sampled Hamiltonian using the anisotropic

Berendsen barostat [170] to determine the equilibrium box dimensions. The initial QM minimized configuration was then anisotropically expanded to match the average box shape from the NPT equilibration simulation and the centers of mass of each benzene were shifted proportionally. For the benzene I crystal in Designed-AMOEBA, the polymorph was equilibrated isotropically to 1 bar and 200K to maintain the AMOEBA unit cell shape. The shape of benzene I in AMOEBA09 was determined based on the benzene I lattice energy minima.

The resulting configuration from the above process for each polymorph was then energy minimized to a maximum force of < 50 kJ/mol/nm and modified such that all benzenes are planar and have all C-C and C-H bonds at their equilibrium length, thus ensuring that all intramolecular interactions are at their minimum value. This modification is carried out by using the 1, 3, 5 position carbons to define the plane of each benzene, and the new planar benzene is placed at the original center of mass. The resulting configuration is used as the restraint structure and initial structure for all subsequent NVT PSCP simulations.

The lattice energy minima for each polymorph in the AMOEBA09 potential was determined using the XTALMIN subroutine in the TINKER simulation package with a maximum RMS gradient per atom of 0.01. The corresponding lattice energy minima for the OPLS-AA and GROMOS54A7 potentials in the GROMACS simulation package were found using an energy minimization to a maximum force of < 2 kJ/mol/nm combined with a perturbation of the box vectors using the Nelder-Mead minimization algorithm as implemented in numpy version 1.8.2.

To compute the relative free energy as a function of temperature, all polymorphs in each potential were simulated in NPT over a range of temperatures between 10K - 250K evenly spaced by 10K. As real crystals approach low temperatures near 0K, the assumptions in the classical potentials will break down and quantum effects such as the zero point value energy will dominate the relative free energy differences. The low temperature classical MD simulation results are presented here to provide a complete comparison of the different classical models rather than as a rigorously accurate description of the physical system.

All point-charge polymorph states used for Hamiltonian reweighting are simulated using

the anisotropic Parrinello-Rahman barostat [174] and temperature coupled using stochastic dynamics [173] with time constants of 1000 ps and 1.0 ps, respectively. For benzene I in the Designed-AMOEBA potential, the simulations were run with the isotropic Parrinello-Rahman barostat to maintain the box shape from the AMOEBA09 Hamiltonian. The large time constant for the pressure coupling is necessary in these small systems to avoid instabilities in the volume oscillations. The volume decorrelation time in these systems is roughly 1 ns based on the peak distances in the volume time series.

All GROMACS trajectories were run for 20 ns, with the first 2 ns discarded for equilibration. One exception is the simulations in the PSCP to turn off the intermolecular interactions, which were run for 2 ns with the first 0.2 ns ignored because they converged significantly faster than the application of restraints. Each step of the PSCP for each polymorph and each potential was run with three independent trajectories and the resulting configurations from the three trajectories were combined during the MBAR free energy estimation. The NPT trajectories in AMOEBA09 were simulated with the TINKER 7.0 simulation package [137] for 5 ns (with the first 1 ns ignored for equilibration) with anisotropic Berendsen pressure coupling and Andersen temperature coupling [169] with a time constant of 2.0 ps and 0.1 ps, respectively. Long-range Lennard Jones interactions were calculated with an analytical dispersion correction and a cutoff distance of 1.5 nm. Long-range electrostatic interactions were treated with Particle Mesh Ewald summation [168] with a cutoff of 1.5 nm and a grid size of 16. The long-range cutoffs in TINKER were chosen to be larger than in GROMACS because the simulation package has no implementation of Particle Mesh Ewald for van der Waals.

The individual trajectories were decorrelated using the method presented by Chodera and others. [211–214] and uncertainties for each step were determined using the variance of 200 independent bootstrap samplings of the uncorrelated data.

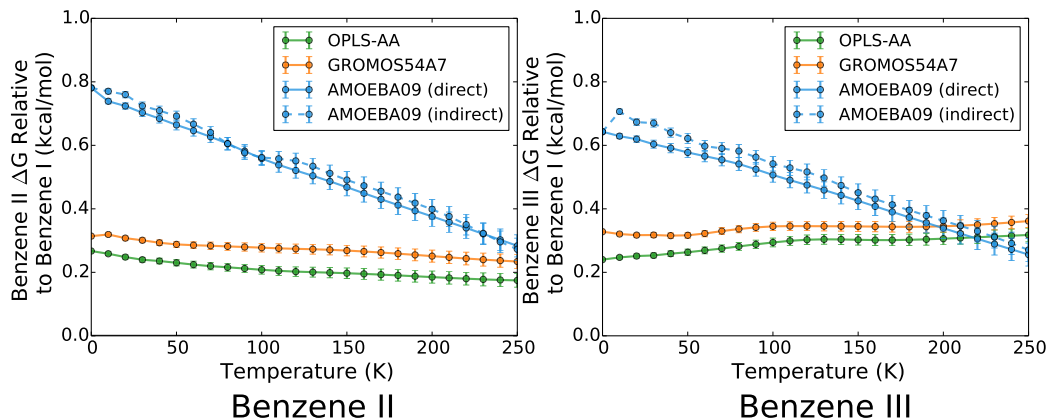


Figure 4.4: The OPLS-AA, GROMOS54A7, and AMOEBA09 potentials all predict the experimentally observed benzene I structure to be the most stable crystal at all temperatures between 0K and 200K. The metastable benzene II and benzene III structures are predicted to become more stable at higher temperatures due to higher entropic favorability. The added effect of polarization over the point-charge potentials is less significant at higher temperatures due to a compensation between the energy and entropy changes between potentials. The indirect AMOEBA09 free energy curve without target sampling is essentially the same as the curve with target sampling at all temperatures above 50K. The relative stability is shown for A) benzene II and B) benzene III.

4.4 Results and Discussion

We simulated the three benzene polymorphs shown in Figure 4.2 in the OPLS-AA, GROMOS54A7, and AMOEBA09 Hamiltonians, and the relative free energies of each polymorph between 0K-250K are shown in Figure 4.4. All free energies in Figure 4.4 were computed from equation 4.6 using the PSCP calculation at 200K. The OPLS-AA and GROMOS54A7 free energy estimates at 100K from this method are indistinguishable from a direct PSCP calculation at 100K as shown in Figure 4.5. The numerical free energy estimates along each step of the PSCP for OPLS-AA at 200K are shown in table 4.1. Tables and figures for all PSCP calculations in all potentials can be found in Appendix C sections C.4 and C.5.

In all three potentials, benzene I is identified as the most stable crystal structure with the

Step	benzene I	benzene II	benzene III
Add Restraints	-1.700 ± 0.014	-1.729 ± 0.010	-1.710 ± 0.008
Remove Interactions	-9.377 ± 0.000	-9.159 ± 0.000	-9.061 ± 0.000
PV Correction	0.326 ± 0.000	0.322 ± 0.000	0.325 ± 0.000
Total	-10.751 ± 0.014	-10.566 ± 0.010	-10.445 ± 0.008

Table 4.1: Polymorph free energy differences in OPLS-AA at 200K for each step in the PSCP path

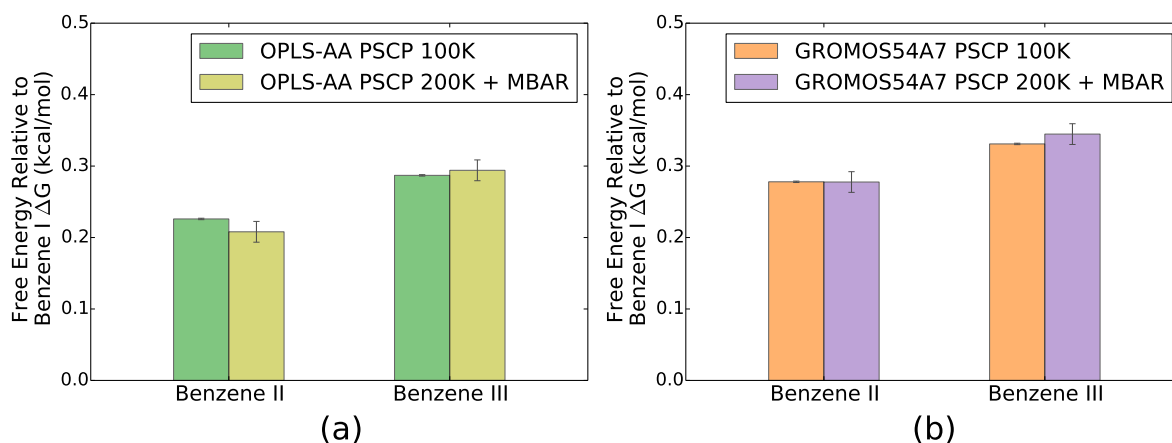


Figure 4.5: The relative free energy estimates of the polymorphs at 100K are statistically indistinguishable when calculated directly from the PSCP at 100K and from the PSCP at 200K combined with MBAR down to 100K.

lowest free energy over all temperatures examined. This is consistent with experimental observations of benzene I at ambient temperature and pressure [32, 76]. The metastable benzene II crystal has a higher entropy than benzene I as evidenced by the decrease in relative free energies at higher temperatures. This decrease in free energy at higher temperatures is consistent with experimental polymorph phase diagrams which show benzene II appearing at higher temperatures [144].

The nearly linear change in free energy indicates that the entropy differences between polymorphs are roughly independent of temperature in all potentials. The essentially constant entropy difference with temperature suggests that for this system the full free energy as a

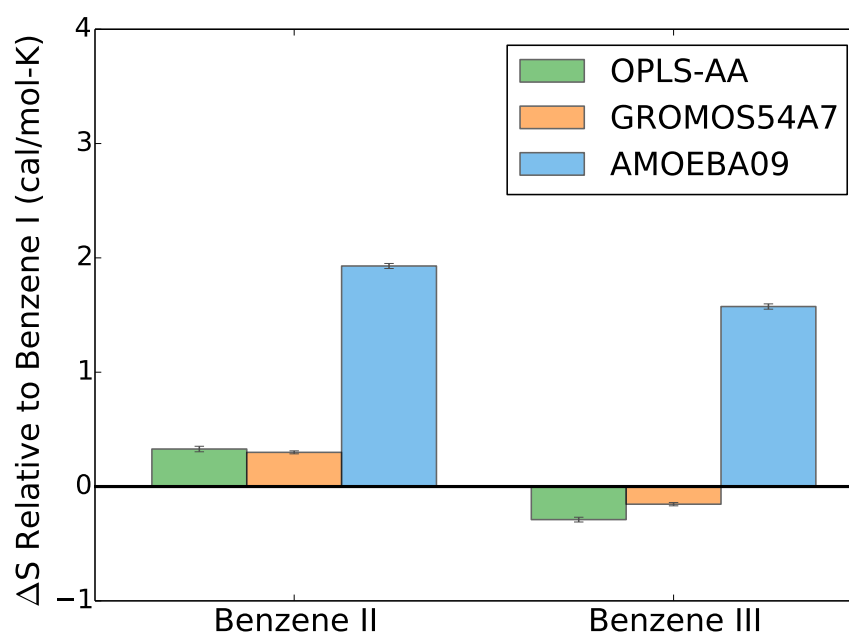


Figure 4.6: The metastable benzene II has a higher entropy than the globally stable benzene I structure in all three potentials. The entropy differences to benzene I in the polarizable AMOEBA09 potential are significantly larger in magnitude than in the two point-charge potentials. Uncertainties in the entropy were calculated by bootstrapping the individual points used to regress the linear fit.

function of temperature could in principle be well approximated using the lattice energy difference at 0K and the entropy difference at a single finite temperature. The relative entropies for benzene II and benzene III are shown in Figure 4.6 using the slope of a linear fit to the free energy curves. In all three potentials, benzene II is observed to have a higher entropy than benzene I. However, the entropy difference in the polarizable AMOEBA09 potential is significantly larger than the entropy difference in the two point-charge potentials.

The observed effect of including polarizability on the relative free energy estimates is highly dependent on the temperature. The 0K lattice energies and 200K free energies of all three polymorphs in all three potentials are shown together in Figure 4.7. In the limit of 0K, the relative stability of benzene II and benzene III decreases by roughly 0.52 and 0.40 kcal·mol⁻¹, respectively, when moving from the point-charge OPLS-AA potential to the polarizable AMOEBA09 potential. This decrease in stability is noticeably larger than the corresponding change of 0.05 and 0.09 kcal·mol⁻¹ when moving from the OPLS-AA potential to the point-charge GROMOS54A7 potential. At 250K, the change in free energy of the metastable polymorphs is 0.08 and 0.12 kcal·mol⁻¹ when moving from the OPLS-AA to AMOEBA09 potential. This decrease is closer to the decrease when moving to GROMOS54A7 of 0.06 and 0.04 kcal·mol⁻¹. The relatively minor change in free energy between point-charge and polarization potentials at 250K despite the large difference in the lattice energy minima suggests that the increased energetic favorability in the polarizable potential comes with a corresponding decrease in entropic favorability. This lower entropy decreases the difference between potentials as temperature increases leading to similar relative free energies in the three potentials at 250K despite large differences in the lattice energy at 0K.

An interesting result from the benzene system is that the entropy difference between benzene II and benzene III are essentially the same in the point-charge and polarizable potential. The relative free energies of the three polymorphs using benzene II as the reference structure are shown in figure 4.8. The free energy against temperature for Benzene III is essentially the same in OPLS-AA and AMOEBA09 except that the y-intercept has changed. This indicates

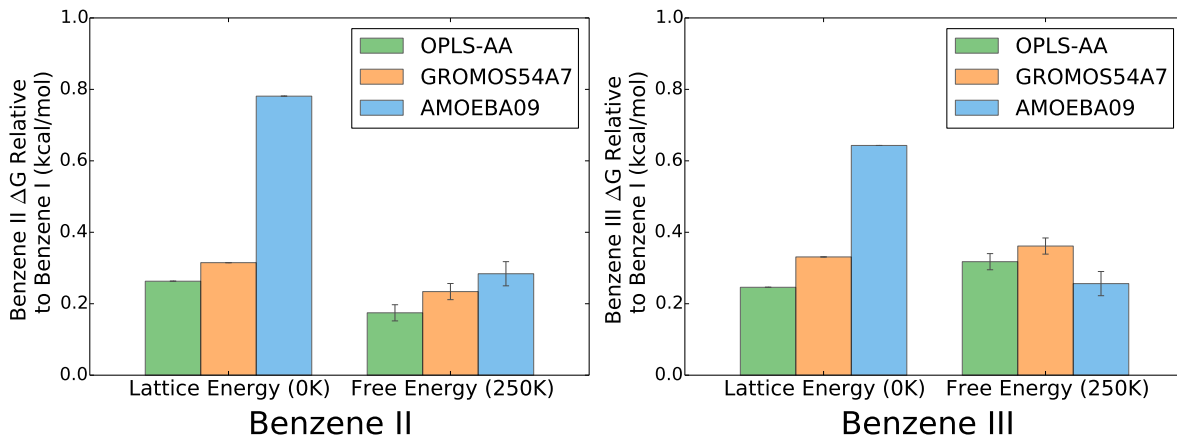


Figure 4.7: The influence of polarization in relative polymorph stability is smaller for 250K free energy differences than 0K lattice energy differences. The free energy difference in the polarizable AMOEBA09 at 250K potential is more closely approximated by the 250K free energy difference in the point-charge OPLS-AA potential than by the minimized lattice energy difference in the AMOEBA09 potential. The relative stability is shown for A) benzene II and B) benzene III for the OPLS-AA, GROMOS54A7, and AMOEBA09 potential using both full MD free energy and minimized lattice energy.

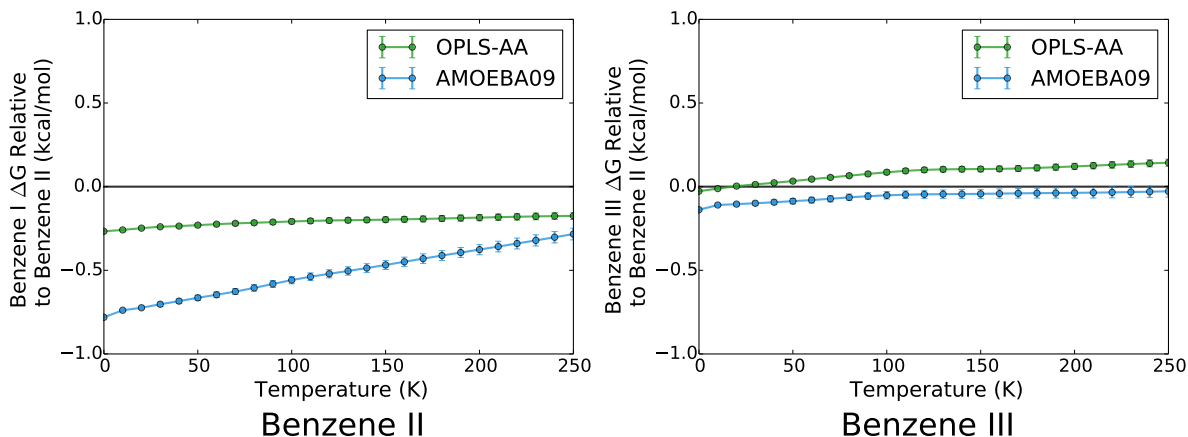


Figure 4.8: The free energy of benzene III relative to benzene III as a function of temperature is essentially the same shape in OPLS-AA and AMOEBA09 except for a constant offset. The free energy of benzene I relative to benzene II is noticeably different in OPLS-AA and AMOEBA09 due to the difference in unit cells of benzene I in the two potentials.

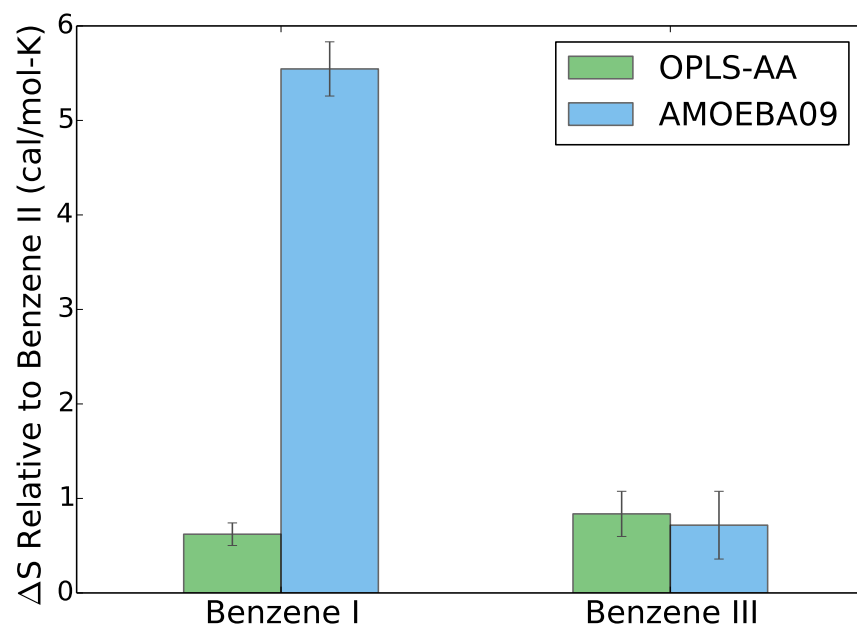
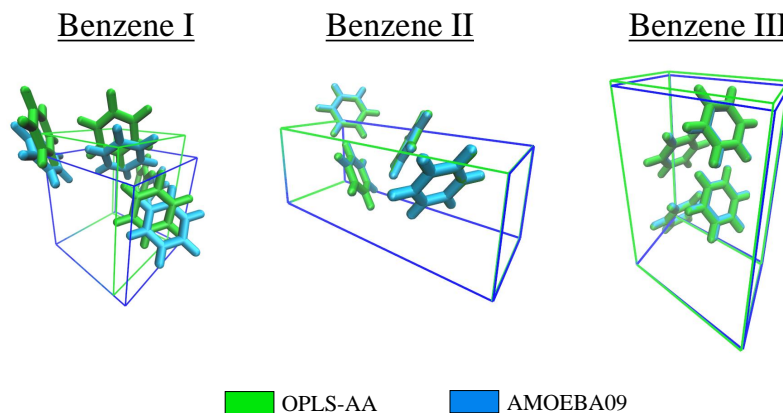


Figure 4.9: The entropy difference between benzene III and benzene II are statistically indistinguishable in OPLS-AA and AMOEBA09. The entropy difference between benzene I and benzene II is significantly different in OPLS-AA and AMOEBA09 due to the difference in unit cells of benzene I in the two potentials.



18

Figure 4.10: The unit cells for benzene II and benzene III are essentially indistinguishable in OPLS-AA and in AMOEBA09. The unit cell for benzene I is significantly different in OPLS-AA and in AMOEBA09.

that the estimated enthalpy difference has changed when moving to AMOEBA09 but the entropy remains essentially the same. Indeed, the entropy difference between benzene II and benzene III are statistically indistinguishable in the two potentials (shown in figure 4.9).

The similar entropy difference between Benzene II and benzene III in the two potentials is in stark contrast with that between benzene II and benzene I. Both the enthalpy and the entropy difference to benzene I change significantly when moving from OPLS-AA to AMOEBA09. Benzene I undergoes significant structural changes (and space group changes) when comparing the lattice minima for OPLS-AA to AMOEBA09 as shown in 4.10. Benzene II and benzene III, however, remain essentially the same. It is likely that this change in structure leads to the noticeably different entropy relative to benzene II. Therefore, we conclude that the entropy differences appear to be less affected than the enthalpy differences in cases where the unit cell and space group are essentially the same. However, in cases where the lattice minima and the space group change when going to a more expensive potential, both the enthalpy and entropy differences change significantly.

The relative free energies of benzene II and benzene III in the polarizable AMOEBA09

potential at higher temperatures are more closely approximated by the relative free energies in the point-charge OPLS-AA potential than by the minimized lattice energy difference in the AMOEBA09 potential as can be seen in Fig. 4.7. This result highlights the potential importance of computing full free energies directly rather than approximating it by the lattice energy when estimating finite temperature polymorph stability, even if only lower-level energy functions are computationally feasible. This is, of course, a result for only one molecule. A further rigorous assessment of the effect of polarization across a large range of chemical functionalities will be required to form a more general conclusion.

Reweighting from the designed point charge potential to the polarizable potential is an almost entirely sufficient replacement for direct sampling in the polarizable potential. The relative free energy curves of the benzene crystals in the AMOEBA09 potential are nearly identical when including or excluding target state sampling in the computation. The curves with and without 5 ns of AMOEBA09 sampling are shown as the ‘direct’ and ‘indirect’ pathway, respectively, in Figure 4.4. The indirect curve for benzene II falls within the uncertainty of the direct curve for all points above 50K. For benzene III, there is a slight bias in the indirect pathway at all temperatures. However, this bias is well below the free energy difference to the most stable benzene I structure and is not large enough to reorder the polymorph stability ranking at any temperature. We hypothesize that the bias appears due to a difference in box shape between the Designed-Amoeba potential and AMOEBA09 in the more complex monoclinic crystal shape of benzene III. The overlap in the box shapes could in principle be improved by running the designed potential isotropically with the finite temperature AMOEBA09 box shape or by using an applied stress tensor, though such approaches are not essential for this system. We note that although all 5 ns of AMOEBA09 sampling was used to generate the direct curve in Figure 4.4, as little as 2ns of simulation is sufficient to converge the free energies in this system. See Appendix C section C.2 for more details.

The accurate free energy estimates without target sampling shown Figure 4.4 are computationally feasible because of the high overlap between the designed potential and the target

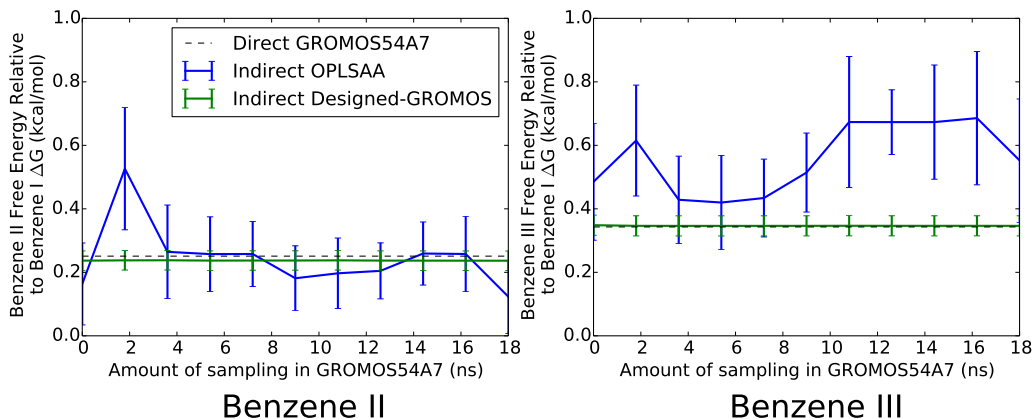


Figure 4.11: The relative free energy estimates for benzene II and benzene III using reweighting along the indirect pathway are significantly more accurate when using the Designed-GROMOS potential rather than OPLS-AA and do not require any sampling in the GROMOS54A7 potential to converge. The dashed line indicates the free energy difference estimate using direct sampling in GROMOS54A7.

Hamiltonian. The same indirect free energy estimates are not feasible with the OPLS-AA potential because of poor overlap to the target state. We examine in more detail the reweighting approach using OPLS-AA and the designed potentials as the sampled state and using both GROMOS54A7 and AMOEBA09 as the target potentials.

The two point-charge potentials OPLS-AA and GROMOS54A7 have relatively poor overlap at 200K leading to a high uncertainty and bias in the indirect free energy difference estimate between the three polymorphs of benzene in the GROMOS54A7 potential when sampling with OPLS-AA (Figure 4.11). A significant amount of sampling is needed in the target GROMOS54A7 potential to converge the free energy difference estimate in the indirect pathway when connecting the polymorphs through the OPLS-AA potential.

The poor overlap between OPLS-AA and GROMOS54A7 in this system is caused primarily by differences in the equilibrium benzene geometry and differences in the intramolecular parameters in the two potentials. The overlap to GROMOS54A7 is significantly increased when sampling with the Designed-GROMOS potential which has intramolecular parameters that

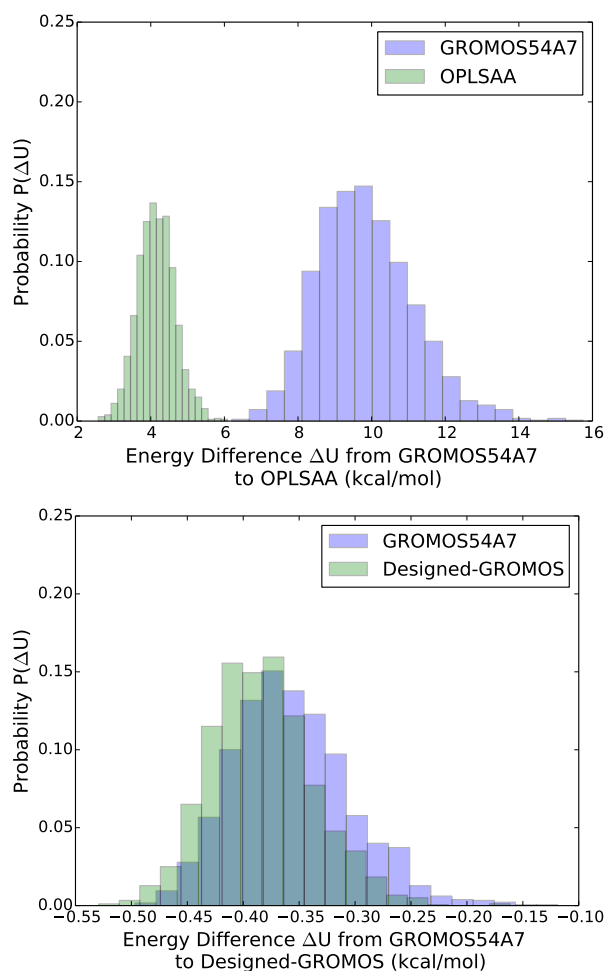


Figure 4.12: The OPLS-AA potential has relatively poor overlap with GROMOS54A7, while the Designed-GROMOS potential has higher overlap with the GROMOS54A7 potential in the solid benzene system. The measure of overlap, O_{ij} , from OPLS-AA to GROMOS54A7 using equation 4.3 is 0.000095 while the overlap from Designed-GROMOS to GROMOS54A7 is 0.391. The distribution sampled from AMOEBA09 is shown in blue and the distribution from the point-charge potentials is shown in green.

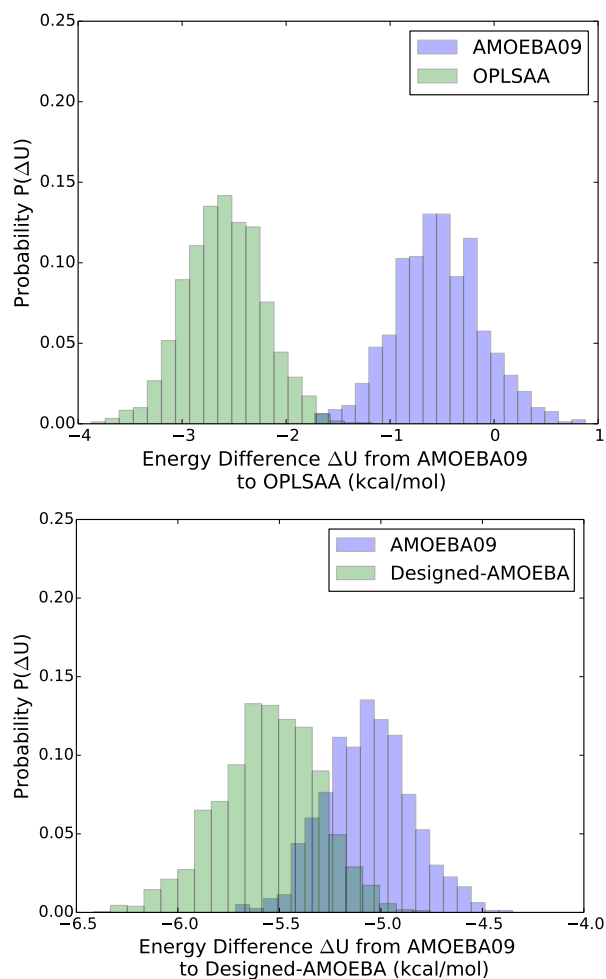


Figure 4.13: The OPLS-AA potential has relatively poor overlap with AMOEBA09, while the Designed-AMOEBA potential has higher overlap with the AMOEBA09 potential in the solid benzene system. The overlap from OPLS-AA to AMOEBA09 is 0.017 while the overlap from Designed-AMOEBA to AMOEBA09 is 0.178.

match those in GROMOS54A7 potential. Figure 4.12 shows the overlap in the energy difference to GROMOS54A7 using both OPLS-AA and Designed-GROMOS as the sampled point-charge potential. The overlap in the configuration space of benzene I increases from 0.00032 to 0.437 when using the Designed-GROMOS potential instead of the OPLS-AA potential. A comparison of numerical overlap for each polymorph in both potentials as well as histograms of energy difference can be found in Appendix C section C.1).

The free energy estimates of the polymorphs converge even without the samples collected directly in GROMOS54A7 as a direct result of the higher overlap to GROMOS54A7 with the designed potential (Figure 4.11). In addition, the uncertainty in the free energy difference between the sampled and target potential is over an order of magnitude smaller for Designed-GROMOS, showing that reweighting to GROMOS54A7 from the Designed-GROMOS potential is 2-3 orders of magnitude more efficient than reweighting with the OPLS-AA potential. This assumes we are in the limit that the squared uncertainties are inversely proportional to the number of reweighted configurations. The bias introduced by using the indirect free energy pathway with Designed-GROMOS to connect the polymorphs was indistinguishable from zero relative to direct simulation in the GROMOS54A7 potential. The bias was determined by connecting the polymorphs directly in the GROMOS54A7 potential and comparing the results to the indirect calculation through the Designed-GROMOS potential (Figure 4.11).

Reweighting from a point-charge potential to the AMOEBA09 polarizable model adds an extra degree of complexity to the reweighting process because the charges in the sampled and unsampled state are no longer identical as in the case of OPLS-AA and GROMOS54A7. However, the overlap between the two potentials is again observed to be most sensitive to differences in the intramolecular parameters. The overlap and histograms of the energy differences to AMOEBA09 are shown in Figure 4.13 for benzene I when both OPLS-AA and Designed-AMOEBA are used as the sampled point-charge potential for the indirect reweighting scheme. The overlap between the OPLS-AA and AMOEBA09 potential is relatively low at $O_{ij} = 0.0037$ due to differences in the intramolecular parameters. The overlap increases to 0.178 when using the Designed-AMOEBA potential which has intramolecular parameters matching those in

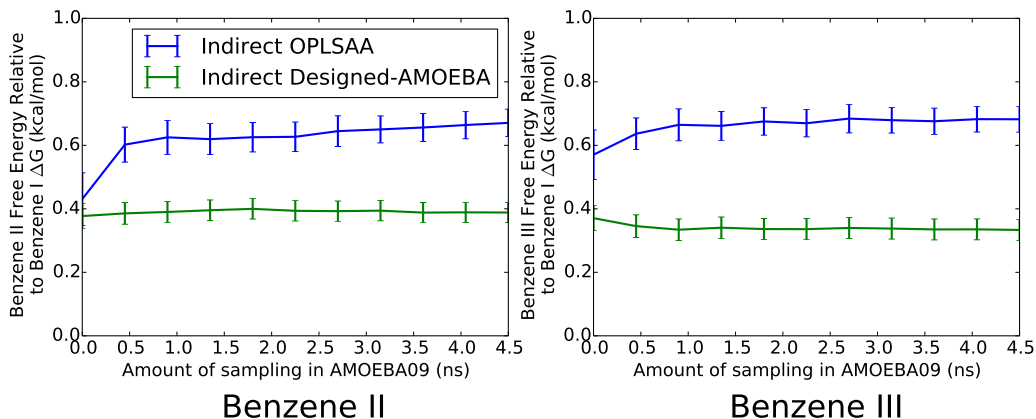


Figure 4.14: The relative free energy estimates for benzene II and benzene III along the indirect pathway are more accurate when using the Designed-AMOEB A potential instead of OPLS-AA and do not require any sampling in the AMOEBA09 potential to converge. The free energy estimates with OPLS-AA have a substantial bias with and without AMOEBA09 sampling as a consequence of the benzene I box shape being different in the point-charge and polarizable potentials.

the AMOEBA09 potential. The uncertainties in the free energy differences between the polymorphs using the indirect pathway with the Designed-AMOEB A potential are roughly $0.03 \text{ kcal}\cdot\text{mol}^{-1}$ which is an order of magnitude lower than the free energy differences between the polymorphs in the AMOEBA09 potential. Although the uncertainties along the indirect path with OPLS-AA are also lower than the polymorph free energy differences, a substantial bias of roughly $0.3 \text{ kcal}\cdot\text{mol}^{-1}$ appears in the free energy estimate with and without AMOEBA09 sampling as a direct consequence of the difference in benzene I box shape between the point-charge and polarizable potential. This bias further emphasizes the need to utilize information from the target potential when reweighting from a cheap to an expensive Hamiltonian as has been done here with the Designed-AMOEB A potential.

The small remaining lack of overlap between the Designed-AMOEB A and AMOEBA09 potential can be further reduced by tuning the charges in the designed point-charge potential. The overlap between the two potentials is already sufficient in this system to estimate the benzene free energies without sampling in AMOEBA09 (Figure 4.14). However, adjusting the

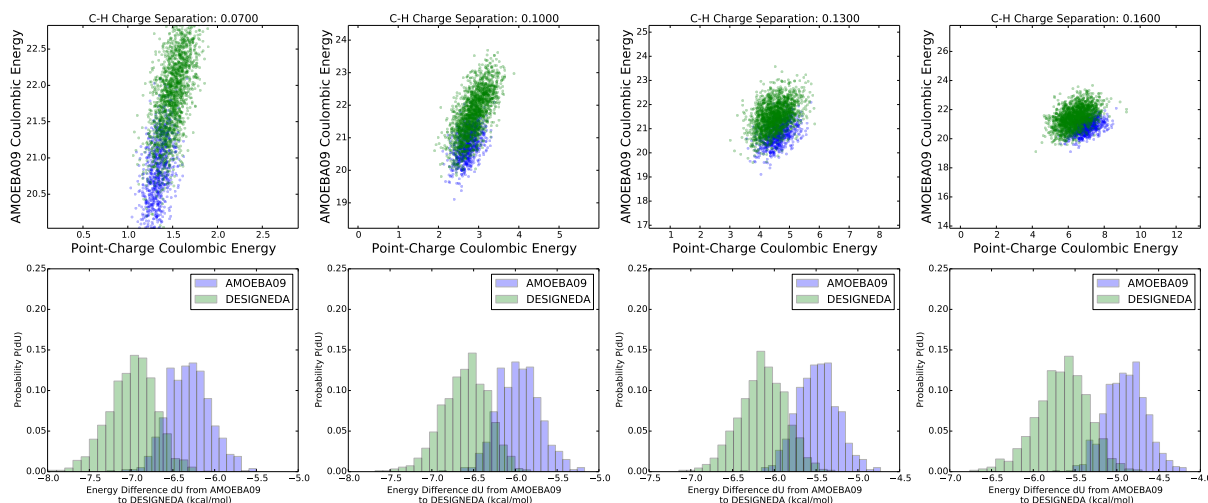


Figure 4.15: The configuration space overlap between the Designed-AMOEBA and AMOEBA09 potential can be increased by tuning the charges on the atoms in the benzene molecules in the designed potential. In the Benzene II crystal the highest overlap occurs at a hydrogen charge $q_H = 0.100$ (and a corresponding negative carbon charge of -0.100), compared to the original $q_H = 0.115$ in OPLS-AA. The distribution sampled from AMOEBA09 is shown in blue and the distribution from the designed point-charge potential is shown in green.

charges may be critical in creating a potential with sufficient overlap in systems with more polarizable molecules. Because of the symmetry and net neutrality of the benzene molecule, the complete charge distribution can be described in a single dimension, the positive charge q_H on all hydrogen atoms, since changes in q_H require a unique change in the q_C to maintain neutrality. The overlap and energy difference histograms for Benzene II are shown in Figure 4.15 for 4 different q_H values ranging from 0.070 - 0.160. The highest overlap in the Benzene II structure occurs with a $q_H = 0.100$ as opposed to the OPLS-AA value of 0.115 which has the highest overlap for Benzene I and Benzene III. This suggests that the point-charge potential in the Benzene II structure slightly overpredicts the average magnitude of the charges that are on the benzene atoms with the polarizable AMOEBA09 Hamiltonian. A rigorous optimization of the point-charges to maximize overlap with AMOEBA09 is outside the scope of the present work, and perturbing the charges will be more complicated in larger molecule systems because there are more than one dimension of charge defining the electrostatic interactions. However, the increase in overlap resulting from perturbing the point charges suggests that it may be possible to explore the configuration space of polarizable systems using a properly selected point-charge potential or a combination of multiple different point-charge potentials.

We can probe the convergence of the free energy differences of the benzene polymorphs by comparing the results obtained with different values of the hydrogen charge. Because free energy is a state function, the free energy difference estimate between two states will necessarily be identical for two independent thermodynamic pathways connecting the states if sufficient overlap of the relevant configuration space has been achieved in both paths. This charge perturbation technique was previously used by Essex and co-workers to test the convergence of reweighting from an MM to QM/MM Hamiltonian in protein-ligand binding systems [247]. The relative free energies of the benzene polymorphs are shown in Figure 4.16 using a range of different hydrogen charge values in the indirect thermodynamic path with Designed-AMOEBA. The free energy difference estimates are statistically indistinguishable for all hydrogen charges between 0.070-0.160 with no bias and little change in the uncertainty. Outside of this charge range, a bias begins to appear in the free energy estimates as a result

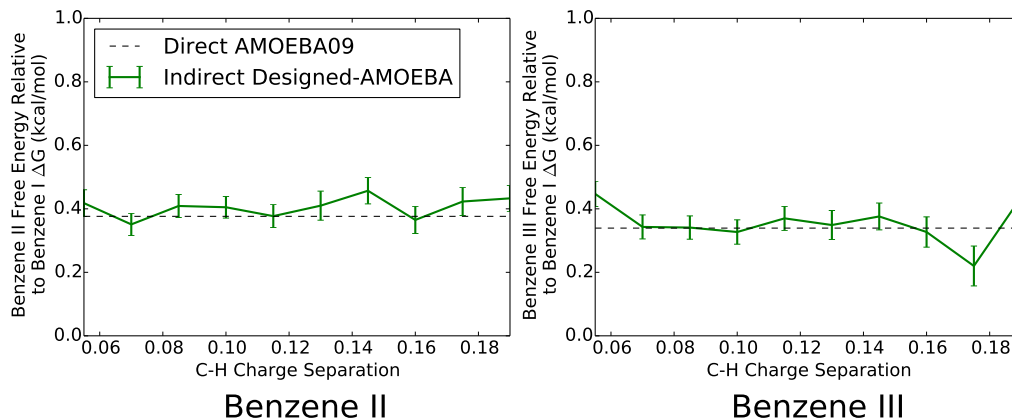


Figure 4.16: The relative free energy difference estimates of the benzene polymorphs are statistically the same for all values of hydrogen charge between 0.070-0.160. In all estimates, the PSCP value is used with a charge of 0.115 and the end states are first reweighted to the new charge in the Designed-AMOEBA potential and then reweighted to the AMOEBA09 potential.

of the box volume not having good overlap in the point-charge and polarizable potentials. There is no method to absolutely ensure that the relevant configuration space of an unsimulated target Hamiltonian has been fully sampled, however the path independence of the free energy estimate (Figure 4.16) as well as the insensitivity of the estimate to adding direct sampling in AMOEBA09 (Figure 4.14) strongly suggests that the reweighted free energy difference estimates presented in Figure 4.4 have converged to the correct values.

The advantage of using the indirect thermodynamic path depicted in Figure 4.1 to calculate the differences in polymorph free energies is that no sampling needs to be done in the expensive polarizable potential. The same AMOEBA09 free energies presented in Figure 4.4 could, in theory, be calculated by direct brute-force sampling in the polarizable potential. The expense of running these simulations in AMOEBA09 can be approximated by assuming that the same total amount of sampling must be done along the PSCP to achieve sufficient convergence. The estimated expense of simulating in AMOEBA09 and the expense incurred by using

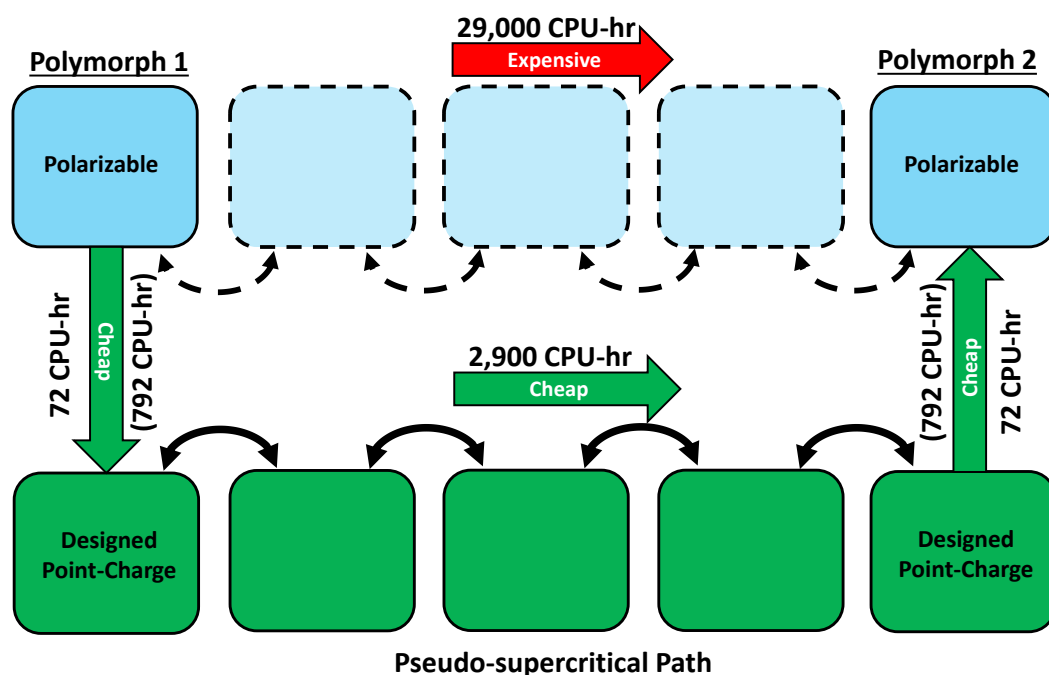


Figure 4.17: The total expense of calculating the free energy difference of polymorphs in the AMOEBA potential is cheaper when using the indirect pathway through a point-charge potential rather than using direct sampling in AMOEBA09 by roughly a factor of 10. The direct pathway to connect polymorphs using the pseudo-supercritical path in AMOEBA requires approximately 29,000 cpu-hours. The indirect pathway using the designed-point charge potential requires approximately 3,000 cpu-hours. The indirect pathway using end state sampling in AMOEBA (shown in parenthesis) requires roughly 4,500 cpu-hours and is still cheaper than the direct pathway because it avoids the need to run the intermediate states of the pseudo-supercritical path in the AMOEBA09 Hamiltonian.

the indirect approach with Designed-AMOEBA09 are shown in Figure 4.17. The indirect pathway is cheaper by roughly an order of magnitude, which reflects the fact that direct sampling in point-charge potentials is 10x faster than sampling in polarizable potentials. In quantum-based Hamiltonians, such as QM/MM, DFT, or MP2, where the sampling is multiple orders of magnitude slower, the cost savings of using this indirect pathway rather than direct sampling will increase proportionally. It is worth noting that although sampling in AMOEBA09 is only 10x slower than non-polarizable force fields, the current implementation of AMOEBA09 cannot be efficiently parallelized onto multiple cores, meaning that the overall wall-clock time can be significantly more than 10x longer for an AMOEBA09 simulation relative to an efficiently parallelized non-polarizable model.

Finally, it is worth noting that the indirect pathway presented here would still be faster than the direct sampling method even if some sampling in the expensive target potential is necessary at the end states to converge the indirect free energy differences. The most expensive step in the overall indirect pathway involves running the PSCP calculations between the two polymorphs. The NPT simulations for each polymorph (used to calculate the free energy difference between the cheap and expensive Hamiltonian) require only a single simulation at each temperature and represent a small fraction of the overall computational expense. Although the overlap between the point-charge and AMOEBA09 potential is high in the benzene system, it is conceivable that other systems may require some small amount of sampling in the AMOEBA09 potential in order to converge the free energy estimates. However, the total cost of the indirect pathway including end state sampling in AMOEBA09 is still considerably lower than the direct pathway because it avoids the expensive AMOEBA09 sampling along every intermediate state in the PSCP.

4.5 Conclusions

The relative Gibbs free energies of three polymorphs of benzene were calculated between 0K-250K at ambient pressure in the point-charge OPLS-AA and GROMOS54A7 potentials as well

as the polarizable AMOEBA09 potential to determine the added effect of including polarizability in Hamiltonians used for polymorph free energy calculations. In all three potentials, the experimentally observed benzene I crystal was correctly identified as the most stable crystal structure with the lowest free energy at all temperatures examined.

The addition of multipoles and polarization to the electrostatic description were found to have a noticeable influence on the relative stability of the crystal structures at 0K. In addition, the unit cell shape of benzene I changes substantially when simulating in the polarizable AMOEBA09 potential rather than the point-charge potentials. However, the effects of polarization were less noticeable at higher temperatures due to an apparent compensation effect between the energy and entropy change when polarization is included. As a result of this compensation, the free energies in all three potentials are similar near room temperature. Indeed, the free energies in the polarizable AMOEBA09 potential are more closely approximated by the free energies in the point-charge potential than by the lattice energies in the polarizable potential. This result motivates the use of full molecular dynamics approaches when estimating the free energy of finite temperature polymorphs. However, the generalization of the effect polarization and thermal motion on small molecule organic crystal free energy landscapes requires additional sampling in other molecular systems.

The free energies in the AMOEBA09 potential were calculated indirectly by simulating in a cheaper point-charge potential and Boltzmann reweighting the resulting configurations, thus avoiding the need to directly simulate within the more expensive polarizable potential. The Boltzmann reweighted path connecting the polymorphs reduces the total simulation cost by roughly a factor of 10 by interconverting the polymorphs in the cheaper MM potential rather than in the polarizable model. The reweighted pathway has no loss of accuracy relative to direct evaluation of polymorph free energies in the target potential based on similar free energy estimates with and without direct sampling in the target potential.

Previous investigations employing the Boltzmann reweighting process were hindered by a lack of overlap in the intramolecular degrees of freedom between the point-charge and polarizable potentials [247,253]. The necessary overlap to carry out the reweighting process between

the point-charge and polarizable AMOEBA09 force-field is created by sampling from point-charge potentials which contained intramolecular parameters matching those in AMOEBA09 Hamiltonian. This procedure, which requires no parameterization or optimization, substantially increases the overlap to the AMOEBA09 potential by over 2 orders of magnitude and removes the need to sample directly in the polarizable potential, which would be necessary to reweight directly from OPLS-AA to AMOEBA09.

We further observe that the overlap can be improved in some systems by adjusting the charges in the benzene molecule within the designed point-charge potential. This step is not necessary in the current benzene system due to the relatively small polarity. However, the approach presented in this work has been developed in order to efficiently elucidate the molecules where polarization strongly influences polymorph stability. In highly polar systems, adjusting the point-charges may be critical to sample the correct configuration space.

Although the Boltzmann reweighted pathway was applied here for polymorph evaluation in organic solids, the methodology of creating point-charge potentials to overlap with expensive Hamiltonians can be straightforwardly applied to other free energy estimation problems in polarizable potentials such as solvation free energies or ligand binding free energies. Additionally, we anticipate that the computational savings will be larger when reweighting to quantum-based Hamiltonians because the savings are roughly proportional to the difference in expense between the designed and target potential.

Chapter 5

Lattice Minima Interconversion at Ambient Temperatures

5.1 Introduction

Computational crystal structure prediction studies search for minima on the lattice energy surface in order to predict the experimental crystal form [11]. The energy-minimized structure with the lowest energy is considered the most stable polymorph. All other lattice minima within a modest energy of the global minima (usually around 0.5 kcal at ambient conditions) are considered to be potential metastable polymorphs.

Static descriptions of crystalline materials, such as lattice energy minimization protocols, do not account for the thermal motions that exist in all solids at ambient conditions. One direct impact of this thermal motion is entropic contributions to the finite temperature stability. Entropic differences between crystal structures account for the emergence of temperature-mediated crystal restructuring presented in Chapter 3.

The static approaches to modeling crystal structures also neglect all kinetic aspects of finite temperature crystalline solids. At room temperature, crystals visit many configurations other than the precise lowest-energy configuration as a result of thermal motion. Crystals near the global lattice minima may ultimately visit multiple other local minima on the lattice energy landscape if the transition barriers between basins are sufficiently small.

CSP studies of solids based on lattice energy minima are well known to suffer from polymorph overprediction [11,202,268]. These CSP studies identify many more low-energy lattice minima than there are experimentally observed crystal structures. A frequent proposal put forward to explain this overprediction is that many of these low-energy lattice minima have small transition barriers to each other and are therefore not all unique structures at ambient temperature. For example, a recent CSP study of pyridine produced three similar lattice minima which the authors argued would interconvert into a single ambient temperature crystal form [269].

Lattice minima exploring other minima at ambient conditions can occur in two different cases: Lattice minima *conversion* and lattice minima *interconversion*. Lattice minima *conversion* occurs when a high-energy metastable form escapes from its local basin into a significantly more stable set of configurations. These transitions will typically be irreversible and monotropic such that the metastable form will convert into a more stable lattice minima upon heating and cooling. However, the more stable lattice minima will not transition back into the original lattice minima if additional heating and cooling cycles are performed.

This type of lattice minima interconversion was first demonstrated for acetic acid using so called ‘short MD shakeup’ simulations [270]. In this study, the lattice minima from a CSP study were run with a short 60 ps MD simulation to allow unstable lattice minima to restructure. The investigators observed that less unique lattice minima remained after the MD simulation because higher energy minima transformed and cooled into nearby more stable minima. This type of lattice minima conversion was also observed in many of the systems shown in Chapter 3, and those will be presented later in this chapter.

A second type of lattice transition is lattice minima *interconversion*. Lattice minima interconversion occurs when a lattice minima escapes its local basin and explores another locally stable set of configurations with similar energies. Because the lattice minima in this scenario have similar stabilities, a single trajectory will explore all of the nearby local basins through successive escaping and entering of the basins regardless of the initial structure. This minima interconversion differs from the irreversible minima conversion discussed above because

any heating and cooling cycle starting from one lattice energy minima is likely to result in a different energy minimized structure.

A common example of lattice minima interconversion is the phenomena of dynamically disordered crystals, also referred to as ‘plastic’ or ‘rotator’ phases. In these plastic phases, the molecules have 3-dimensional translational order and rotational *disorder* around their crystallographic positions [67]. Each degenerate structure of these disordered crystals represents a stable lattice minima. Molecules that form plastic crystals typically have a spherical or disk-like shape. Examples include caffeine [271], cyclopentane [100], ethanol [272], succinonitrile [209], methane [67], C₆₀ [4], and alkali hydroxides [273]. These molecules at ambient temperature will maintain an average spatial position within the crystal lattice, but will freely rotate about one or more axes due to the near-symmetric nature of the electrostatic structure. Each of the continuously visited and rotationally-related configurations corresponds with a distinct lattice energy minima.

Both lattice minima conversion and interconversion have been scarcely studied in molecular models because they require full molecular dynamics simulations. Entropy differences between crystals can in some cases be captured with a quasi-harmonic approximation. Therefore, effects such as temperature-mediated restructuring can in certain systems be detected without molecular dynamics, as shown for select rigid molecules in Chapter 3. However, a quasi-harmonic approximation considers all lattice minima as unique free energy minima and cannot show when two energy minima correspond with configurations from the same polymorph ensemble [31]. Therefore, the most common molecular models of crystals cannot be used to search for minima interconversion or the rise of disordered crystals. Indeed, Desiraju recently stated “A drawback of the presently available molecular crystal structure modeling and prediction software is that we are unable to simulate disordered structures. There is a definite need to be able to handle this lacuna”. [274].

Molecular dynamics has the ability to detect minima conversion, interconversion, and the emergence of disordered crystals. In this chapter, we demonstrate minima conversion in a number of different organic crystals using molecular dynamics. The unstable minima that

convert in this way can be viewed as ‘dead-end’ minima that will not be uniquely observed experimentally. This effect of unstable minima being eliminated with molecular dynamics has the potential to substantially streamline computational crystal structure prediction.

The ultimate challenge for molecular dynamics models is to show which of the lattice energy minima from a CSP convert into another minima at ambient conditions and therefore do not correspond with a unique free energy minima. An example of minima from a CSP of resorcinol are discussed at the end of this chapter. However, the process of heating up the structures and determining the fraction that convert is left for future research.

5.2 Methodology

Simulation trajectories that convert to the same crystal ensemble will have identical averaged properties at equilibrium. One could, in theory, use any averaged system property as an order parameter to show that the ensembles are equivalent. However, many averaged system properties could yield false positives for matching ensembles. For example, the equilibrium bond length of a particular bond in a rigid molecule may be indistinguishable between two polymorphs that are otherwise markedly different. Therefore the bond length would be a poor choice of order parameter to distinguish polymorphs.

In this work, we use the average energy, average box vectors, and free energies to determine whether two or more simulation trajectories have converged to the same equilibrium ensemble. In essence, two simulation are considered to be the same polymorph if they have an indistinguishable equilibrium average energy, box vectors, and free energy.

Both minima conversion and interconversion will result in two lattice energy minima producing a single crystal ensemble at ambient conditions. The two types of transformations are distinguished in each case by quenching the ensemble down to 0K. More specifically, the set of saved configurations from the high temperature trajectory are all subjected to crystal minimization down to 0K. The configurations resulting from an irreversible minima conversion

will all freeze down to the single most stable initial lattice energy minima. Lattice minima interconversion, by contrast, will produce a mixture of different lattice energy minima after the trajectory is quenched.

The 12 polymorphic systems presented in Chapter 3 are used as the test systems to explore polymorph (inter)conversion. The initial structures were taken from the Cambridge Structural Database and expanded to supercells with dimensions larger than 1.6 nm as described earlier.

Each of the initial crystal structures were heated to a range of temperatures from 10K up to 300K. The average energies and volumes of each polymorph along the temperature range were examined to identify where multiple trajectories produce the same averaged properties. In the cases where two or more structures converged to the same equilibrium ensemble, the configurations from the 300K trajectories were each crystal minimized to determine which type of transformation occurred.

Minima conversions are identified by annealing the energy minimized structures from 0K up to 300K and then reminimizing the resulting trajectory. If the lattice minima has not escaped into a new basin, the initial and final structure from this annealing process will correspond with the same energy minimized configuration. In the instances where the simulation has escaped from a metastable form into a new energy basin, the reminimized configuration will be different than the original starting structure.

5.3 Minima Conversion

A number of the minimized experimental crystal structures from the CSD underwent an irreversible monotropic transformation into a new crystal form upon heating to high temperatures, which we refer to as ‘minima conversion’. The restructuring events appear as sharp discontinuities on plots of both average energy and average volume against temperature. An

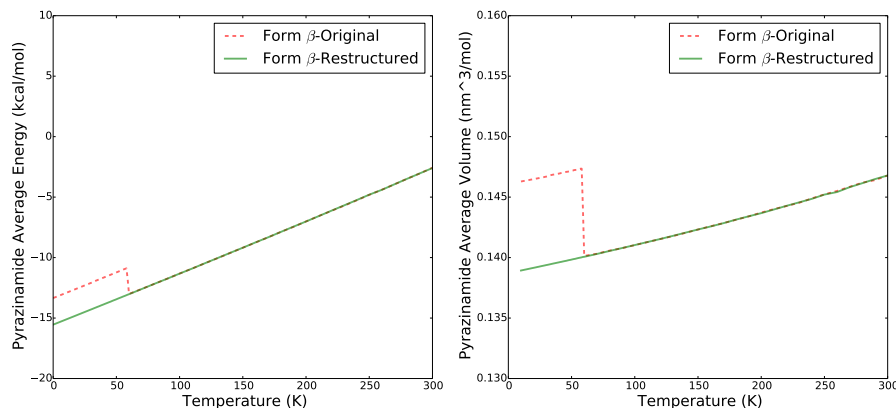


Figure 5.1: Pyrazinamide β restructures to a new crystal form between 50 K and 60 K. The forms have indistinguishable average energy and unit cell volume above the restructuring temperature.

example of this discontinuity is shown for pyrazinamide β in Figure 5.1. The minimized structures for the original and restructured forms of pyrazinamide β are shown in Figure 5.2. Similar transformations occurred for two polymorphs of tolbutamide, chlorpropamide, and aripiprazole. The average energy and average volume for the original and restructured forms of those compounds are shown in Figures 5.3–5.5.

The new crystal structures formed from the restructuring process at high temperatures were found to be stable and do not spontaneously restructure back to the original form at any temperature between 0 K and 300 K. The minimized lattice energies of the restructured forms are significantly lower than that of the initial lattice minima in all of these cases. The reminimized configurations at 300 K from both trajectories therefore all collapsed to the single significantly more stable lattice minima.

It is extremely unlikely that the original and restructured lattice energy minima each correspond to a different experimentally observable structure at ambient conditions. Polymorphs that are metastable at all temperatures can in some cases be observed experimentally as long as the kinetic barriers to restructure into a more stable form are prohibitively high. A rigorous estimate of the restructuring kinetics in bulk crystal is outside the scope of this study because

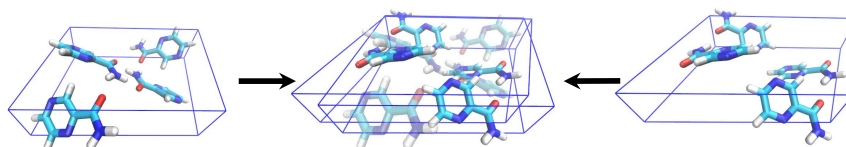
Form II (Original)Form II (Restructured)

Figure 5.2: Pyrzin form II restructures at temperatures above 60K. The minimized experimental Form II structure (left) has lattice parameters $a=14.12$, $b=4.13$, $c=10.28$, $\beta=102.59$. The restructured form II above 60K (right) has lattice parameters $a=14.92$, $b=4.08$, $c=10.06$, $\beta=114.98$. Both minima have the spacegroup $P2_1/c$.

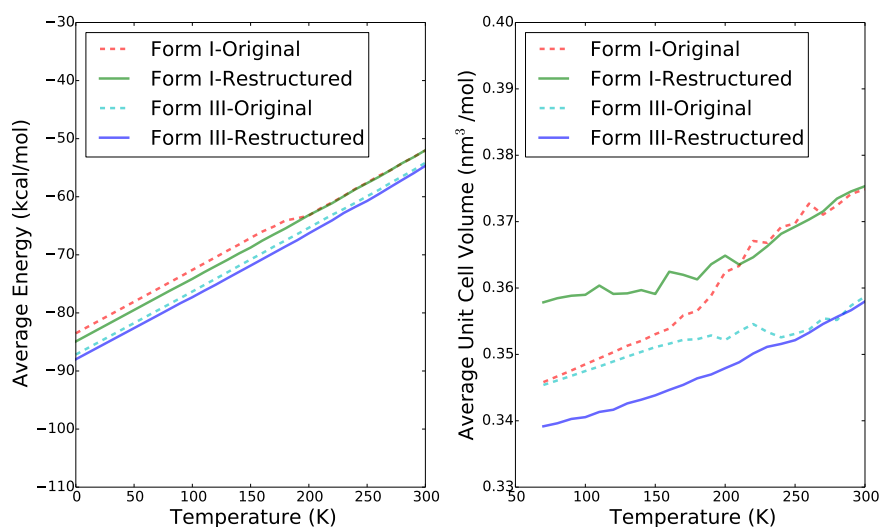


Figure 5.3: Tolbutamide form I and form III restructure to new crystal forms at ambient temperatures. The forms have indistinguishable average energy and unit cell volume above the restructuring temperature.

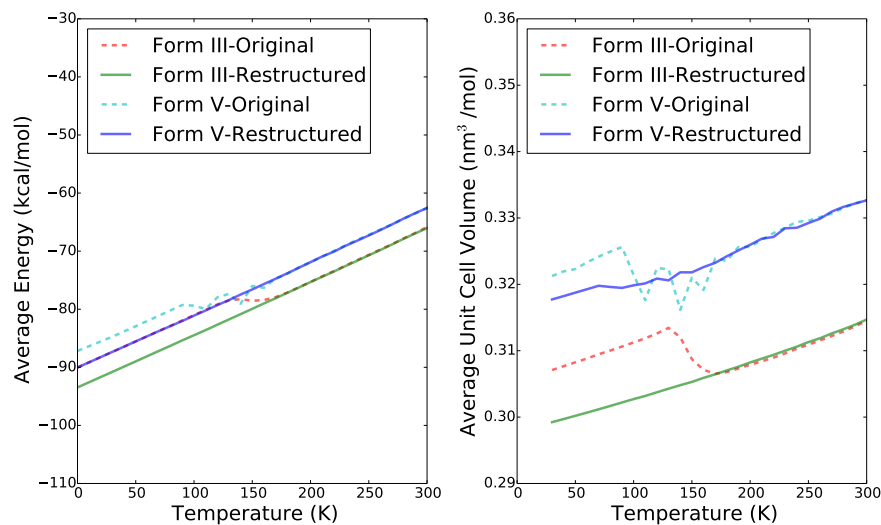


Figure 5.4: Chlorpropamide form III and form V restructure to new crystal forms at ambient temperatures.. The forms have indistinguishable average energy and unit cell volume above the restructuring temperature.

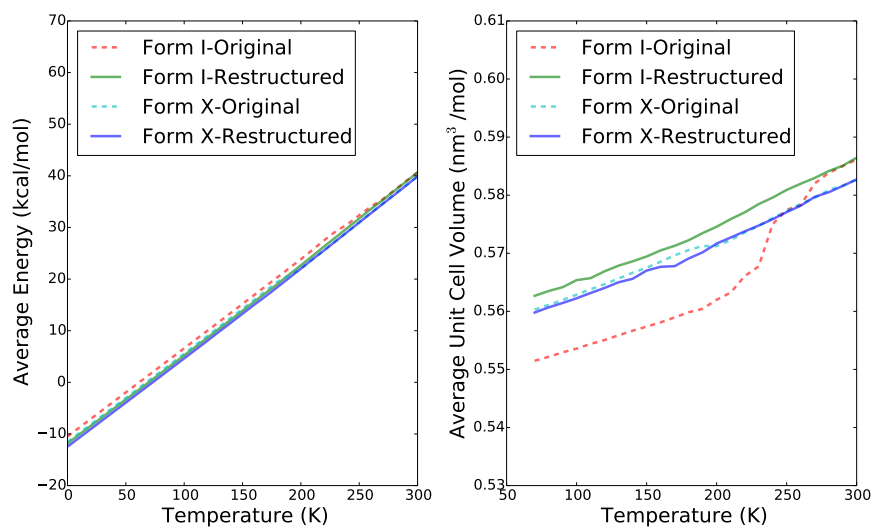


Figure 5.5: Aripiprazole form I and form X restructure to new crystal forms at ambient temperatures. The forms have indistinguishable average energy and unit cell volume above the restructuring temperature.

it requires significantly larger supercells. However, the restructuring events in this study occur on the time scales of nanosecond molecular dynamics simulations, suggesting that the barriers would be easily crossed at ambient conditions on experimental time scales. Therefore, although the original minimized structure corresponds with a locally stable point on the lattice energy surface, it does not correspond with a polymorph that is distinct from the restructured lattice minima at ambient conditions. Instead, the two lattice energy minima both lead to the same crystal form at ambient temperature. In other words, the two structures rapidly converge to the same equilibrium ensemble at ambient temperature. Indeed, the average energy, average volume, and relative free energy of pyrazinamide β at 300K are identical when starting from the original and restructured lattice minima (Figure 5.1).

The occurrence of multiple lattice minima having the same ambient temperature form appears particularly prevalent in larger molecules with high flexibility. All three systems in this study with >30 atoms per molecule and > 1 rotatable bond had at least two crystal forms undergo restructuring. By contrast, only 2 of the rigid systems had a polymorph undergo minima conversion. These examples of converting lattice minima illustrate an additional area where molecular dynamics can improve current crystal structure predictions by eliminating “degenerate” lattice minima that are not also unique free energy minima at room temperature.

5.4 Minima Interconversion

Both cyclopentane and succinonitrile had disordered crystal forms that were observed to experience lattice minima interconversion. The energy and density of all the crystal minimized configurations from the 150K trajectory of form I cyclopentane and the 300K trajectory of form II succinonitrile are shown in Figure 5.6. This plot reveals many lattice energy minima that are of similar energy. Additionally, the energy basins corresponding to each of these minima were explored from a single molecular dynamics trajectory. These clearly demonstrate the phenomena of lattice minima interconversion described in this chapter.

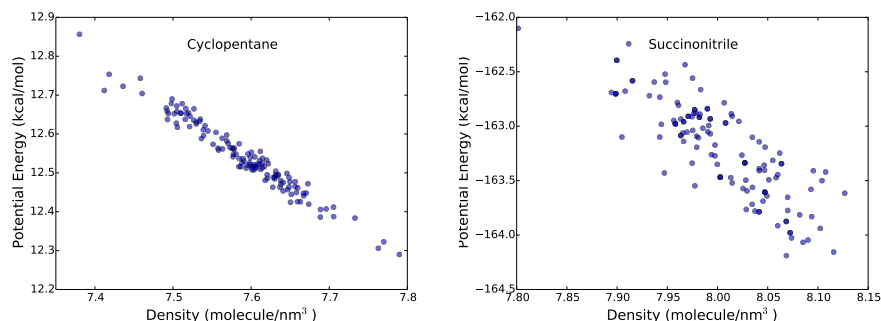


Figure 5.6: The trajectories of form I cyclopentane (left) and form II succinonitrile (right) sample the configuration space around numerous lattice energy minima rather than around single low energy configuration. These lattice energy minima therefore interconvert at 150K and 300K respectively.

5.5 Conclusions and Future Outlook

Numerous crystal polymorphs were simulated at ambient temperature with molecular dynamics to search for conversion of lattice minima. A majority of the crystal structures simulated did not undergo any lattice minima conversion at ambient temperature. However, irreversible conversion from an energetically unstable lattice minima to a more stable configuration was observed in 7 of the 36 crystal polymorphs examined. 6 of the 7 polymorphs that underwent irreversible conversion were from flexible molecules with at least 1 rotatable degree of freedom. This suggests that minima conversion is more prevalent in larger and more flexible molecules rather than small rigid molecules.

Additionally, 2 crystal forms were found to be disordered and undergo lattice minima interconversion. In these systems, the ambient temperature trajectory explores configurations that correspond with many lattice energy minima rather than a single low energy structure. Thus, freezing the trajectory down to 0K leads to numerous lattice energy minima.

The results presented here demonstrate that molecular dynamics can effectively identify

when two or more initial crystal configurations produce indistinguishable polymorph ensembles at ambient temperatures. These methods would be invaluable as a final step for lattice-energy based crystal structure prediction studies. CSP studies produce far more minima than there are experimental structures. Molecular dynamics has the potential to identify which of these lattice minima convert into other structures at ambient conditions. Any minima that converts at room temperature can be effectively ignored from any future refinement step or targeted crystallization experiment.

Figure 5.7 depicts the lattice energy landscape of resorcinol. This landscape includes 140 distinct structures. It is very likely that many of these lattice minima will undergo conversion or interconversion at ambient temperatures. Heating up configurations produced from a CSP such as this to eliminate dead-end structures represents a promising future direction for molecular dynamics in crystal structure prediction. Specifically, Chapters 3, 4, and 5 each proposed a method for eliminating dead-end structures through temperature-mediated transformations, the effects of more accurate potentials, and conversion of lattice minima, respectively. The exact fraction of dead-end lattice minima that can be explained away by each of these three proposals would be invaluable information for the crystal structure prediction field.

Another promising future avenue for molecular dynamics in crystal structure prediction is estimating the added value of using quantum mechanical, rather than classical, descriptions of crystal polymorphs. In Chapter 4, three crystal polymorphs of benzene were modeled with a cheaper point-charge potential as well as a more expensive polarizable potential. In two of the three benzene polymorphs examined, the polarizable hamiltonian had only a minor effect on crystal shape and relative entropy. However, the polarizable Hamiltonian produced a restructured benzene I polymorph that had significantly different relative entropy and enthalpy estimates compared with those from the point-charge potential. Although the polarizable Hamiltonian is more expensive and likely to be more physically accurate, both of these Hamiltonians ignore any quantum mechanical effects on the system. It is plausible that quantum-based Hamiltonians will produce crystal polymorphs with significantly different enthalpies and entropies which could lead to more accurate comparisons to experimental

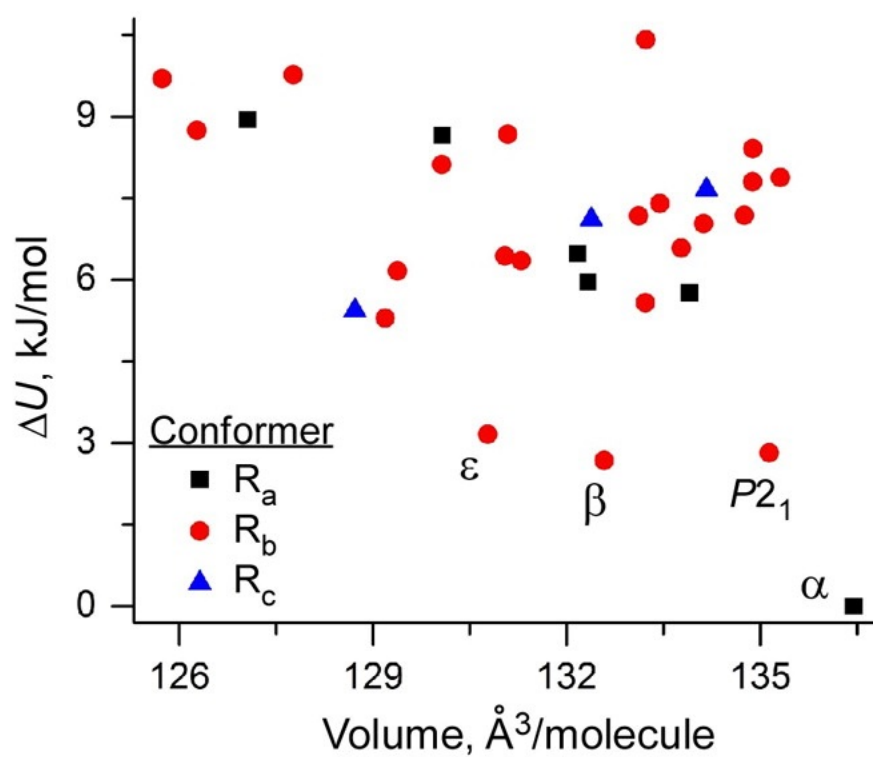


Figure 5.7: The lattice energy landscape of resorcinol, reproduced from ref 198)

measurements. Thus, crystal free energy models in quantum mechanical Hamiltonians represents a ripe area for future research.

The polarizable free energies in Chapter 4 were computed with little or no direct sampling in the expensive potential using a novel method of Hamiltonian reweighting. This indirect Hamiltonian reweighting procedure significantly reduces the cost of estimating free energies in the polarizable Hamiltonian by avoiding the need for direct simulation in the expensive potential. Hamiltonians with quantum mechanical effects are often multiple orders of magnitude more expensive than the polarizable Hamiltonian used in Chapter 4. Therefore, any investigation into the effects of quantum Hamiltonians on crystal stabilities will likely require some type of indirect sampling and reweighting scheme.

The next chapter briefly departs from crystal polymorphism to compare different reweighting methods from a classical to a hybrid QM/MM potential in estimating solution phase solvation free energies. The results show promise for being able to accurately and efficiently compute free energy differences in more expensive QM/MM potentials in solution phase systems. Such reweighting protocols may also be effective in free energy differences between solid polymorphs. Additionally, the comparison definitively demonstrates that MBAR is more efficient than other commonly utilized multistate reweighting techniques from reweighting from MM to QM/MM potentials.

Chapter 6

Evaluating Methods to Reweight from MM to QM/MM Potentials

6.1 Introduction

Molecular dynamics and Monte Carlo simulations have become ubiquitous tools in the field of computational chemistry and allow researchers to calculate numerous thermodynamic properties of interest such as solvation free energies [249,257,275–277], ligand binding [278–281], protein folding [282–285], and now also temperature-dependent crystal polymorph free energies [119]. These computational techniques continue to experience increasing success in reproducing and even predicting experimentally measured properties, as demonstrated by recent blind prediction studies. [109–113,249,257,286–296] The two primary challenges in matching computed estimates with experimental values are the efficiency of the simulation in sampling relevant system configurations and the accuracy of the Hamiltonian used to generate the samples [297].

The ability to estimate free energy differences using expensive energy models, up to and including quantum mechanics, is often desirable because these *ab initio* potentials automatically incorporate degrees of electronic freedom, such as polarizability, which may be critical to describe the system of interest. Quantum mechanical methods also avoid the need to parameterize classical partial charges and intramolecular parameters on a per-atom basis. The

importance of capturing the proper polarization, charges, and bonded parameters were re-emphasized in the recent SAMPL blind prediction studies. [286–295].

In the context of solid materials, quantum effects may play a crucial role in crystal stabilities. Reilly and Tkatchenko recently showed that the stability of the observed form I of acetylsalicylic acid (Aspirin) results from a coupling of the electronic fluctuations and lattice vibrations in higher level energy models rather than a mechanical instability of form II as previously proposed [234]. DFT studies of PVF have also shown that the dipole-dipole interactions are essential in reproducing the correct vibrational spectra, which affects both the ferro- and piezo-electric properties and also the free energy [236].

Polarity can also influence the kinetics and synthesis process toward metastable polymorphs. Recent studies have shown that the use of polar solvent [298] and adjusting the polarity of the crystallization catalyst [299] can change the observed structure. Polarity can also affect the symmetry group of the crystal structure [119,300] and the growth of thin films [301]. Polarization is also suggested to explain the high pressure-high temperature polymorphs of multiple common compounds [78–80,302]. Capturing these effects with a computer model will require a Hamiltonian that incorporates at least polarization and possibly quantum effects as well.

Even in cases where an accurate result was obtained with a classical potential, the cheap potential is often first fitted to previously collected experimental or high-accuracy QM data [118,142,198,235]. In a ‘true’ crystal structure prediction study of a totally new and unstudied material, this data will not exist and a series of QM calculations will need to be run anyway. Thus, the ability to compute crystal free energies directly in expensive *ab initio* potentials is a key frontier for unlocking true blind crystal prediction capabilities. However, simulating directly in complex Hamiltonians with quantum mechanical effects is often prohibitively expensive, which further aggravates the challenge to obtain a sufficient number of samples for converged thermodynamic quantities.

In Chapter 4, we saw that thermodynamic quantities like free energy can be computed in a more expensive potential indirectly using simulations from a cheaper potential. The equations

for computing these indirect quantities are [191]:

$$\Delta A_{12} = -\beta^{-1} \ln \left\langle e^{-\beta \Delta U_{12}(x)} \right\rangle_1 \quad (6.1)$$

$$\langle O \rangle_2 = \left\langle O(x) e^{\beta \Delta A_{12} - \beta \Delta U_{12}(x)} \right\rangle_1 \quad (6.2)$$

Where $\Delta U_{12}(x) = U_2(x) - U_1(x)$ is the energy difference between energy function 1 and energy function 2, $\Delta A_{12} = A_2 - A_1$ is the free energy difference between the two potentials, and $O(x)$ is any observable of the system. These equations require only energy evaluations in the target Hamiltonian which avoids completely the computational expense of direct sampling in the expensive potential. When the target Hamiltonian is an expensive quantum-based potential, the rate limiting step in computing a free energy estimate using eq 6.2 is the energy evaluations of the cheap sampled MM configurations in the target potential. Although a large quantity of samples can easily be drawn from the cheap MM state, each configuration must be re-evaluated in the expensive potential before being used in equation 6.2. It is therefore critical to efficiently utilize the information from the MM sampling, as only a few configurations may have significant probability in the expensive potential. As can be seen in eq 6.2, if $U_{12}(x)$ varies significantly from one configuration to another, the average will be dominated by only a small set of configurations. All ensemble averages in the expensive potential, as well as the free energy estimate, will have large uncertainties if only a small number of configurations with high probability in the expensive potential are sampled.

A range of different alternate approaches to eq 6.2 have been used previously to improve the efficiency of reweighting from a cheap to expensive potential. As we review below, each of these approaches has a different bias and statistical efficiency. There is currently no clear consensus on which approaches are most efficient and have the lowest systematic error (bias), as evidenced by the numerous studies in the past decade which use a variety of different Boltzmann reweighting techniques, and we attempt to resolve this situation here.

The simplest method to obtain free energy differences via Hamiltonian reweighting involves reweighting from one sampled to one unsampled state, generally referred to as exponential reweighting or the Zwanzig equation [191,238,240,242–244,247–250,252], which is simply eq 6.1. This method produces asymptotically unbiased free energy differences in the limit of large sampling. However, this method suffers from the known issue of being numerically noisy and highly sensitive to the degree of phase space overlap and amount of sampling [259].

In recent years, a series of reweighting techniques have been developed which use configurations from multiple sampled states simultaneously to more efficiently calculate the free energy difference between a cheap and expensive potential [188,255,303]. Multistate reweighting techniques generally outperform Zwanzig approaches [257,304] because information from two or more sampled states is used to probe the configuration space of the unsampled state. Two of the most widely used multistate reweighting methods are the Non-Boltzmann Bennett (NBB) method developed by König et al. [255–257,291] and the Multistate Bennett Acceptance Ratio (MBAR) method developed by Shirts and Chodera [188,189]. MBAR leads to free energy estimates identical to those from the weighted histogram analysis method in the limit of zero width bins [305], and also includes accurate asymptotic error estimates [304]. With both MBAR and NBB, an estimate of the free energy difference between two states in the expensive potential is calculated by evaluating the energy in the expensive potential of all configurations from some set of cheap sampled states. Although NBB and MBAR use identical inputs, the differences in the formulation of each statistical estimator potentially lead to differences in the bias and uncertainty in the final free energy estimate. To obtain the optimal estimate of the free energy difference to an unsampled state given a set of configurations in a sampled state one should use the reweighting method which produces the lowest biases and uncertainties, and the decision between NBB and MBAR requires a direct comparison of the biases and uncertainties produced by both estimators.

In this chapter, we directly compare the NBB and MBAR methods, as well as several variants, on the same dataset to determine the most efficient ways to obtain an estimate of the free energy in an expensive Hamiltonian given sampling from a cheaper classical potential. A

recent study suggests that the traditional exponential averaging method outperforms the multistate NBB method when reweighting between an MM and QM/MM potential [306]. Here we extend this comparison by including the Multistate Bennett Acceptance Ratio method as well as logical variants of both MBAR and NBB. All methods are compared for an equal computational expense (an equal number of QM/MM energy evaluations, since all methods use the same MM simulations), and the statistical uncertainty is used as the benchmark for method efficiency. We first examine, in depth, a simple system of methanol and ethane in water and calculate the difference in solvation free energy between these two compounds using both NBB and MBAR. The same comparative analysis of the uncertainties between MBAR and NBB is then applied to the solvation free energy of the molecules in the SAMPL4 small molecule solvation dataset. In all cases, the molecules are simulated in a classical fixed-charge potential and the system is reweighted to a QM/MM Hamiltonian. A previous study using NBB showed that the QM/MM Hamiltonian provides a small but statistically significant improvement over standard MM force fields at matching experimental data [249, 257, 296]. We evaluate whether the magnitude of this benefit changes when different methods are used to reweight the same MM configurations to the same QM/MM Hamiltonian.

With all multistate reweighting methods, the energy re-evaluations in the expensive potential are performed with configurations drawn from multiple different sampled states. A ‘state’ in the context of this chapter consists of a temperature and volume as well as a specific Hamiltonian. This Hamiltonian at each state can correspond to a physical system, referred to as the end state, or can be one of many unphysical modified Hamiltonians, referred to as intermediate states. One possible allocation of these QM/MM calculations is to place an equal number of energy re-evaluations at each sampled end state and intermediate state. However, one could in theory place a larger fraction of the energy re-evaluations in a particular sampled state in order to lower the uncertainty in the final free energy estimate. Thus, we also test a variety of different allocations and assess whether we minimize the variance in the free energy difference estimate by evenly dividing the energy re-evaluations among all sampled states or placing more of the re-evaluations in a single state.

It is worth noting that a number of previous studies by Essex and coworkers have made modifications to the original Zwanzig relationship in order to overcome the large statistical uncertainty in eq 6.1 by using only the intermolecular interaction energies for ΔU_{12} [247,249,250]. The resulting free energy produced by eq 6.1 with this substitution reflects the difference in free energy between the pure MM state and an unsampled state with MM intramolecular interactions and QM/MM intermolecular interactions. The additional step to convert the intramolecular interactions from MM to QM/MM at each endstate is assumed to contribute negligibly to the overall free energy difference. The bias introduced by this approximation is in many cases small, [247,249,250] but there is no rigorous bounds or systematic prediction of the size of the approximation. Understanding the limits of this uncontrolled approximation is beyond the goals of the present study and we instead focus on reweighting variants that use rigorously correct approaches with asymptotically unbiased free energy differences.

6.2 Reweighting Methods

6.2.1 Non-Boltzmann Bennett

Bennett showed [190] that the minimum variance estimate of the free energy between two states, given configurations sampled from both states, can be calculated with

$$\Delta A_{0,1} = \beta^{-1} \ln \left(\frac{\langle f(\beta(U_0 - U_1 + C)) \rangle_1}{\langle f(\beta(U_1 - U_0 - C)) \rangle_0} \right) + C \quad (6.3)$$

where $\Delta A_{0,1}$ is the free energy difference between states 1 and 0, U_i is the energy of a configuration in state i , f denotes the Fermi function

$$f(x) = (1 + \exp\{\beta x\})^{-1} \quad (6.4)$$

and

$$C = \Delta A_{0,1} + \beta^{-1} \ln \left(\frac{N_1}{N_0} \right) \quad (6.5)$$

where N_i denotes the number of configurations sampled from state i . The estimator shown in eq 6.3 is valid and asymptotically unbiased for any C , and thus will converge toward the correct free energy with any value, but the choice of C shown in eq 6.5 gives the lowest possible statistical uncertainty. If no samples are drawn from either of the states 0 or 1 then equation 6.3 cannot be used to estimate the free energy difference. We often encounter this situation when one or both of the states are in a quantum-based Hamiltonian where direct sampling is computationally infeasible.

In the Non-Boltzmann Bennett method, the energy difference between the two sampled states 0 and 1 in eq 6.3 is replaced with the energy difference between two *unsampled* states 0 and 1, and the configurations are drawn from two ‘biased’ states 0, b and 1, b . The bias in the potential energy for each sampled configuration in the biased state is

$$\Delta V^{bias} = U^{biased} - U \quad (6.6)$$

where U^{biased} and U represent the potential energy in the sampled and unsampled Hamiltonian, respectively. The ensemble averages in the unsampled states 0 and 1 are recovered from the configurations sampled in the biased states using the relationship originally developed by Torrie and Valleau [307, 308] to calculate the unbiased ensemble average $\langle X \rangle$ of a property X from samples collected in a biased state:

$$\langle X \rangle = \frac{\langle X \exp(\beta V^{bias}) \rangle_b}{\exp(\beta V_b^{bias})} \quad (6.7)$$

Calculating the ensemble averages in equation 6.3 with configurations in biased states reweighted with equation 6.7 leads to the NBB equation for calculating the free energy difference between two unsampled states:

$$\Delta A_{0,1} = \beta^{-1} \ln \left(\frac{\langle f(U_0 - U_1 + C) \exp\{(\beta V_1^{bias})\} \rangle_{1,b} \langle \exp\{(\beta V_0^{bias})\} \rangle_{0,b}}{\langle f(U_1 - U_0 - C) \exp\{(\beta V_0^{bias})\} \rangle_{0,b} \langle \exp\{(\beta V_1^{bias})\} \rangle_{1,b}} \right) + C \quad (6.8)$$

Note that the minimum variance constant C is defined in terms of the number of samples from each of states 0 and 1. However, there are no samples from state 0 and 1. Since we cannot use Bennett's original formula for the minimum variance choice of C , the choice used in previous NBB papers was to implicitly set $N_0 = N_1$. Since both N_0 and N_1 are in fact equal to each other (besides being equal to zero), this seems a natural choice, but the original derivation of the minimum variance C in eq 6.5 is no longer applicable. However, the formula will still converge to the correct answer in the limit of sufficient samples for all finite choices of C .

In many instances, the initial and final unsampled states of interest are connected through a series of intermediate states $\lambda_0, \dots, \lambda_K$ and the full free energy difference is computed through the sum of the free energy differences between each step along the entire thermodynamic path.

$$\Delta A_{tot} = \sum_{k=0}^{K-1} \Delta A_{k,k+1} \quad (6.9)$$

The computational overhead of computing free energies with eq 6.8 generally depends most sensitively on the number of energy evaluations in the expensive unsampled potential. Nominally, the QM energy of every configuration in each sampled intermediate state should be computed for NBB, and the total free energy difference is calculated as the sum across all unsampled intermediate states with eq 6.9. The NBB method that uses all intermediate unsampled states is referred to as 'NBB-direct' and a diagram for this method is presented in Figure 6.1. The term 'direct' in this context does not reflect that sampling is done directly in the expensive potential. Rather, this follows the naming convention used in a recent study by Min and coworkers in which 'direct' implies that the alchemical path connecting the end states includes intermediate states in the expensive potential [309]. If a total of N samples are drawn from each of K intermediate states, the total number of QM energy evaluations with

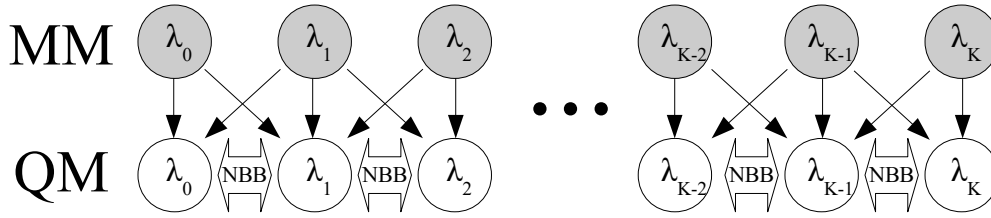


Figure 6.1: Illustration of the NBB-direct method in which each intermediate MM state is connected to an intermediate unsampled QM state. Grey nodes indicate sampled states and white nodes indicate unsampled states. Thin black arrows indicate where configurations from the sampled state are re-evaluated in the target state. Thick arrows indicate where free energy differences are evaluated. The NBB-direct method requires QM energies for all configurations from all sampled MM states.

NBB-direct is KN .

In the systems studied in this work, the intermediate states along the thermodynamic path represent a mixed Hamiltonian of the two physical end states. All intermediate states are therefore unphysical and the QM/MM energies of configurations in these states are not well-defined. We instead assign energies to these intermediate states by first evaluating the QM/MM energy at both physical end states and then interpolate the energy in the intermediate state λ_k with:

$$U_k^{QM/MM}(x) = \lambda_k^y U_K^{QM/MM}(x) + (1 - \lambda_k)^y U_0^{QM/MM}(x) \quad (6.10)$$

where $U_k^{QM/MM}(x)$ is the QM/MM energy of configuration x in the intermediate state λ_k , interpolating between the physical end states λ_0 and λ_K , and y denotes the order of the mixing rule between the two states (generally assumed to be 1). Because a QM/MM energy evaluation is computed at both end states for each configuration, there are a total of $2KN$ QM/MM

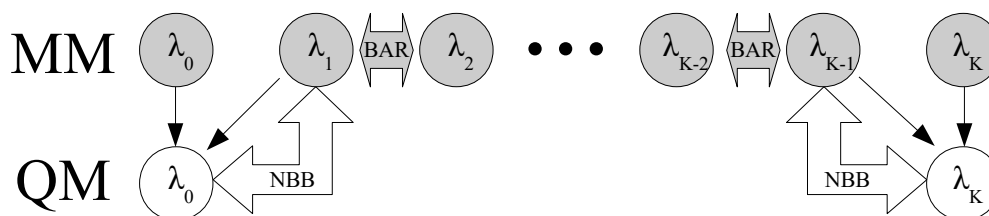


Figure 6.2: Illustration of the NBB-indirect method in which each intermediate MM state is connected through BAR except the final two MM states which are reweighted to the unsampled QM state. Grey nodes indicate sampled states and white nodes indicate unsampled states. Thin black arrows indicate that configurations from the sampled state are re-evaluated in the target state. Thick arrows indicate where free energy differences are evaluated. The NBB-indirect method requires QM energies for all configurations from λ_0 , λ_1 , λ_{K-1} , λ_K .

evaluations necessary for NBB-direct.

The number of QM evaluations can be reduced below $2KN$ by setting the bias, $V_i^{bias} = 0$ for all $i \in \{1..K-1\}$. The intermediate unsampled QM states collapse to the sampled MM states and no QM energies need to be calculated at these states. The thermodynamic cycle in this reweighting scheme, referred to as ‘NBB-indirect’, is illustrated in Figure 6.2 and was adopted in all previously published studies employing NBB to save computational expense [255–257, 291]. NBB-indirect requires the QM energies of all N configurations sampled from states $\lambda_0, \lambda_1, \lambda_n - 1$, and λ_n for a total of $4N$ QM energy evaluations. No interpolation from equation 6.10 is needed for NBB-indirect because only the QM energies at one of the physical end states is needed for every reweighted configuration. Although NBB-indirect is the most commonly used variant of the Non-Boltzmann Bennett method, we have included the NBB-direct method for completeness because it utilizes additional QM/MM calculations which gives the potential for NBB-direct to outperform NBB-indirect.

6.2.2 Multistate Bennett Acceptance Ratio

The Multistate Bennett Acceptance Ratio (MBAR) method was developed as an extension of the widely used Weighted Histogram Analysis Method (WHAM) in the limit of zero-width bins, and a generalization of the minimum variance Bennett Acceptance Ratio (BAR) method for more than 2 states. MBAR calculates the free energy differences between any number of sampled and unsampled states simultaneously and has been proven to have the lowest statistical error among all reweighting estimators [188,189]. The MBAR estimate of the relative free energy of state i is

$$A_i = -\beta^{-1} \ln \sum_{j=1}^K \sum_{n=1}^{N_j} \frac{e^{-\beta U_i(x_{jn})}}{\sum_{k=1}^K N_k e^{\beta A_k - \beta U_k(x_{jn})}} \quad (6.11)$$

where x_{jn} is the n^{th} configuration sampled in state j , and U_i is the energy of the configuration in state i . A notable difference between the NBB estimator (eq 6.8) and the MBAR estimator (eq 6.11) is that MBAR sums over all sampled configurations in the ensemble averages at each

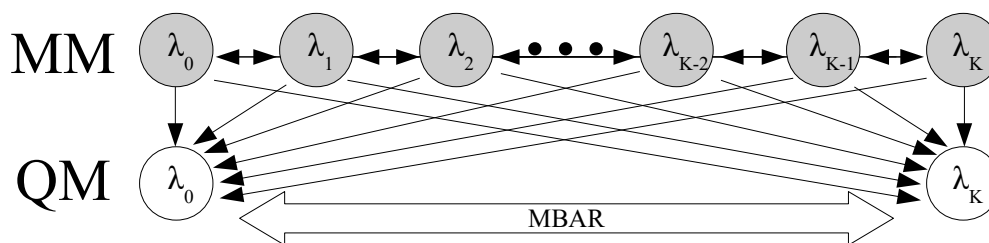


Figure 6.3: Illustration of the full MBAR method in which all sampled MM states are connected together and the free energies of all sampled and unsampled states are calculated simultaneously. Grey nodes indicate sampled states and white nodes indicate unsampled states. Thin black arrows indicate that configurations from the sampled state are re-evaluated in the target state. Thick arrows indicate where free energy differences are evaluated. The full MBAR method requires QM energies for all configurations from all sampled MM states.

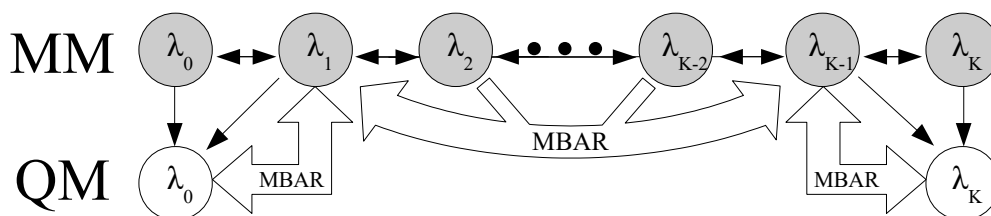


Figure 6.4: Illustration of the three-stage MBAR method in which the states are divided into three groups that are each solved simultaneously with MBAR. Grey nodes indicate sampled states and white nodes indicate unsampled states. Thin black arrows indicate that configurations from the sampled state are re-evaluated in the target state. Thick arrows indicate where free energy differences are evaluated. The three-stage MBAR method requires QM energies for all configurations from the MM states included in the groups with the unsampled QM states. The three-stage MBAR method with two reweighted states at each end is shown here, however the method can generally be applied to any number of reweighted states.

state, while NBB uses only configurations sampled from state i when computing ensemble averages in state i (see Appendix D section D.2). The MBAR estimator therefore uses extra information about state i that is not utilized in the NBB estimator. To calculate the difference in free energy between two unsampled QM end states using samples from K intermediate states, MBAR requires the energy of every sampled configuration in each of the two unsampled states for a total of $2KN$ QM calculations, which is the same expense as NBB-direct. Unlike NBB-direct, there are no unphysical QM/MM intermediate states necessary for the full MBAR method. The sampled configurations need to be re-evaluated only in the two physical QM/MM end states. Potential energies of configurations at each MM state must also be calculated in all other MM states, but this expense is usually negligible compared to the expense of recomputing energies in the QM potential. The reweighting process for MBAR is illustrated in Figure 6.3.

The number of QM calculations necessary to estimate the free energy difference between the two unsampled states along a thermodynamic path with MBAR can be reduced by partitioning the thermodynamic path into three stages and computing the free energy difference across each section individually with MBAR. This reweighting scheme, which we call three-stage MBAR, is illustrated in Figure 6.4. The number of intermediate states used in the MBAR calculation on either end of the thermodynamic path in Figure 6.4 can vary anywhere between 1 and $K/2$ states. The total number of QM energy evaluations with three-stage MBAR can therefore be set anywhere between $2N$ and KN depending on the number of reweighted MM intermediate states. In the limit that 1 state is reweighted at each end, the three-stage MBAR method collapses to a Zwanzig reweighting between the physical QM and MM end states followed by MBAR across the entire range of MM states. For simplicity, we henceforth refer to the three-stage MBAR method with N reweighted states at each end of the thermodynamic cycle as three-stage MBAR(N). We note that for more complex calculations, what we describe as the three-stage MBAR method may contain more than three actual stages. For example, in the SAMPL4 dataset calculations described in Section 6.3.2 there are a total of four states, with the gas phase and solution phases each having one stage connecting the MM state to the

QM/MM state and one stage connecting the MM state to an analytical reference.

6.2.3 Effective Number of Samples

A final concept we will use in this chapter is the number of effective samples. The number of effective samples provides a rough quantitative metric of the information gained in an unsampled state using samples from another state. If we have a weighted average $\langle A \rangle = \sum_{n=1}^N A_n W_n$, then Kish's formula for calculating the number of effective samples is [310]

$$\Delta N_{eff} = \frac{(\sum_{i=1}^n W_i)^2}{\sum_{i=1}^n W_i^2} = \frac{1}{\sum_{i=1}^n W_i^2} \quad (6.12)$$

where the last step assumes normalized weights $\sum_n W_n = 1$. As can be seen from this formula, if one weight is significantly larger than all others, the number of effective samples approaches 1; if all weights are similar, the number of effective samples approaches N .

Ensemble averages calculated via MBAR can be expressed in terms of a weighted sum over the samples where the weights in state i are calculated as

$$W_{in} = \frac{e^{\beta A_i - \beta U_i(x_n)}}{\sum_k N_k e^{\beta A_k - \beta U_k(x_n)}} \quad (6.13)$$

where A is the Helmholtz free energy of each state (or Gibbs free energy for NPT simulations), U is the energy of each state, N_k is the number of samples from state k , and states $k = 1 \dots K$ are sampled states. If state i has high overlap with the other K states, then the number of effective samples in state i will approach $\sum_k N_k$. If state i , has high overlap with only one state k , then the number of effective samples in state i will approach N_k . Finally, if state i has very little overlap with any state, then the number of effective samples in state i will approach 1.

6.3 Simulation Details

6.3.1 Ethane-Methanol

All simulations were conducted with CHARMM [163,311], using the CHARMM22 [312] force field. The QM and QM/MM calculations were performed with Q-Chem [313] based on the CHARMM/Q-Chem interface [314].

The implementation of the ethane-methanol free energy calculations follows the one presented in Reference 256. Solvation free energy differences between ethane and methanol were calculated using the standard thermodynamic cycle. Gas phase simulations were conducted with Langevin dynamics, using a friction coefficient of 5 ps^{-1} on all atoms and random forces according to a target temperature of 300 K. In solution, we used 862 water molecules and an octahedral box that was cut from a cube with a side length of 32.166 Å. The temperature was maintained at about 300 K by a Nosé-Hoover thermostat [172]. Lennard-Jones interactions were switched off between 10 and 12 Å, while electrostatic interactions were computed with the Particle Mesh Ewald method [168]. The time step was 1 fs and SHAKE was applied to all hydrogen atoms. In the gas phase, the cut-off radius was set to 998 Å.

Free energy differences were calculated based on simulations of 5 ns in gas phase and 1 ns in solution. Trajectories were written every 100 steps in gas phase and 20 steps in solution for a total of 10,000 samples from each phase. For the free energy calculations, eleven λ points were employed in both the gas phase and solution phase (0.0, 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 1.0). Quantum potential energies were calculated for all trajectories. Each simulation was repeated four times, starting with different initial random velocities.

The dual topology hybrid of ethane and methanol was implemented using the MSCALE module of CHARMM [315]. Energy evaluation was divided into three tasks:

- Calculating the energy contributions of all bond, bond angle and Urey-Bradley terms of the full hybrid molecule, U_{bonded}^{MM} (this was done in the main MSCALE process to maintain the connectivity of the hybrid molecule).

- dihedral angle and non-bonded contributions corresponding to state 0, $U_{only\ 0}^{MM}$ (i.e., all atoms that are not part of ethane or the common environment were deleted)
- dihedral angle and non-bonded contributions corresponding to state 1, $U_{only\ 1}^{MM}$ (i.e., all atoms that are not part of methanol or the common environment were deleted). The λ states were generated by mixing those three energy contributions according to $U_{\lambda}^{MM} = U_{bonded}^{MM} + (1 - \lambda)U_{only\ 0}^{MM} + \lambda U_{only\ 1}^{MM}$

The quantum-mechanical potential energies (U^{QM}) were calculated with CHARMM and Q-Chem based on input files generated with the CHARMM-Q-Chem interface. B3LYP/6-31G* was used for the solute in both gas phase and the explicit solvent QM/MM calculations, and in solution, the solvent was treated classically. The standard QM/MM electrostatic embedding scheme was used such that the solute-solvent interactions were calculated with QM treating the water atoms as point charges.

To evaluate the potential energy of each frame of the gas phase trajectories, each calculation was divided into two tasks:

- removing all atoms not corresponding to the initial state 0 (ethane) and calculating the potential energy, $U_{only\ 0}^{QM}$ and
- removing all atoms not corresponding to the final state 1 (methanol) and calculating the potential energy, $U_{only\ 1}^{QM}$. To calculate the potential energy of λ states, the two terms were mixed according to eq 6.10.

The corresponding explicit solvent QM/MM potential energy calculations ($U^{QM/MM}$) were divided into two steps:

- Calculating the QM/MM potential energy using a single box of water molecules that were centered around the solute for each frame of the trajectory ($U_{box}^{QM/MM}$). The quantum region consisted of the solute and the MM region consisted of the water molecules. Only the intramolecular solute-solute interactions and the solute-solvent interactions were considered. No cutoff was used. To generate the end states of the alchemical

transformation, all atoms not corresponding to the end state were removed from the calculation, leading to $U_{box,0}^{QM}$ for ethane and $U_{box,1}^{QM}$ for methanol.

- The solvent-solvent interactions were added by deleting the solute from the trajectory and calculating the potential energy using periodic boundary conditions (U_{pbc}^{MM}). Thus, this approach neglects the long range solute-solvent and solute-solute interactions with the image atoms of the periodic system. However, the box size is larger than twice the standard CHARMM cutoff of 12 Å and therefore the non-bonded contribution from the periodic image with this setting is exactly zero.

To generate potential energies for each λ intermediate state, $U^{QM/MM}$ was evaluated once for the initial state 0 ($U_{box,0}^{QM/MM}$) and once for the final state 1 ($U_{box,1}^{QM/MM}$), leading to $U_{\lambda}^{QM/MM} = U_{bonded}^{MM} + (1 - \lambda)U_{box,0}^{QM/MM} + \lambda U_{box,1}^{QM/MM} + U_{pbc}^{MM}$ for simulations in solution.

6.3.2 SAMPL4 subset

The data for the blind subset of the SAMPL4 hydration free energy challenge are based on the potential energy data used in reference 257. The blind subset originally consisted of 24 molecules, but 3 out of the 24 molecules (molecules 7, 8 and 18) had to be withdrawn from the challenge. As described in Ref. 257, the original simulations were conducted with CHARMM using the CHARMM General force field (CGenFF), version 0.9.6 beta. The thermodynamic cycle follows precisely the one described in Ref 257. Briefly, the gas and solvent phase molecules were driven to a non-interacting ideal gas state through a two step alchemical mutation. First, all charges of the solute were set to zero. Second, all Lennard-Jones interactions of the solute were set to zero. A total of 6 and 7 lambda points were used for each respective step of the transformation in the gas phase, and 12 and 13 lambda points were used for the transformation in the solution phase. The original QM and QM/MM calculations were performed with Q-Chem based on the CHARMM/Q-Chem interface. The QM and QM/MM data consists of calculations with B3LYP/6-31G* that were only performed for the first λ step of the MM alchemical uncharging transformation. This corresponds to simulations with normal CGenFF

charges ($\lambda=0.00$) and with CGenFF charges that were scaled by a factor of 0.95 ($\lambda=0.05$) in both gas phase and solution. The gas phase data consisted of 45000 data points for each λ point, corresponding to a simulation time of 4.5 nanoseconds. The solvent phase data consisted of 10000 data points for each λ point, corresponding to a simulation times between 0.5 and 1 nanoseconds. To improve sampling, Hamiltonian replica exchange was used between the λ points. Other details of the alchemical simulations are described in Ref. 257.

6.3.3 Harmonic Oscillators

A toy system of 1-D harmonic oscillators is used in addition to the ethane-methanol system and the SAMPL4 molecule system to probe the uncertainties in free energy estimates using the various reweighting methods described in section 6.2. Harmonic oscillators provide a valuable testing ground for free energy estimators because the free energy difference between two oscillators can be calculated analytically, and samples from harmonic oscillators can be easily drawn. The potential energy and free energy of a 1-D harmonic oscillator are:

$$U(x) = k(x - \bar{x})^2 \quad (6.14)$$

$$A = -\beta^{-1} \int_{-\infty}^{\infty} \exp[-\beta U(x)] dx = \frac{1}{2\beta} \ln \frac{\pi}{\beta k} \quad (6.15)$$

where $U(x)$ is the potential energy of the harmonic oscillator at position x , k is the spring constant, \bar{x} is the equilibrium distance of the oscillator, and A is the absolute free energy. In this work, two harmonic oscillators are sampled (satisfying the NBB-indirect method requirements of two sampled states) and used to estimate the free energy difference between the second oscillator and a variety of different unsampled harmonic oscillators. In all simulations, the arbitrary parameter β is set to unity and 1000 configurations are drawn from each of the sampled harmonic oscillators.

The first oscillator test system studied here is a series of 11 oscillators which all have a mean position of $\bar{x} = 0$ and have evenly spaced harmonic spring constants ranging from $k = 1.0$ to

$k = 6.0$. An illustration of the potential energy surface of the sampled and unsampled oscillators in this system are plotted in Figure 6.14. The increasing spring constants decrease the overlap with the sampled oscillators and mimic the physical effect a degree of freedom that is stiffer in the target state than the sampled state, such as a higher frequency bond vibration. The statistical estimators described in section 6.2 are applied to this oscillator system in two different cases. In the first case, the two oscillators with the smallest spring constants are sampled and reweighted to the stiffer unsampled oscillators. In the second case, the two harmonic oscillators with the largest spring constants are sampled and reweighted to the unsampled oscillators with weaker spring constants.

The second test system examined here is another series of 11 oscillators which all have a harmonic spring constant of $k = 1.0$ and evenly spaced mean offsets ranging from $\bar{x} = 0$ to $\bar{x} = 9$. An illustration of the potential energy surfaces for these sampled and unsampled oscillators are shown in Figure 6.15. The varying offset for the unsampled oscillators mimics the physical effect of a degree of freedom with a different mean value in the target and sampled states, such as a longer average bond length.

6.3.4 Uncertainty Analysis

All uncertainties in the free energy estimates produced by NBB and MBAR were calculated by computing the variance of 200 independent bootstrap iterations. The configurations were decorrelated before bootstrapping using the method of Chodera et al. [213]. Uncertainties in the uncertainties for the SAMPL4 dataset were estimated assuming a normal distribution using $\text{Var}[\sigma^2] = \frac{2\sigma^4}{n-1}$ where σ^2 is the variance in the free energy estimate and n represents the number of bootstrap iterations [304].

6.4 Results and Discussion

6.4.1 Ethane-Methanol

The estimates of the difference in solvation free energy between ethane and methanol are shown in Figure 6.5, calculated with only the MM energies as well as with the QM/MM energies using both NBB and MBAR. The uncertainties with MBAR ($0.022 \text{ kcal}\cdot\text{mol}^{-1}$) are 4 times smaller than the uncertainties with the NBB-direct method ($0.088 \text{ kcal}\cdot\text{mol}^{-1}$) when applied to the same set of data. These uncertainties suggest that for this transformation, the full MBAR method is approximately 16x more efficient at producing the free energy estimate. This assumes we are in the limit that the squared uncertainties are inversely proportional to the number of reweighted configurations, which we generally do arrive at for most reweighting methods in the first few hundred samples [304]. Under this relatively mild assumption, then, the NBB-direct method would require 16x more QM/MM energy evaluations in order to produce the same magnitude of uncertainties as the full implementation of MBAR.

In addition to larger uncertainties, the free energy estimate with the NBB-direct method is also dependent on the path connecting the unsampled states. The path dependence of the NBB-direct method occurs because the method requires an estimate of the free energies of all intermediate states in the QM/MM Hamiltonian. The free energy between the end state QM/MM states is calculated as a sum of the pairwise free energies between each adjacent intermediate QM/MM state shown numerically in equation 6.9 and graphically in Figure 6.1. The full MBAR method can in principle be used to estimate the free energies of these intermediate QM/MM states. However, the relative free energies of unsampled states are not dependent on the energies or free energies of any other unsampled state as can be seen from equation 6.11. Therefore, the estimate of the free energy difference between the physical ethane and methanol end states with full MBAR is not affected by the inclusion or exclusion of intermediate QM/MM states.

As one example of the path dependence in the NBB-direct method, a noticeable systematic error (bias) is introduced in the free energy difference estimate when the mixing rule used to

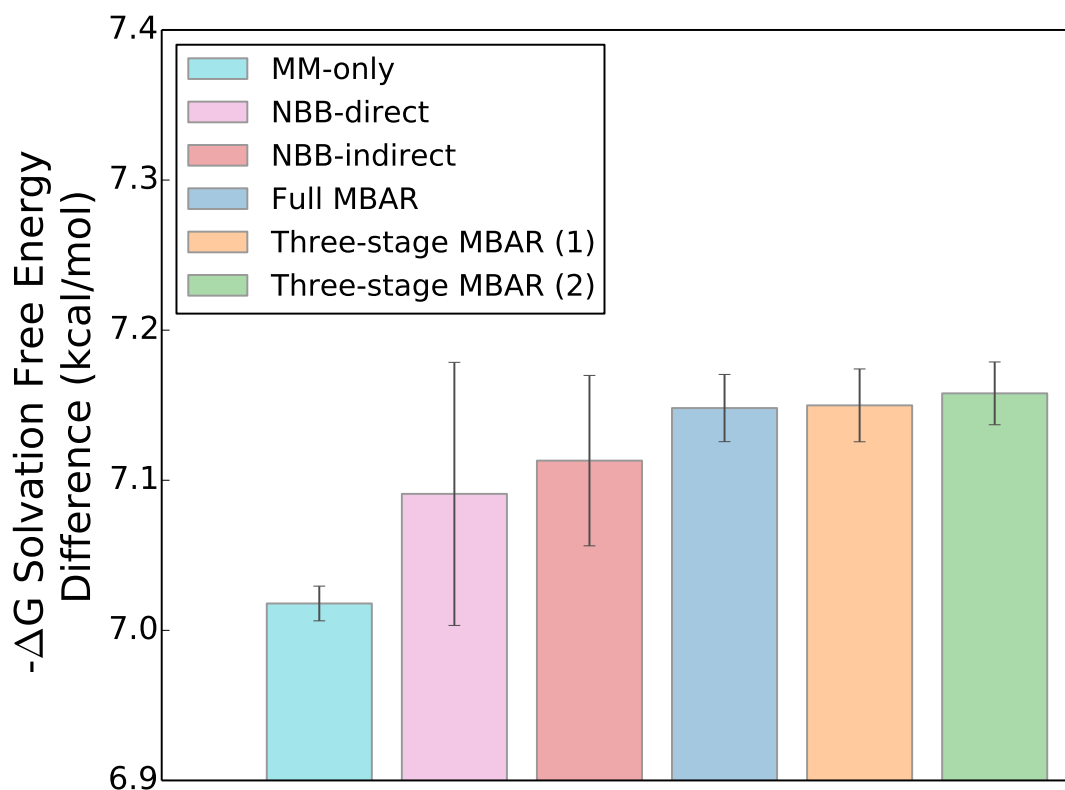


Figure 6.5: The solvation free energy difference estimate between ethane and methanol in the QM/MM Hamiltonian is statistically distinguishable from the estimate in the MM Hamiltonian for all methods except NBB-direct. The bootstrapped uncertainties for the MBAR-derived methods are smaller than the uncertainties using the NBB-derived methods. The parenthetical values for the three-stage MBAR method indicate the number of sampled states reweighted to each unsampled state.

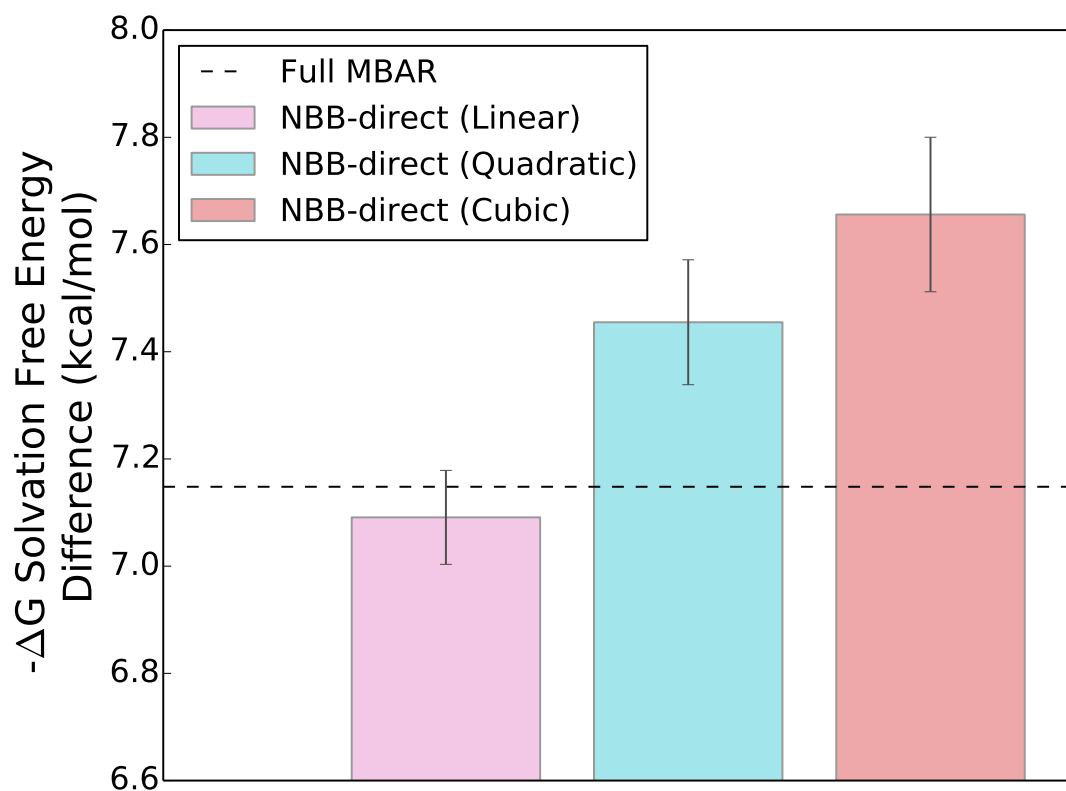


Figure 6.6: The free energy difference estimate with NBB-direct is dependent on the thermodynamic path connecting the unsampled states. Higher order mixing rules for the unsampled QM/MM intermediate states lead to statistically significant biases in the estimate. The dashed line indicates the free energy estimate using the full MBAR method.

calculate the intermediate QM/MM energies (eq 6.10) is increased above first order (shown in Figure 6.6). The free energy difference estimate with the quadratic and cubic mixing rule is statistically different from the estimate with the linear NBB-direct and the full MBAR method based on the bootstrapped uncertainty values. This suggests that using a non-linear mixing rule with NBB-direct method introduces a systematic error into the free energy estimation and not just an increase in the statistical uncertainty. With MBAR, the weights and free energy estimates depend only on the sampled states, and therefore the unsampled state free energies are independent of the number or location of the unsampled intermediate states. The larger uncertainties and the dependence on unsampled pathway for NBB-direct shown in Figure 6.5 and 6.6 demonstrate a clear advantage in using MBAR instead of the NBB-direct method to estimate the solvation free energy difference between methanol and ethane.

The NBB-indirect and three-stage MBAR reweighting methods require significantly less computational expense than MBAR and NBB-direct because the QM/MM energy re-evaluations only need to be conducted at a subset of states near the physical end states of the thermodynamic cycle rather than at every ‘mixed’ intermediate state (see Figure 6.2 and 6.4). The three-stage MBAR method can be run with any number of reweighted states on either end of the thermodynamic path ranging from 1 to $\frac{K}{2}$ where K is the total number of sampled states. However the free energy difference estimate and uncertainty in this system are insensitive to the number of reweighted states (Figure 6.7), and the estimates with more than one reweighted state are statistically identical to the estimates with only the physical end states evaluated in the QM/MM potential.

The solvation free energy difference estimates with the NBB-indirect and three-stage MBAR methods are in good agreement with the estimate using the full QM/MM dataset generated with MBAR (Figure 6.5). The three-stage MBAR(2) method with 2 reweighted states at each end has the most natural comparison to NBB-indirect because both methods use exactly the same number of QM/MM energy evaluations. The uncertainty with three-stage MBAR(2) is $0.021 \text{ kcal}\cdot\text{mol}^{-1}$, while with NBB-indirect the uncertainty is $0.057 \text{ kcal}\cdot\text{mol}^{-1}$. This suggests that three-stage MBAR(2) is roughly 7x more efficient at producing the free energy estimate

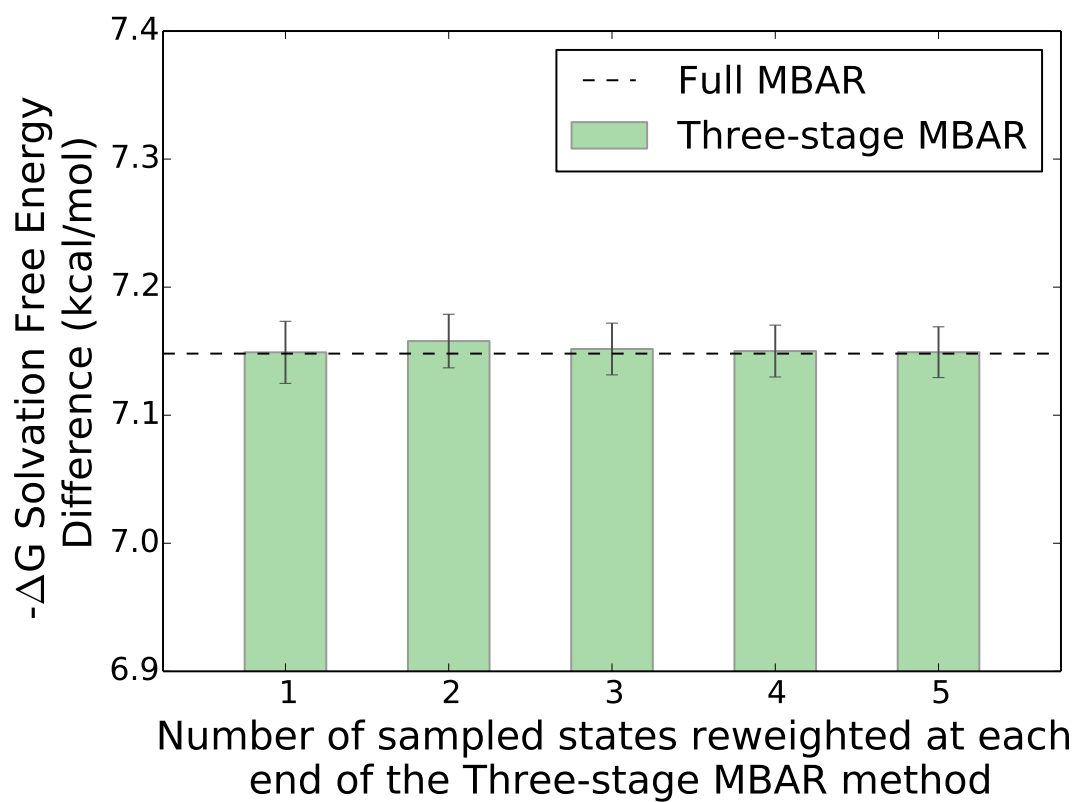


Figure 6.7: The free energy difference estimate with three-stage MBAR in the ethane-methanol system is insensitive to the number of states reweighted at each end of the thermodynamic path. The free energy estimate using any number of reweighted states is within the uncertainty of the value estimated by the full MBAR method which uses all of the available configurations.

with the abridged dataset, again assuming that the squared uncertainties are inversely proportional to the number of reweighted configurations. Overall, the data suggests that for transformations similar to the methanol-ethane system, the QM/MM evaluations should be conducted primarily at the end states of the alchemical path and that the reweighting analysis be conducted with three-stage MBAR (i.e. Zwanzig reweighting of end states to QM/MM plus MBAR for free energies of the intermediate states) rather than NBB-indirect to minimize the uncertainty in the free energy estimate while simultaneously maximizing efficiency.

For the methanol-ethane system, the solvation free energy difference estimate and uncertainty in the QM/MM potential with three-stage MBAR is statistically the same when re-evaluating configurations from all of the sampled states versus only re-evaluating configurations in the physical MM end states. The results from this system suggest that in free energy studies where two states in an expensive Hamiltonian are connected through a set of cheaper intermediate states, the energy re-evaluations in the expensive potential should likely be done primarily with configurations sampled at the end states of the thermodynamic path where the overlap to the expensive potential is the highest. In most cases, the end states are likely to have more overlap because they both represent physical states, but these results do not rule out exceptions, such as where a partially charged MM end state from an over-polarized force field resembles the QM potential more closely than the fully charged MM end state.

We compare in more detail this approach of only using end state samples for the reweighting step from the cheap to expensive potential with the alternative of using samples from multiple intermediate steps. We evaluated the free energy difference estimate and uncertainty for the ethane-methanol system in the QM/MM potential using three-stage MBAR for a fixed total of 15000 energy evaluations evenly drawn from N reweighted states at each end of the thermodynamic path. When $N = 1$, all 15000 samples are drawn from the physical end states and re-evaluated in the QM/MM potential along with equation 6.11 to estimate the solvation free energy difference between the physical QM/MM and MM state. When $N = 5$, 3000 samples are re-evaluated from each of the 5 MM intermediate states closest to the physical end state. Figure 6.8 shows that the uncertainty in the free energy estimate to move from the unsampled

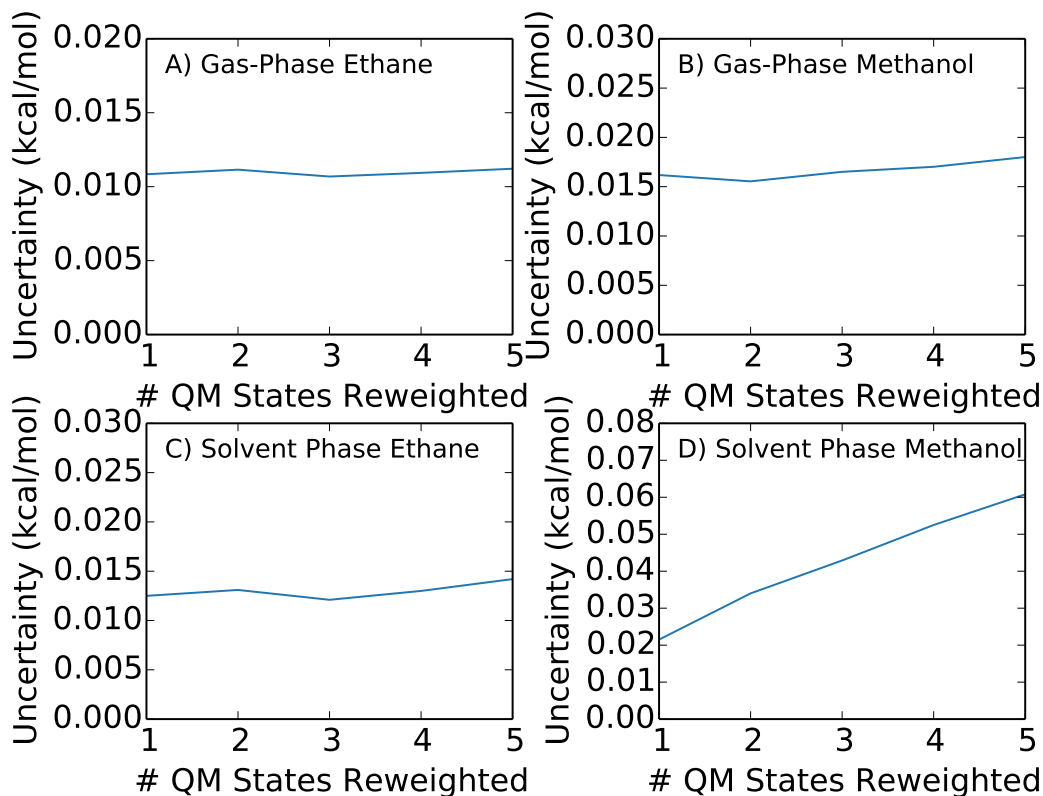


Figure 6.8: The bootstrapped uncertainty in the solvation free energy estimate with three-stage MBAR are lowest when all QM/MM re-evaluations are done only at the physical MM end states where the overlap to the QM states is highest. The uncertainties increase when re-evaluated samples are moved away from the physical MM end states. The increase is most pronounced for solvent phase methanol, because only the physical MM end state has high overlap to the QM/MM state. A total of 15,000 samples are taken evenly from each reweighted intermediate state and the uncertainties as a function of the number of reweighted states are shown for A) gas phase ethane, B) gas phase methanol, C) solvent phase ethane, and D) solvent phase methanol.

QM/MM state to the farthest reweighted MM state is lowest when all energy evaluations are done at the physical end states of the thermodynamic pathway.

We note that the approach of reweighting only a single MM state to the QM/MM Hamiltonian is just the traditional Zwanzig reweighting scheme. Therefore, in this system of ethane and methanol, the solvation free energy difference is best estimated using the Zwanzig reweighting method rather than a multistate approach such as NBB or full MBAR for a fixed amount of energy re-evaluations. This is consistent with the recent study by Jia et al. in which the Zwanzig approach was seen to outperform the NBB approach [306]. When samples are removed from the physical end states and spread out into the intermediate states, the uncertainty in the free energy estimate tends to increase because the physical end states have the highest overlap to the QM/MM potential. This effect is most pronounced in the solvent phase methanol state because only the physical MM end state has high overlap to the QM/MM state. A rigorous minimization of the uncertainty as a function of the distribution of re-evaluated configurations is outside the scope of the present study. However, this analysis demonstrates that to minimize the uncertainty and maximize computational efficiency, the energy re-evaluations should be conducted at the state with the highest overlap, which in this system is the end states of the thermodynamic pathway. It is difficult to determine the highest overlap state without running any simulations. However, a set of short simulations could in theory be used to identify the state with the highest overlap, and the subsequent full production simulations and QM energy re-evaluations can be run from this state alone.

We briefly note that if the mixed Hamiltonian MM states do not have good overlap with adjacent intermediate states, the overall uncertainty in the QM/MM free energy will be large regardless of the choice of state to reweight to the QM/MM Hamiltonian.

6.4.2 SAMPL4 subset

The analysis of the solvation free energy difference in the ethane-methanol system indicates that the uncertainties using MBAR are lower than the uncertainties using NBB for the same set of data. However, this relative uncertainty comparison could in theory reflect features specific

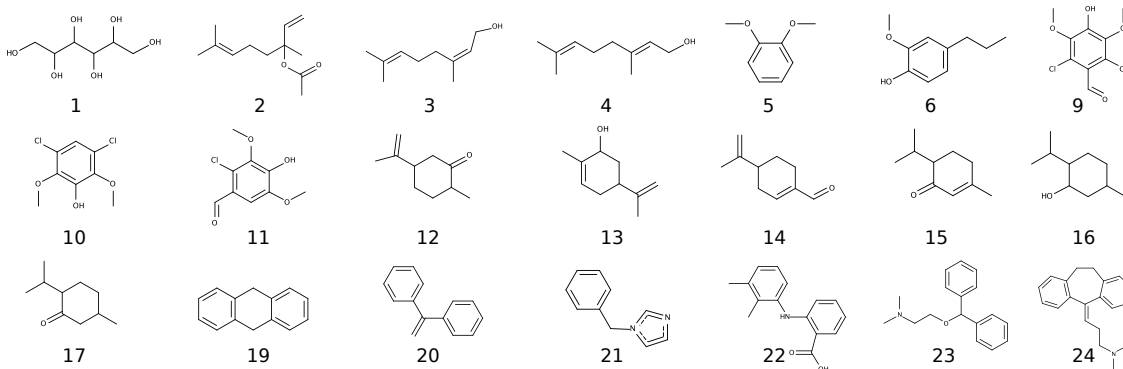


Figure 6.9: The molecules in the blind SAMPL4 solvation dataset examined in this work.

to the ethane-methanol system rather than intrinsic properties of the analysis methods. To further generalize the results, we applied the same solvation free energy difference estimates to the SAMPL4 small molecule solvation dataset (shown in Figure 6.9) with both three-stage MBAR and NBB-indirect. To save computational expense, only the final two states of the path connecting the SAMPL4 molecules were re-evaluated in the QM/MM potential and therefore it is not possible to calculate the free energy difference using either NBB-direct or MBAR. However, the prior analysis with the ethane-methanol system suggests that reweighting with only 2 states is likely to be statistically equivalent to the full MBAR calculation (Figure 6.7).

The solvation free energies and uncertainties for each SAMPL4 molecule with MM-only, NBB-indirect, three-stage MBAR(1), and three-stage MBAR(2) are shown in Figure 6.10 and Figure 6.11. The bootstrapped uncertainty values using three-stage MBAR(2) are lower than the corresponding uncertainties with NBB-indirect for 14 of the 21 molecules in the SAMPL4 molecule set. The uncertainty with NBB is on average 15% larger than the uncertainty with three-stage MBAR(2), suggesting that three-stage MBAR(2) is 28% more efficient at estimating the free energy difference using the same set of data. It is worth noting that the uncertainties with three-stage MBAR(1) are similar to those with three-stage MBAR(2) despite having half as many energy re-evaluations in the QM/MM Hamiltonian. The average uncertainty with three-stage MBAR(1) is 21% smaller than the average uncertainty with three-stage MBAR(2)

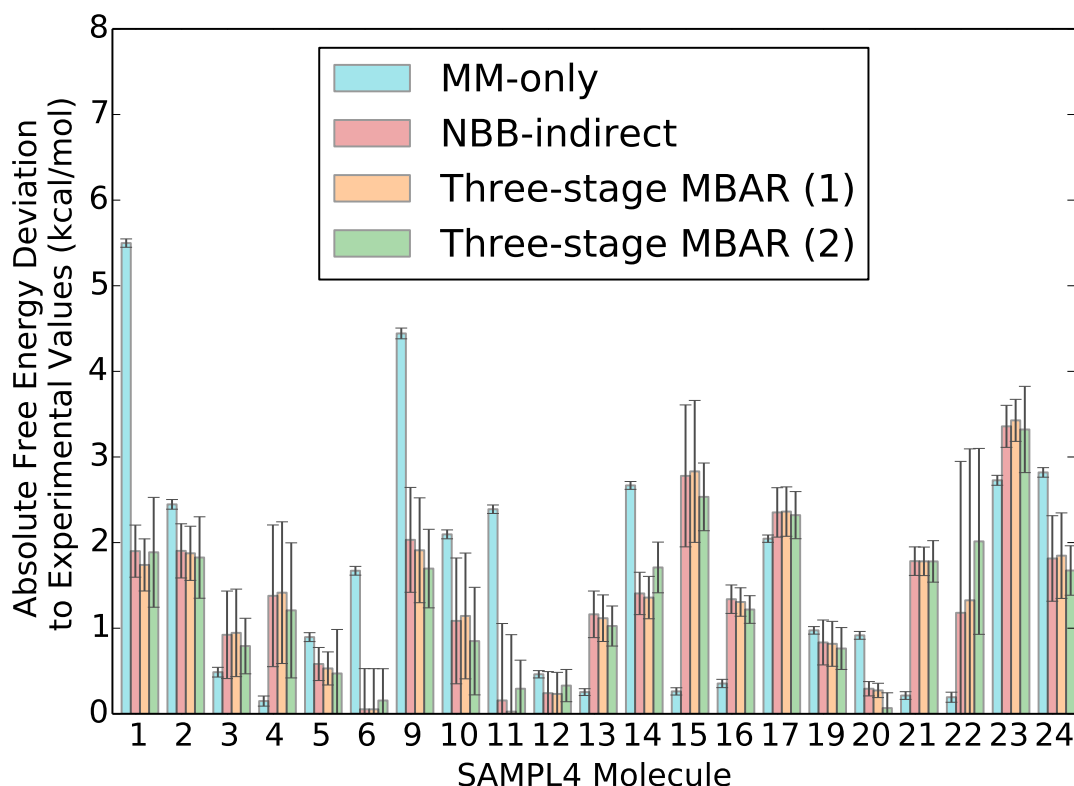


Figure 6.10: The absolute deviation between the experimental solvation free energy and the theoretical estimate are lower on average in the QM/MM potential than in the MM potential. The deviations in the estimates with MBAR are slightly lower on average than those with NBB. Uncertainties for each reweighting method are bootstrapped variances of the solvation free energy estimate. The parenthetical values for the three-stage MBAR method indicate the number of sampled states reweighted to each unsampled state.

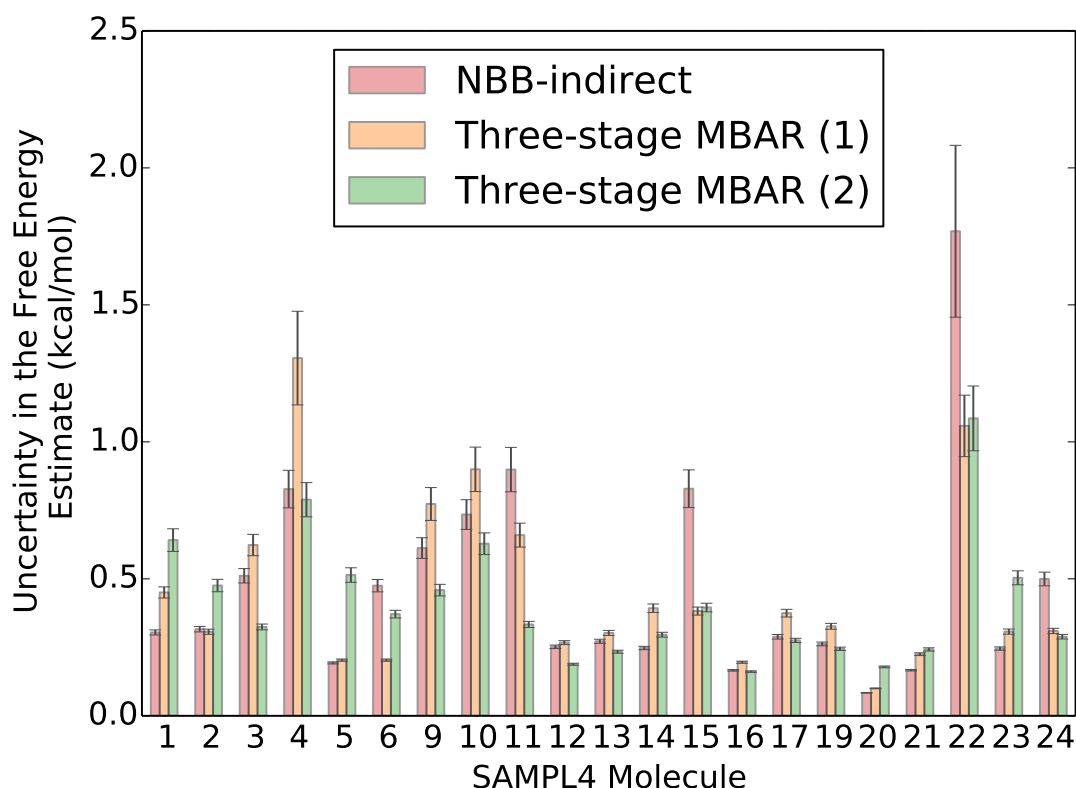


Figure 6.11: The uncertainties in the solvation free energy estimate of the molecules in the SAMPL4 dataset calculated with bootstrapping are larger on average with the NBB-indirect method ($0.47 \text{ kcal}\cdot\text{mol}^{-1}$) than with three-stage MBAR ($0.41 \text{ kcal}\cdot\text{mol}^{-1}$). Uncertainties in the uncertainties were estimated through standard error propagation assuming a normal distribution. The parenthetical values for the three-stage MBAR method indicate the number of sampled states reweighted to each unsampled state.

for a fixed number of energy re-evaluations assuming that the uncertainties are proportional to the square root of the number of QM/MM calculations. This suggests that in the SAMPL4 dataset the Zwanzig approach is best for a fixed amount of QM/MM energy calculations, similar to the results shown earlier for the ethane-methanol system.

The larger uncertainties with NBB-indirect relative to three-stage MBAR(2) are qualitatively consistent with the trend observed in the ethane-methanol system. However the uncertainties with the two methods in the SAMPL4 dataset are noticeably closer in magnitude relative to the uncertainties in the ethane-methanol system. The MM and QM/MM Hamiltonians have significantly less overlap in the SAMPL4 molecules than in the ethane and methanol systems based on the larger uncertainties in the free energy estimate (Figure 6.11) as well as the smaller number of effective samples (Figure 6.12 and 6.13). The average number of effective samples in the solvent phase for the SAMPL4 molecules is 8.7, which is significantly lower than the solvent phase effective samples for ethane (5426) and methanol (779). Correspondingly, the average uncertainty in the solvation free energy for the SAMPL4 molecules with three-stage MBAR(2) is $0.41 \text{ kcal}\cdot\text{mol}^{-1}$, which is more than an order of magnitude larger than the uncertainty in the solvation free energy difference between ethane and methanol ($0.021 \text{ kcal}\cdot\text{mol}^{-1}$). The bootstrapped uncertainty values for the SAMPL4 molecule solvation free energy estimates were on average 28% larger than the analytical uncertainty estimates calculated with MBAR. This discrepancy in error estimates suggests that the sampling in the SAMPL4 molecule dataset has not reached the point where asymptotic statistical variances apply. We therefore conclude that both NBB-indirect and three-stage MBAR(2) provide similar variances for systems with low overlap, but MBAR and three-stage MBAR(2) provide significantly lower variance estimates than NBB-indirect for systems with high overlap to the expensive potential.

We also note that in all molecules where NBB has a noticeably lower variance than three-stage MBAR(2) in Figure 6.11, the number of effective samples for these molecules decreases when using samples from both the physical state 0 and the charge-scaled (by 0.95) state 1, rather than just state 0 (Figure 6.12 and 6.13). A *decrease* in the number of effective samples resulting from the *addition* of configurations can only be caused by a significant shift in the

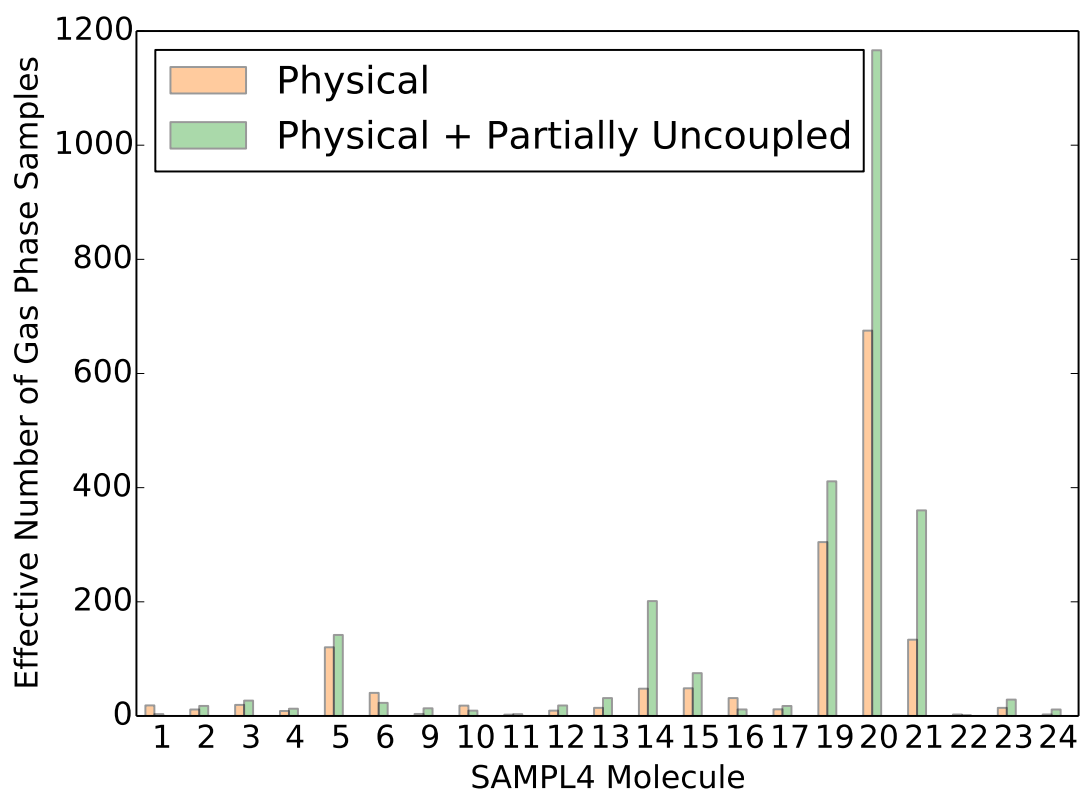


Figure 6.12: The number of effective samples in the QM/MM potential for the SAMPL4 molecules in the gas phase vary significantly and are all lower on average than the number of effective samples for ethane (6413) and methanol (7622). The numbers of effective samples were calculated using just the configurations from state 0 as well as using all configurations from state 0 and state 1.

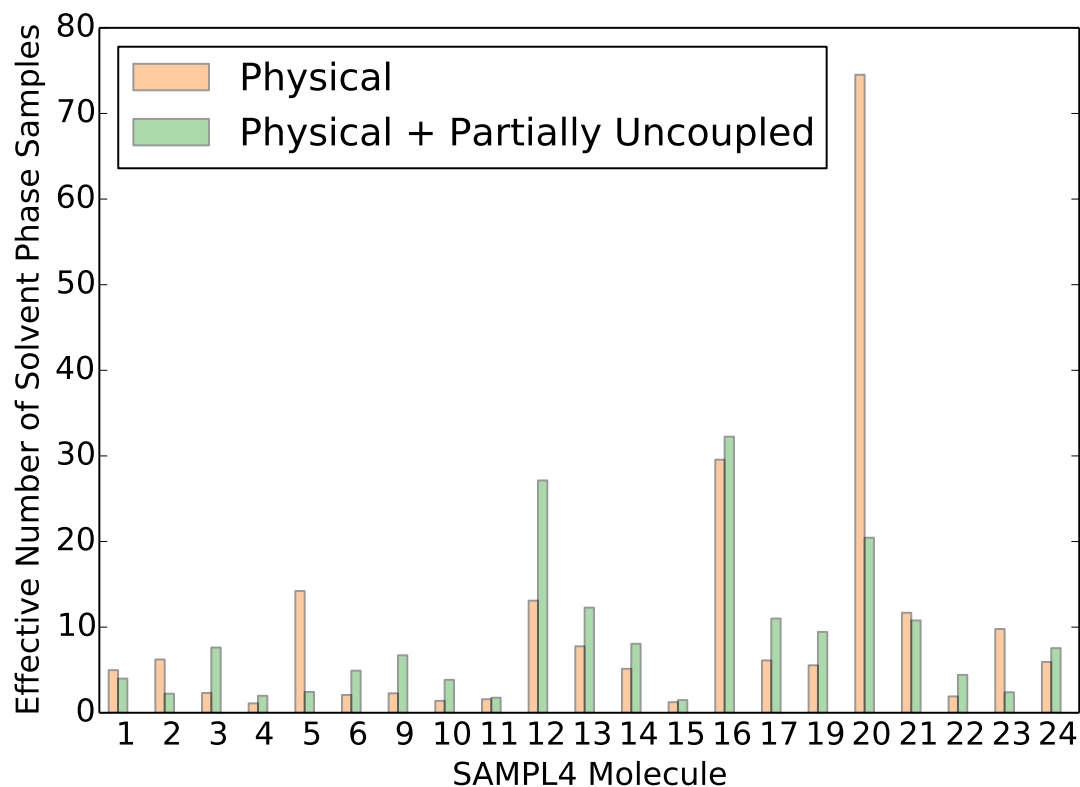


Figure 6.13: The number of effective samples in the QM/MM potential for the SAMPL4 molecules in the solvent phase are significantly lower on average than the number of effective samples for ethane (5426) and methanol (779). The numbers of effective samples were calculated using just the configurations from state 0 as well as using all configurations from state 0 and state 1. In 6 of the 21 molecules, the number of effective samples decreases when adding in the samples from state 1, suggesting that the QM/MM free energy estimates need significantly more samples before converging.

free energy estimate, as seen by inspection of the formulas for the weights (equation 6.13) and effective samples (equation 6.12). The free energy estimates in these molecules have therefore not yet well converged with the current sampling, and any estimate of the free energies or the uncertainties are likely to be spurious.

The large uncertainties and change in the number of effective samples can in theory be used as a quick test of whether the target state is being effectively sampled by the cheaper point-charge potentials. One can simply reweight using three-stage MBAR(1) as well as three-stage MBAR(2). If the number of effective samples decreases when adding in an additional state in the reweighting process or the free energy estimate changes beyond the uncertainty, this strongly suggests that the free energy estimate has not converged. If the number of effective samples increases and the free energy estimate remains unchanged, this suggests that the overlap between the point-charge and target potential is reasonable.

The overlap between the MM and QM/MM Hamiltonians in the SAMPL4 molecules could in theory be improved by designing new MM potentials for each molecule which more closely resemble the QM/MM potential energy surface. Poor overlap was observed even in the gas phase for some of the SAMPL4 molecules, as indicated by less than 10 gas phase effective samples for 4 of the molecules in the dataset. If gas phase simulations already have low overlap, it is unlikely that the overlap will be high when the solvent water molecules are added. This lack of configurational overlap even without the effects of the water molecules suggests that there are significant differences in even the intramolecular interactions between the MM and QM/MM Hamiltonians for these molecules. Indeed, the 5 SAMPL4 molecules with the highest number of gas phase effective samples have rigid aromatic and double-bonded structures which enforce similar intramolecular geometries at both the MM and QM/MM level of theory. Adjusting the intramolecular parameters in the MM potential to match the QM/MM Hamiltonian is likely to increase the configurational overlap in the remaining molecules. In Chapter 4, we showed that the process of designing cheaper potentials to match more expensive polarizable energy functions can indeed be carried out successfully in reweighting crystal polymorphs [119]. Thus it is very likely that such efforts will succeed in gas phase reweighting,

and are reasonably likely to improve overlap in solution phase systems as well.

The solvation free energy differences in the QM/MM Hamiltonian for the SAMPL4 molecules agree favorably with the experimental values, as shown in Figure 6.10. When the QM/MM Hamiltonian with three-stage MBAR(2) is used instead of the pure MM potential, the RMSD between the theoretical and experimental solvation free energy differences decreases from 2.17 kcal·mol⁻¹ to 1.57 kcal·mol⁻¹ and the mean average deviation decreases from 1.62 kcal·mol⁻¹ to 1.33 kcal·mol⁻¹. This is consistent with the findings of König et al. [257] and Essex et al. [249] who saw that the deviation to experimental values decreased when reweighting to the QM/MM Hamiltonian. It is worth noting that although the deviation to experiments is smaller on average with the QM/MM potential, the estimates with only the MM energies were closer in 9 of the 21 SAMPL4 molecules as well as in the ethane-methanol system.

Both NBB-indirect and three-stage MBAR are asymptotically unbiased estimators, and for the SAMPL4 dataset both methods provide essentially equal improvement over the MM potential with respect to the experimental values. With three-stage MBAR(2), the average RMSD is 1.57 kcal·mol⁻¹ and with NBB-indirect and three-stage MBAR(1) the average RMSD is 1.60 kcal·mol⁻¹. The mean average deviation with three-stage MBAR(1) and three-stage MBAR(2) is 1.35 kcal·mol⁻¹ and 1.33 kcal·mol⁻¹, respectively, and the mean average deviation with NBB-indirect is 1.36 kcal·mol⁻¹.

6.4.3 Harmonic Oscillators

As a final comparison of the NBB, three-stage MBAR(1), and three-stage MBAR(2) estimators (which in this test case is equivalent to MBAR), all three methods were used to calculate the free energy difference between a sampled and unsampled harmonic oscillator through an intermediate sampled harmonic oscillator. In the first test case, sampled harmonic oscillators with weak spring constants are reweighted to stiffer harmonic oscillators with larger spring constants. In the second test case, the two stiffest oscillators are sampled and reweighted to the oscillators with weaker spring constants. In the final test case, all oscillators have the same

spring constant, but the two sampled oscillators and the unsampled oscillators have increasingly offset equilibrium positions.

In the first two cases, for all values of the spring constant the uncertainty with MBAR is lower than the uncertainty with NBB-indirect (Figure 6.14). The uncertainty is highest with the three-stage MBAR(1) approach. However this mostly reflects the fact that half as many sampled configurations are re-evaluated in the target state with the Zwanzig approach, corresponding with only one oscillator being used in the reweighting rather than both oscillators for NBB-indirect and three-stage MBAR(2). When the three-stage MBAR(1) approach is run with twice as many configurations sampled from the first oscillator, the free energy difference estimate has a lower uncertainty than both NBB-indirect and three-stage MBAR(2). This is consistent with the results from the ethane-methanol system and the SAMPL4 system, where the lowest variance free energy estimate (for a fixed amount of re-evaluations) came from placing all the re-evaluations into the one state with the highest overlap to the target unsampled state, which is equivalent to the Zwanzig approach (Figure 6.8). For all methods examined with these harmonic oscillators, the bias in the free energy estimate was negligible.

In the final case, when the offset of the unsampled oscillators is nonzero, a systematic error appears in the free energy estimate for all methods, and this bias becomes more significant as the offset increases (Figure 6.15). For a sufficiently large offset, the variance in the free energy difference estimate for NBB-indirect is smaller than the variance in the estimate with three-stage MBAR(2). However, this change occurs at an offset where the systematic error is so large that the free energy estimate is essentially meaningless. This is consistent with the findings from the SAMPL4 molecule dataset in which NBB-indirect produced lower variance free energy estimates only for molecules with poor or biased sampling.

Finally, we note that the lower uncertainties for NBB-indirect relative to three-stage MBAR(1) occurs in the limiting case that $N_1 = N_2$, where N_i is the number of samples drawn from sampled oscillator i and oscillator 1 represents the sampled state farther away from the unsampled states. When $N_1 > N_2$ the uncertainty for NBB-indirect is larger than both three-stage MBAR (1) and three-stage MBAR (2) (see Appendix D section D.1, and Fig. D.3).

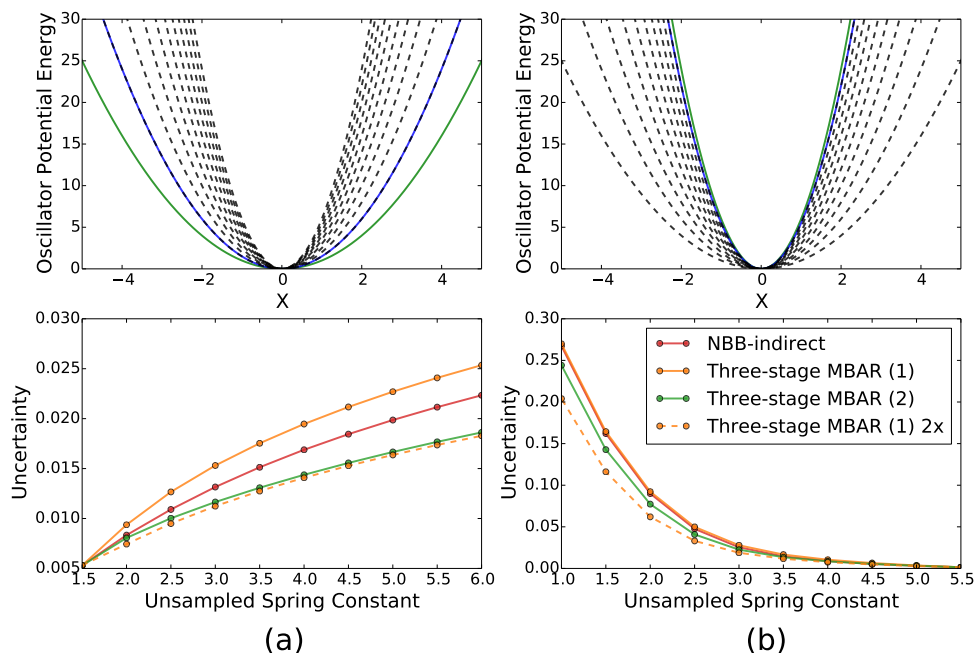


Figure 6.14: The uncertainties in the free energy difference between the sampled and unsampled oscillators increase for all methods as the spring constant of the sampled and unsampled oscillators become farther apart. The uncertainties are lower with three-stage MBAR (2) than with NBB-indirect for all values of the spring constant of the unsampled oscillator in each case. For an equal amount of energy re-evaluations, the method with the lowest variance involves placing all the re-evaluations in the sampled oscillator with the highest overlap and reweighting with the Zwanzig equation. The set of harmonic oscillators are shown above in which the offset is identical and the spring constant of the oscillators changes. The green and blue solid lines represent the first and second sampled oscillator, and the dashed lines indicate the set of unsampled oscillators. The sampled oscillators are reweighted to the unsampled oscillators in two cases (a) the two oscillators with the smallest spring constant are sampled and (b) the two oscillators with the largest spring constant are sampled.

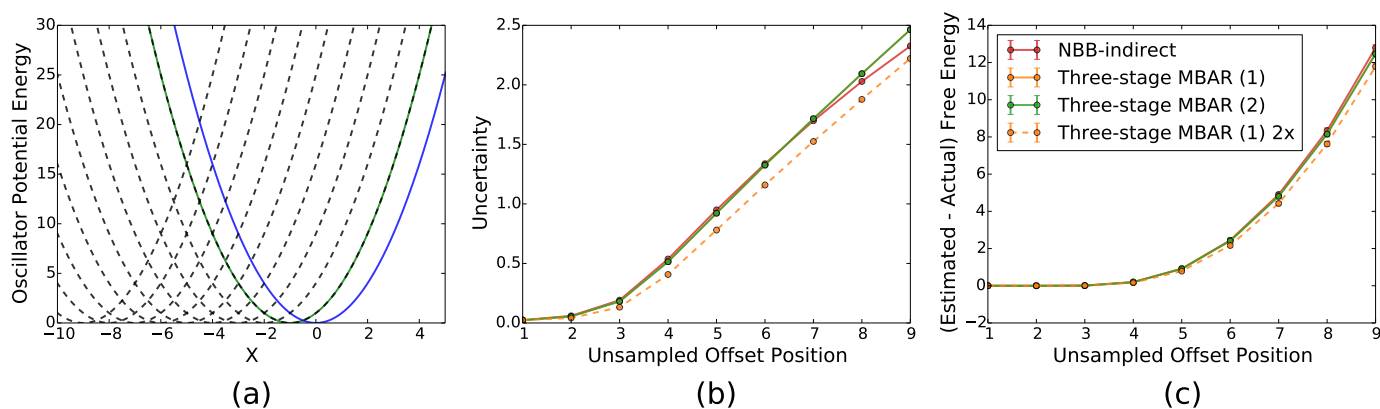


Figure 6.15: The uncertainties in the free energy estimates for offset oscillators are relatively large for all methods. The variance in the estimate for unsampled offsets above 6 are lower for NBB-indirect than three-stage MBAR (2). However, the bias in the estimate for offsets above 6 are much larger than the uncertainties and all estimates with all methods are essentially meaningless. The set of offset oscillators with identical spring constants are shown on the left. The green and blue lines represent the first and second sampled oscillators and the dashed lines indicate the set of unsampled oscillators.

6.5 Conclusions

The solvation free energy difference between methanol and ethane was calculated in a QM/MM Hamiltonian by sampling in a cheap MM potential followed by Hamiltonian reweighting. The Hamiltonian reweighting was carried out using a variety of different algorithms derived from the Multistate Bennett Acceptance Ratio method and the Non-Boltzmann Bennett method, and the MBAR-based methods produce on average smaller uncertainties than the NBB-derived methods for identical sets of data. In addition to lower uncertainties, the MBAR method provides a more natural way to estimate the free energies between two or more unsampled states because the final estimate does not depend on the path connecting the unsampled states. Additionally, one needs to identify which MM state connects to the unsampled QM/MM state in NBB, but no such identification is required with MBAR.

The total solvation free energy difference estimate for the ethane-methanol system is the same, within uncertainty, when reweighting with only the physical end state configurations rather than reweighting all configurations from all intermediate states between ethane and methanol for both NBB and MBAR. The overall computational expense is significantly reduced when reweighting with only the configurations sampled at the end states because it avoids the costly step of re-evaluating the energy of all intermediate configurations in the expensive QM/MM Hamiltonian. Additionally, the lowest variance free energy estimate for a fixed amount of QM/MM energy evaluations occurs when all of the evaluations are conducted at the physical MM end states where the overlap to the corresponding QM/MM states is the highest. Both of these results motivate concentrating the energy re-evaluations on only the configurations sampled at the state with the highest overlap to the expensive potential in order to maximize computational efficiency. Although any intermediate state could, in theory, have the highest overlap with the QM states, it is natural to expect that a fully physical MM state (rather than an unphysical mixed state) would best represent the configurations in a fully physical QM state.

The solvation free energy differences were also estimated for molecules in the SAMPL4

small molecule solvation dataset, and the uncertainties with three-stage MBAR with 2 reweighted states are lower than with NBB-indirect in 14 of the 21 molecules tested. However, the difference in the uncertainties between three-stage MBAR and NBB-indirect were smaller for the SAMPL4 molecule dataset than for the methanol-ethane system. MBAR has been previously shown to have the lowest asymptotic variance possible among all reweighting estimators [188]. We therefore hypothesize that the deviation from these asymptotic limits results from very poor sampling in the SAMPL4 dataset relative to the sampling in the methanol-ethane system. This hypothesis is supported by lower overall uncertainties and higher numbers of effective samples for methanol and ethane when compared with the molecules in the SAMPL4 dataset. It is also supported by the fact that molecules with a smaller statistical error with NBB than with MBAR have low, even single digit numbers of effective samples and free energy estimates that are not fully converged with the available samples as shown in Figures 6.12 and 6.13.

The observation of higher uncertainties for NBB-indirect relative to MBAR is also affirmed by a simple system of 1-D harmonic oscillators in which two sampled oscillators are reweighted to a series of unsampled oscillators. In all examined cases with low bias, the MBAR derived uncertainties are lower than the NBB-indirect uncertainties for a fixed amount of energy re-evaluations. The NBB-indirect method produces lower variance estimates only in cases where the bias is substantially higher than the uncertainty, similar to the results observed in some of the SAMPL4 molecules. In all three systems examined, the uncertainties are lowest for a fixed number of QM/MM calculations when all energy re-evaluations are placed at a single MM state with the highest overlap to the QM/MM Hamiltonian. This single-step reweighting approach is the three-stage MBAR(1) method, which is equivalent to the standard Zwanzig approach plus a multistate reweighting estimate of the cheap intermediate states.

Both NBB and MBAR provide an asymptotically unbiased estimate of the true free energy of the model, and both estimators produce solvation free energy differences in the QM/MM potential that more closely agree with experiments than the MM potential estimates on average. The mean average deviation in the SAMPL4 dataset decreases from $1.62 \text{ kcal}\cdot\text{mol}^{-1}$

to $1.36 \text{ kcal}\cdot\text{mol}^{-1}$ and $1.33 \text{ kcal}\cdot\text{mol}^{-1}$ when moving to the QM/MM potential with NBB-indirect and three-stage MBAR(2), respectively. Given the approximations used to construct the QM/MM Hamiltonian, however, the improvement to experiments relative to the MM potential may be fortuitous.

Based on the results presented in this study, we recommend that to reduce the variance in future free energy difference estimates in unsampled potentials using Hamiltonian reweighting, the sampling and energy re-evaluations should be focused primarily on the thermodynamic states with the highest overlap to the target potential, which are usually the cheaper potential's physical end states. We also recommend that the subsequent reweighting analysis should be carried out with MBAR or the three-stage variant, which is equivalent to the Zwanzig approach at the end state in the limit of one reweighted state. We note, however, that without a sufficiently close match in configuration space overlap between the MM and QM/MM Hamiltonians, no reweighting method works particularly well. The results in this chapter do, however, suggest that reweighting of crystal polymorphs to more expensive potentials with quantum mechanical effects may be within reach in the near future.

Appendix A

Appendix A

A.1 Derivations of Equations

A.1.1 Relating the Helmholtz Free Energy to the Gibbs Free Energy

The statistical mechanics definition of the Gibbs Free Energy G is generally expressed as an integral over phase space:

$$G = -k_B T \ln Z = -k_B T \ln \frac{1}{V_0 N! h^{3N}} \int_{\Gamma} e^{-\beta(U+PV)} d\Gamma \quad (\text{A.1})$$

Where Z is the grand canonical partition function, U is the total system energy, P is the external pressure, V is the current system volume, Γ represents all possible positions and momenta of the particles in the system as well as all possible system volumes, N is the total number of particles in the system, h is a dimensionality constant often taken to be Plank's constant, and V_0 is a reference volume. The multidimensional integral Γ in equation A.1 can be rewritten as an integral over all positions and momenta for system volume V inside of an integral over all possible volumes:

$$G = -k_B T \ln \frac{1}{V_0 N! h^{3N}} \int_V \int_X \int_p e^{-\beta(U+PV)} dp dX dV = -k_B T \ln \frac{1}{V_0} \int_V e^{-\beta PV} Q(V) dV \quad (\text{A.2})$$

where X and p represents all possible positions and momenta of the particles, respectively,

within the system volume V , and $Q(V)$ is the canonical partition function. The canonical partition function is substituted into equation A.2 through the following definition:

$$Q(V) = \frac{1}{N!h^{3N}} \int_X \int_p e^{-\beta U} dp dX = e^{-\beta A(V)} \quad (\text{A.3})$$

where A is the Helmholtz free energy. By combining equation A.2 with A.3 we arrive at the relation between the Gibbs free energy and Helmholtz free energy presented as equation 2.14 in the main text.

$$G = -k_B T \ln \frac{1}{V_0} \int_V e^{-\beta(A(V)+PV)} dV \quad (\text{A.4})$$

A.1.2 Computing free energy differences at multiple temperatures with MBAR

The fundamental quantity computed by the Multistate Bennett Acceptance Ratio (MBAR) is the reduced free energy:

$$G_k = \beta_T f_k \quad (\text{A.5})$$

where G_k is the dimensionalized free energy of state k , f_k is the reduced free energy, and β_T is the Boltzmann factor at the temperature T of state k . The difference in reduced free energy between any two states 1 and 2 at temperatures T_{ref} and T , respectively, can be expressed as:

$$f_2 - f_1 = \beta_T G_2(T) - \beta_{T_{ref}} G_1(T_{ref}) \quad (\text{A.6})$$

The last equation can be rewritten to express the Gibbs free energy of state 2 at temperature T :

$$G_2(T) = \frac{1}{\beta_T}(f_2 - f_1) - \frac{\beta_{T_{ref}}}{\beta_T}G_1(T_{ref}) \quad (\text{A.7})$$

If the two states at different temperatures are of the same polymorph i then $G_2(T)$ and $G_1(T_{ref})$ collapse to $G_i(T)$ and $G_i(T_{ref})$. We can therefore express the free energy change between T and T_{ref} of polymorph i and polymorph j as:

$$G_i(T) = \frac{1}{\beta_T}(f_i(T) - f_i(T_{ref})) - \frac{\beta_{T_{ref}}}{\beta_T}G_i(T_{ref}) \quad (\text{A.8})$$

$$G_j(T) = \frac{1}{\beta_T}(f_j(T) - f_j(T_{ref})) - \frac{\beta_{T_{ref}}}{\beta_T}G_j(T_{ref}) \quad (\text{A.9})$$

By subtracting equation (A.8) from equation (A.9) and rearranging terms we arrive at the free energy difference between two polymorphs at any temperature T shown in equation (2.17) of the main text:

$$\Delta G_{ij}(T) = G_j(T) - G_i(T) = k_B T \left(\Delta f_{ij}(T) - \Delta f_{ij}(T_{ref}) \right) + \frac{T}{T_{ref}} \Delta G_{ij}(T_{ref}) \quad (\text{A.10})$$

Appendix B

Appendix B

B.1 Table of Experimental Polymorph Data

B.2 Extrapolation of Free Energy from Finite Temperature to 0 K

The crystals can not be actually simulated at a temperature of exactly 0 K. Furthermore, the multistate reweighting calculation in equation 11 of ref 119 used to generate the temperature-dependent free energy profiles cannot include configurations at $T = 0$ K. Therefore, the crystals need to be simulated down to a temperature in which the free energy linearly transitions from finite temperature to 0 K. The linearity of the free energy was determined by extrapolating the free energy from finite temperatures down to 0 K using the equation:

$$\Delta G(0K) = \Delta G(T) + T\Delta S(T) \quad (\text{B.1})$$

where $\Delta G(T)$ is the free energy difference at some temperature, T , and $\Delta S(T)$ is the estimated entropy. In all systems, the estimated free energy at 0 K was calculated using equation B.1 for both $T = 10$ K and $T = 20$ K. These estimates were then compared against the true minimized lattice energy difference. Table B.2 below depicts the lattice energy difference as well as the extrapolated free energy for all systems. In all cases, the extrapolated free energy is within

0.005 kcal·mol⁻¹ of the minimized lattice energy difference indicating that the free energy changes linearly in the unsampled temperature range between 0 K and 10 K.

B.3 Complete PSCP Cycle Closure

The relative free energies as a function of temperature for all crystals were computed by using a reference free energy as well as simulations over a range of temperature combined with multi-state reweighting using MBAR. The equations for computing this temperature-dependent free energy are outlined in section 2.5 and in Ref 119. The free energies produced from this process were validated by running multiple reference free energy calculations at different temperatures and examining the cycle closure relative to the uncertainty. The relative free energies at T_1 estimated with a direct PSCP calculation and with a PSCP at T_2 followed by MBAR down to T_1 are shown for all systems below.

Crystal Name	Adenine Form I	Adenine Form II
CCSD Refcode	KOBFUD	KOBFUD01
Stable Temperature	0K - 523K	> 523K
Temperature Ref	92	92
Supercell Dimensions	3x1x4	3x1x2
Crystal Name	Carbamazepine Form I	Carbamazepine Form III
CCSD Refcode	CBMZPN11	CBMZPN02
Stable Temperature	> 435K	0K-435K
Temperature Ref	93	93
Supercell Dimensions	4x2x1	4x2x2
Crystal Name	Paracetamol Form I	Paracetamol Form II
CCSD Refcode	HXACAN01	HXACAN
Stable Temperature	> 108K	0K-108K
Temperature Ref	94	94
Supercell Dimensions	2x2x3	2x1x3
Crystal Name	Resorcinol α	Resorcinol β
CCSD Refcode	RESORA03	RESORA08
Stable Temperature	0K - 337K	> 337K
Temperature Ref	98	98
Supercell Dimensions	3x4x6	6x3x4
Crystal Name	Cyclopentane Form I	Cyclopentane Form III
CCSD Refcode	N/A	ZZZVYE01
Stable Temperature	> 138K	0K-118K
Temperature Ref	100	100
Supercell Dimensions	2x4x2	2x4x2
Crystal Name	Succinonitrile Form I	Succinonitrile Form II
CCSD Refcode	QOPBED	N/A
Stable Temperature	0K-233K	> 233K
Temperature Ref	101	101
Supercell Dimensions	2x3x4	2x3x4
Crystal Name	Diiodobenzene α	Diiodobenzene β
CCSD Refcode	ZZZPRO03	ZZZPRO04
Stable Temperature	0K-326K	> 326K
Temperature Ref	103	103
Supercell Dimensions	2x3x3	2x3x3
Crystal Name	Piracetam Form I	Piracetam Form III
CCSD Refcode	BISMEV03	BISMEV02
Stable Temperature	> 393K	0K-393K
Temperature Ref	102	102
Supercell Dimensions	3x2x3	2x3x3
Crystal Name	Pyrazinamide Form III	Pyrazinamide Form IV
CCSD Refcode	PYRZIN20	PYRZIN16
Stable Temperature	> 428K	0K-298K
Temperature Ref	95	95
Supercell Dimensions	4x3x2	2x3x4
Crystal Name	Aripiprazole Form I	Aripiprazole Form X
CCSD Refcode	MELFIT01	MELFIT05
Stable Temperature	> 350K	0K - 335K
Temperature Ref	18	18
Supercell Dimensions	3x2x2	3x2x2
Crystal Name	Tolbutamide Form I	Tolbutamide Form II
CCSD Refcode	ZZZPUS04	ZZZPUS05
Stable Temperature	> 385K	0K - 385K
Temperature Ref	96	96
Supercell Dimensions	2x3x2	3x2x1
Crystal Name	Chlorpropamide Form I	Chlorpropamide Form V
CCSD Refcode	BEDMIG	BEDMIG04
Stable Temperature	0K - 393K	> 393K
Temperature Ref	97	97
Supercell Dimensions	1x4x2	2x2x1

Table B.1: Experimental polymorph transition temperatures and literature references, supercell dimensions, and reference codes for crystals examined in Chapter 3.

System	$T = 0K$	$T = 10K$	$T = 20K$
1,4-Diiodobenzene	0.258	0.256	0.257
Resorcinol	-0.278	-0.278	-0.279
Adenine	2.255	2.255	2.256
Pyrazinamide	-0.176	-0.176	-0.176
Carbamazepine	0.565	0.564	0.565
Piracetam	-1.459	-1.459	-1.458
Paracetamol	-0.448	-0.447	-0.446
Cyclopentane	0.351	0.355	0.349
Succinonitrile	1.235	1.234	1.234
Tolbutamide	-2.226	-2.237	-2.230
Chlorpropamide	2.789	2.792	2.787
Aripiprazole	0.596	0.596	0.596

Table B.2: Estimated free energy differences at 0 K computed directly and extrapolated from 10 K and 20 K for all systems. The extrapolated estimates are within roughly $0.005 \text{ kcal}\cdot\text{mol}^{-1}$ of the true minimized lattice energy difference indicating that the free energy changes linearly between 0 K and 10 K. Estimated uncertainties for all extrapolations were on the order $0.0001 \text{ kcal}\cdot\text{mol}^{-1}$

System	T_1	$\Delta G(T_1)$	T_2	$\Delta G(T_2) + \Delta\Delta G(T_2 \rightarrow T_1)$
1,4-Diiodobenzene	100 K	0.242 ± 0.002	200 K	0.241 ± 0.001
Resorcinol	100 K	-0.307 ± 0.002	200 K	-0.303 ± 0.000
Adenine	100 K	2.039 ± 0.002	200 K	2.035 ± 0.000
Pyrazinamide	100 K	-0.359 ± 0.005	200 K	-0.345 ± 0.002
Carbamazepine	100 K	-0.398 ± 0.016	200 K	-0.376 ± 0.005
Piracetam	100 K	-1.459 ± 0.005	200 K	-1.465 ± 0.001
Paracetamol	100 K	0.520 ± 0.011	200 K	0.526 ± 0.001
Cyclopentane	30 K	0.272 ± 0.001	50 K	0.250 ± 0.001
Succinonitrile	100 K	1.122 ± 0.001	200 K	1.054 ± 0.005
Tolbutamide	100 K	-1.995 ± 0.041	200 K	-1.897 ± 0.012
Chlorpropamide	100 K	2.274 ± 0.004	200 K	2.307 ± 0.020
Aripiprazole	300 K	0.346 ± 0.016	200 K	0.409 ± 0.019

Table B.3: Pseudo-supercritical Path free energy cycle closures for all polymorph pairs in this study. The free energy difference estimated directly at T_1 is statistically the same as the free energy difference obtained by running a PSCP at T_2 and connecting this state to T_1 with a series of intermediate simulations and multistate reweighting with MBAR.

Appendix C

Appendix C

C.1 End State Hamiltonian Reweighting for All Polymorphs

The sampling along the PSCP in the expensive potentials is avoided by running the full path directly in a cheaper designed point-charge potential and reweighting the end states from the cheap to target Hamiltonian. The free energy estimate to switch from the sampled cheap potential to the unsampled expensive potential requires high overlap between the two states, which is visualized by viewing the overlap in the histograms of the energy difference sampled in both states.

C.2 Indirect AMOEBA09 Free Energy vs Temperature Convergence

The relative free energy in the AMOEBA09 potential as a function of temperature can be calculated without any direct sampling in AMOEBA09, or can be computed with any amount of direct AMOEBA09 included in the computation. Without any direct AMOEBA09 sampling, a slight bias of around $0.05 \text{ kcal}\cdot\text{mol}^{-1}$ appears in the free energy curve of benzene III. This bias

Polymorph	OPLS-AA	Designed-GROMOS
benzene I	0.000095	0.391
benzene II	0.00019	0.242
benzene III	0.00058	0.541

Table C.1: Overlap to GROMOS54A7 for three benzene polymorphs using both the OPLS-AA and Designed-GROMOS potentials

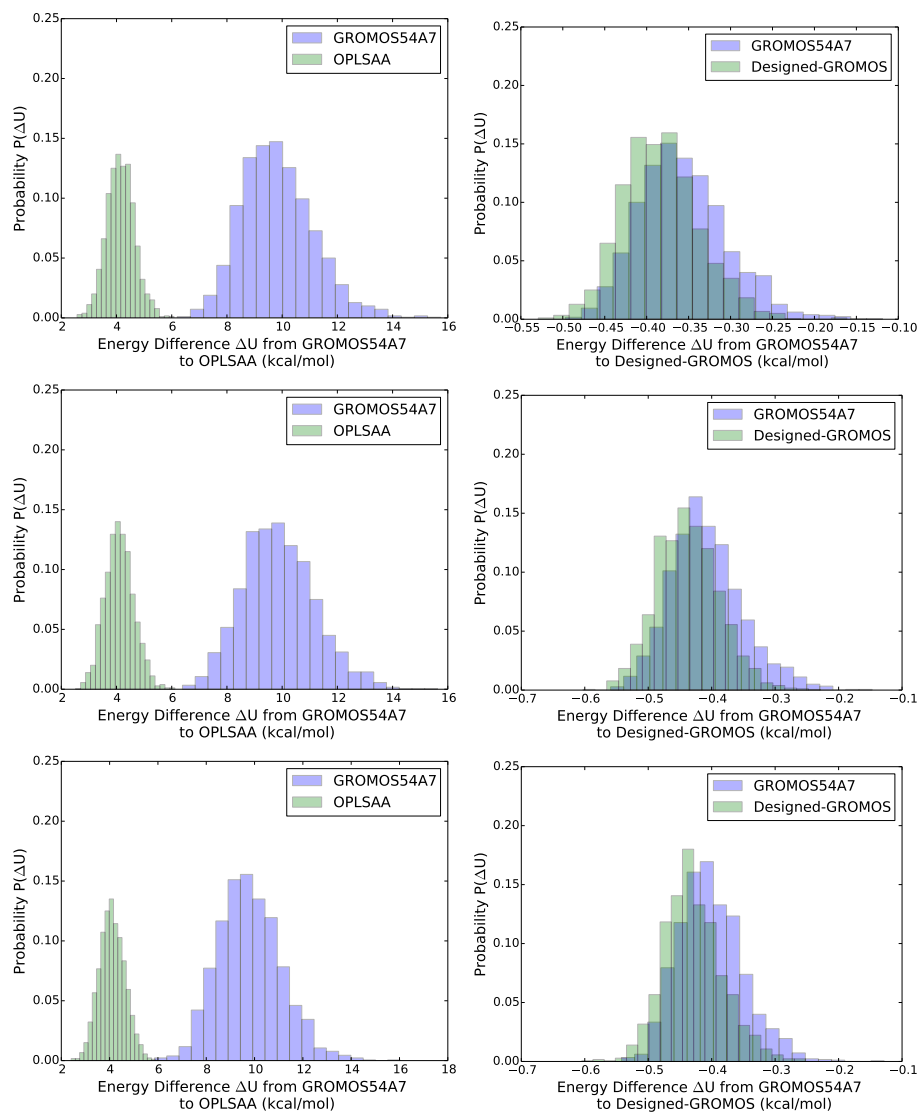


Figure C.1: The OPLS-AA potential has relatively poor overlap with GROMOS54A7, while the Designed-AMOEBA potential has higher overlap with the GROMOS54A7 potential in the solid benzene system for each of the 3 benzene polymorphs.

Polymorph	OPLS-AA	Designed-AMOEBA
benzene I	0.017	0.178
benzene II	0.025	0.089
benzene III	0.019	0.101

Table C.2: Overlap to AMOEBA09 for three benzene polymorphs using various potentials with a hydrogen charge of 0.115

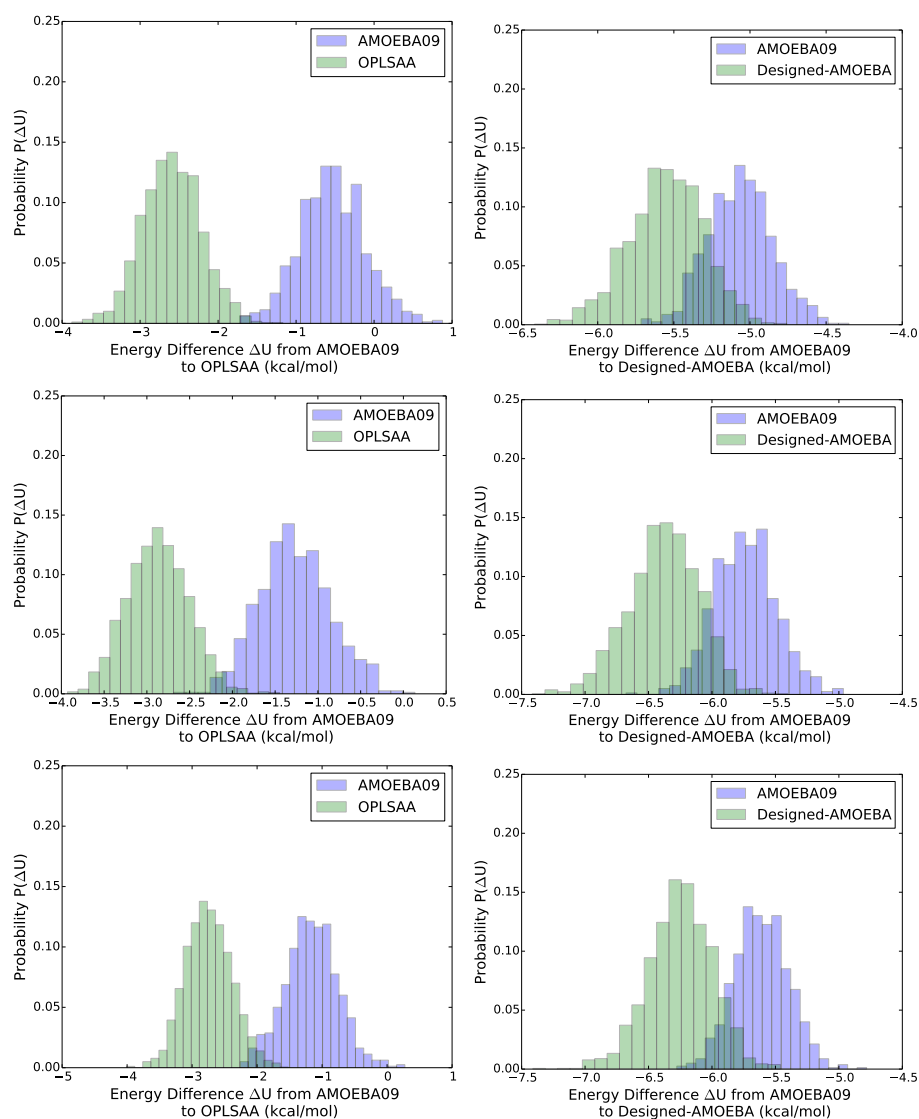


Figure C.2: The OPLS-AA potential has relatively poor overlap with AMOEBA09, while the Designed-AMOEBAA potential has higher overlap with the AMOEBA09 potential in the solid benzene system for each of the 3 benzene polymorphs.

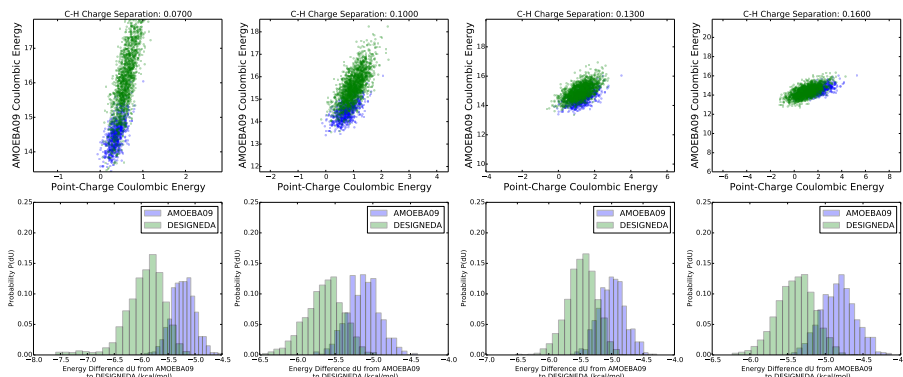


Figure C.3: The configuration space overlap between the Designed-AMOEBA and AMOEBA09 potential can be increased by tuning the charges on the atoms in the benzene molecules in the designed potential. In the Benzene I crystal the highest overlap occurs at the original OPLS-AA hydrogen charge $q_H = 0.115$ (and a corresponding negative carbon charge of -0.115). The distribution sampled from AMOEBA09 is shown in blue and the distribution from the designed point-charge potential is shown in green.

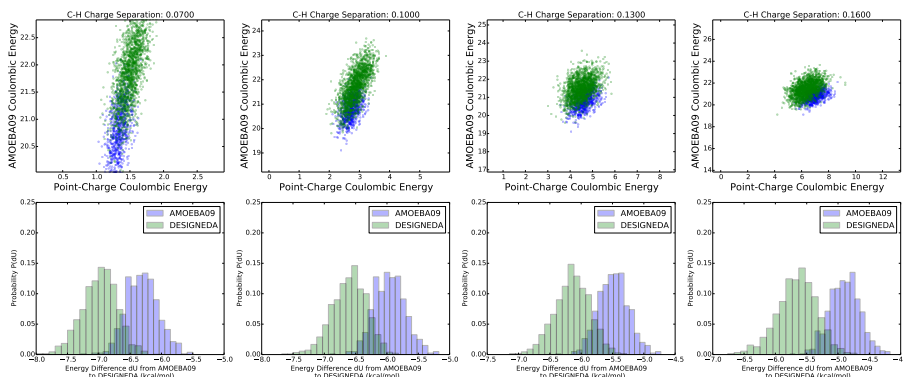


Figure C.4: The configuration space overlap between the Designed-AMOEBA and AMOEBA09 potential can be increased by tuning the charges on the atoms in the benzene molecules in the designed potential. In the Benzene II crystal the highest overlap occurs at a hydrogen charge $q_H = 0.100$ (and a corresponding negative carbon charge of -0.100), compared to the original $q_H = 0.115$ in OPLS-AA. The distribution sampled from AMOEBA09 is shown in blue and the distribution from the designed point-charge potential is shown in green. This figure is reproduced from Figure 4.15 of the main text.

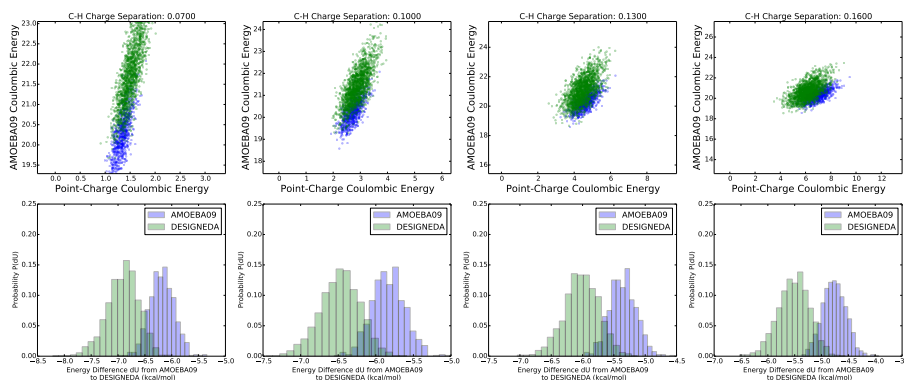


Figure C.5: The configuration space overlap between the Designed-AMOEBA and AMOEBA9 potential can be increased by tuning the charges on the atoms in the benzene molecules in the designed potential. In the Benzene III crystal the highest overlap occurs at at the original OPLS-AA hydrogen charge $q_H = 0.115$ (and a corresponding negative carbon charge of -0.115). The distribution sampled from AMOEBA9 is shown in blue and the distribution from the designed point-charge potential is shown in green.

can be virtually eliminated with just a short amount of sampling in the AMOEBA9 potential. The free energy curve using a range of AMOEBA9 sampling between 2-5ns is shown in Figure C.6. In all cases, the first 1ns of the simulation is discarded for equilibration.

C.3 Potentials and Functional Forms

The OPLS-AA and GROMOS54A7 potentials both follow the traditional point-charge potential form in which the total potential energy is the sum of intramolecular and intermolecular interactions according to:

$$U_{tot} = U_{bond} + U_{angle} + U_{dihedral} + U_{vdw} + U_{elec} \quad (\text{C.1})$$

where U_{tot} is the total potential energy, U_{bond} is the bond stretching term, U_{angle} represents the angle stretching interactions, $U_{dihedral}$ represents the proper and improper dihedral torsions,

U_{vdw} contains the intermolecular van der Waals forces, and U_{elec} contains the Coulombic interactions between all point charges in the system.

In the OPLS potential, the bonded interactions are represented by the following equations

$$U_{bond} = \frac{1}{2} K_b (b - b_0)^2 \quad (C.2)$$

$$U_{angle} = \frac{1}{2} K_\theta (\theta - \theta_0)^2 \quad (C.3)$$

$$U_{dihedral} = \sum_{n=0}^5 C_n (\cos \phi)^n + \sum_{n=1}^4 F_n (1 - (-1)^n \cos(n\phi)) \quad (C.4)$$

where b is the bond distance, b_0 is the equilibrium bond length, K_b is the bond spring constant, θ is the valence angle between 3 atomic centers, θ_0 is the equilibrium angle, K_θ is the angle force constant, ϕ is the dihedral angle between 4 atomic centers, ϕ_0 is the equilibrium dihedral angle, and C_n is the coefficient for the n^{th} term in the dihedral summation. The Van der Waals interactions in OPLS are represented by the sigma-epsilon Lennard-Jones potential with the following form

$$U_{vdw} = 4\epsilon_{ij} \left(\left(\frac{\sigma_{ij}}{r_{ij}} \right)^{12} - \left(\frac{\sigma_{ij}}{r_{ij}} \right)^6 \right) \quad (C.5)$$

where σ_{ij} is the Lennard-Jones sigma, ϵ_{ij} is the potential well depth, and r_{ij} is the interatomic distance between particles i and j .

In the GROMOS54A7 potential, the bonded interactions are described by the following equations

$$U_{bond} = \frac{1}{4} K_b (b^2 - b_0^2)^2 \quad (C.6)$$

$$U_{angle} = \frac{1}{2}K_{\theta}(\cos \theta - \cos \theta_0)^2 \quad (C.7)$$

$$U_{dihedral} = \frac{1}{2}K_{\phi}(\phi - \phi_0)^2 \quad (C.8)$$

where K_{ϕ} is the dihedral force constant. The van der Waals interactions in GROMOS54A7 are represented by the C6-C12 Lennard-Jones potential with the following form

$$U_{vdw} = \frac{C_{ij}^{(12)}}{r_{ij}^{12}} - \frac{C_{ij}^{(6)}}{r_{ij}^6} \quad (C.9)$$

where $C_{ij}^{(6)}$ and $C_{ij}^{(12)}$ are fitted coefficients describing the dispersion and repulsion interactions between particles i and j . Both the OPLS-AA and GROMOS54A7 treat the pairwise Coulombic interactions between two point charges using

$$U_{elec} = \frac{q_i q_j}{4\pi\epsilon_0 r_{ij}} \quad (C.10)$$

where q_i is the point-charge on particle i , ϵ is the permittivity, and r_{ij} is the interatomic distance between the charges.

The AMOEBA force field incorporates the effects of distributed multipoles and induced polarizability in addition to the point-charge electrostatics and intramolecular interactions present in the GROMOS54A7 and OPLS-AA potentials. The total potential energy in AMOEBA is described by the equation:

$$U_{tot} = U_{bond} + U_{angle} + U_{str-bnd} + U_{torsion} + U_{oop} + U_{vdw} + U_{elec}^{perm} + U_{elec}^{ind} \quad (C.11)$$

where $U_{str-bnd}$ characterizes the stretch-bend interactions, $U_{torsion}$ is the dihedral torsion interactions, U_{oop} is the out-of-plane deformation energy, and U_{elec}^{perm} and U_{elec}^{ind} describe the permanent and induced electrostatic interactions, respectively. The AMOEBA potential adopts

the following functional forms for intramolecular interactions:

$$U_{bond} = K_b(b - b_0)^2 \left[1 - 2.55(b - b_0) + \frac{7}{12} 2.55^2(b - b_0)^2 \right] \quad (\text{C.12})$$

$$U_{angle} = K_\theta(\theta - \theta_0)^2 \left[1 - 0.014(\theta - \theta_0) + 5.6 \times 10^{-5}(\theta - \theta_0)^2 - 7 \times 10^{-7}(\theta - \theta_0)^3 + 2.2 \times 10^{-8}(\theta - \theta_0)^4 \right] \quad (\text{C.13})$$

$$U_{str-bnd} = K_{str-bnd}(b - b_0)(b' - b'_0)(\theta - \theta_0) \quad (\text{C.14})$$

$$U_{torsion} = \sum_{n=1}^3 C_n(1 + \cos(n\phi - \phi_0)) \quad (\text{C.15})$$

$$U_{oop} = K_\chi \chi^2 \quad (\text{C.16})$$

where $K_{str-bnd}$ is the force constant for the correlated bond stretch and bend motion, χ is the out of plane dihedral angle, and K_χ is the force constant maintaining planar geometry for planar molecules. The van der Waals interactions in AMOEBA are described by a buffered 14-7 potential originally proposed by Halgren [316]:

$$U_{vdw} = \epsilon_{ij} \left(\frac{1 + \delta}{\frac{r_{ij}}{r_{ij}^0} + \delta} \right)^7 \left(\frac{1 + \gamma}{\left(\frac{r_{ij}}{r_{ij}^0} \right)^7 + \gamma} - 2 \right) \quad (\text{C.17})$$

where ϵ_{ij} is the potential well depth, r_{ij} is the distance between atoms i and j , r_{ij}^0 is the minimum energy distance, and δ and γ are buffering constants which are set to 0.007 and 0.12, respectively.

The electrostatic interactions in AMOEBA09 consist of both permanent and induced interactions between particles. A full description of the interactions and functional forms for the AMOEBA potential can be found in the potential-specific publications [264–266]. However, a brief overview is provided here. Given the set of permanent monopoles (q), dipoles (μ), and quadrupoles (Q) around an atom center, permanent and induced electrostatic interactions are calculated with

$$M_i = [q_i, \mu_{ix}, \mu_{iy}, \mu_{iz}, Q_{ixx}, Q_{ixy}, Q_{ixz}, \dots, Q_{izz}]^T \quad (\text{C.18})$$

$$T_{ij} = \begin{bmatrix} 1 & \frac{\partial}{\partial x_j} & \frac{\partial}{\partial y_j} & \frac{\partial}{\partial z_j} & \cdots \\ \frac{\partial}{\partial x_i} & \frac{\partial^2}{\partial x_i \partial x_j} & \frac{\partial^2}{\partial x_i \partial y_j} & \frac{\partial^2}{\partial x_i \partial z_j} & \cdots \\ \frac{\partial}{\partial y_i} & \frac{\partial^2}{\partial y_i \partial x_j} & \frac{\partial^2}{\partial y_i \partial y_j} & \frac{\partial^2}{\partial y_i \partial z_j} & \cdots \\ \frac{\partial}{\partial z_i} & \frac{\partial^2}{\partial z_i \partial x_j} & \frac{\partial^2}{\partial z_i \partial y_j} & \frac{\partial^2}{\partial z_i \partial z_j} & \cdots \\ \vdots & \vdots & \vdots & \vdots & \ddots \end{bmatrix} \left(\frac{1}{r_{ij}} \right) \quad (\text{C.19})$$

$$U_{elec}^{perm} = M_i^T T_{ij} M_j \quad (\text{C.20})$$

$$\mu_{i,\alpha}^{ind} = \alpha_i \left(\sum_{\{j\}} T_{\alpha}^{ij} M_j + \sum_{\{j'\}} T_{\alpha,\beta}^{ij'} \mu_{j',\beta}^{ind} \right) \text{ for } \alpha, \beta = 1, 2, 3 \quad (\text{C.21})$$

$$U_{elec}^{ind} = -\frac{1}{2} \sum_i (\mu_{i,\alpha}^{ind})^T E_i \quad (\text{C.22})$$

$$(\text{C.23})$$

where T_{ij} is the matrix for the interaction of sites i and j separated by a distance r_{ij} , and α_i is the polarizability of atom i . The set $\{j\}$ is all atomic sites that are outside the molecule containing atom i . The set $\{j'\}$ is all atomic sites other than i itself. E_i is the direct electric field due to the permanent multipoles in set $\{j\}$.

The designed point-charge potentials, labeled Designed-AMOEBA and Designed-AMOEBA, were made to have high overlap with the target potentials by using intramolecular interactions that matched those in GROMOS54A7 and AMOEBA09. The designed point charge potentials follow the general structure of the OPLS-AA point charge potential. The total potential energy

function for both designed potentials is

$$U_{tot} = U_{bond} + U_{angle} + U_{dihedral} + U_{vdw} + U_{elec} \quad (C.24)$$

where U_{tot} is the total potential energy, U_{bond} is the bond stretching term, U_{angle} represents the angle stretching interaction, $U_{dihedral}$ represents the proper and improper dihedral torsions, U_{vdw} is the intermolecular van der Waals forces, and U_{elec} is the Coulombic interactions between all point charges in the system. The bond and angle interactions take the following functional form

$$U_{bond} = \frac{1}{2}K_b(b - b_0)^2 \quad (C.25)$$

$$U_{angle} = \frac{1}{2}K_\theta(\theta - \theta_0)^2 \quad (C.26)$$

where the parameters K_b , b_0 , K_θ , θ_0 , were selected to match the corresponding bond and angle interactions in the GROMOS54A7 and AMOEBA09 potential.

The dihedral interactions in the Designed-AMOEBA potential were taken directly from the GROMOS54A7 forcefield and follow the form:

$$U_{dihedral} = \frac{1}{2}K_\phi(\phi - \phi_0)^2 \quad (C.27)$$

The dihedral interactions in the Designed-Amoeba potential were constructed to match the torsion and out of plane interactions in the AMOEBA potential. The total dihedral interaction takes the form:

$$U_{torsion} = K_\chi \chi^2 \sum_{n=1}^3 C_n (1 + \cos(n\phi - \phi_0)) \quad (C.28)$$

where the first term in equation C.28 matches the out-of-plane interactions and the second term represents the torsion terms in AMOEBA.

The OPLS-AA and Coulombic interactions were used in both the Designed-AMOEBA and Designed-AMOEBA potentials.

The exact parameters for the intermolecular and intramolecular interactions in each of the designed potentials is included as a set of *.itp files which are compatible with the GROMACS simulation package.

C.4 Pseudo-Supercritical Path Calculations

The pseudo-supercritical path calculates the Helmholtz free energy to transform a fully interacting solid crystal into a noninteracting ideal gas at constant volume. Strictly speaking, the noninteracting benzene molecules also contain intramolecular interactions, which are not present for ideal gases, however these contributions cancel between polymorphs. The two steps involved in the transformation process are restraining the atoms to the crystalline position, and turning off all intermolecular interactions. The total potential energy function (reproduced from equation 2.12 is:

$$U_{tot}(\lambda, \gamma) = U_{intra} + \gamma^2 U_{inter} + (1 - \lambda)^4 \frac{k}{2} (x - \bar{x})^2 \quad (\text{C.29})$$

The total Helmholtz free energy change in the PSCP reflects the free energy to drive the harmonic restraint coupling parameter, λ , to unity and driving the intermolecular interaction coupling parameter, γ , to zero.

The free energy profiles for each of the polymorphs along each step of the PSCP in each potential is listed below. Numerical values of the free energy difference in each step of the PSCP are provided in section C.5.

C.5 Benzene Polymorph Free Energies using 4- and 72- Benzene Model

The effect of using a 4-benzene unit cell as the representative crystal model was probed by calculating the relative free energy of the three benzene crystal structures in the point-charge potentials using a 3x3x2 benzene super cell comprised of 72 independently moving benzene molecules. The relative free energies were calculated using the same method presented in the main text of interconverting the polymorphs within a potential using a Pseudo-supercritical Path (PSCP) and switching potentials using both the Zwanzig (EXP) method and the Bennett Acceptance Ratio (BAR) method.

The relative free energies along the PSCP of the three benzene polymorphs studied in the main text for the 72 benzene model in OPLS-AA, GROMOS54A7, and Designed-GROMOS are presented in Tables C.3-C.5. For comparison, the corresponding free energies using the 4-benzene model are also included in the Tables.

The relative free energies to transform the benzene polymorphs from the GROMOS54A7 potential to the Designed-GROMOS potential for both the 4 and 72 benzene supercell models are shown in Table C.9. The free energy changes are statistically the same within uncertainty. This result demonstrates that the entropy loss to go from a 72 benzene super cell to a 4 benzene system are essentially the same in both the GROMOS54A7 and Designed-GROMOS potentials. This result also suggests that in theory the PSCP could be done in with larger system and the free energies to transform between potentials could be calculated with a smaller system where the overlap is higher. We note that the relative PV correction factors were found to be negligible for the 4 benzene system at 200K and were therefore ignored at 100K and in the 72 benzene system.

Step	Benzene I (4)	Benzene II (4)	Benzene III (4)	Benzene I (72)	Benzene II (72)	Benzene III (72)
Add Restraints	-1.700 ± 0.013	-1.729 ± 0.018	-1.709 ± 0.017	-1.900 ± 0.006	-1.975 ± 0.007	-2.074 ± 0.010
Remove Interactions	-9.377 ± 0.001	-9.159 ± 0.001	-9.061 ± 0.001	-9.316 ± 0.000	-9.050 ± 0.000	-8.970 ± 0.000
PV Correction	0.326 ± 0.000	0.322 ± 0.000	0.325 ± 0.000	0.000 ± 0.000	0.000 ± 0.000	0.000 ± 0.000
Total	-10.751 ± 0.014	-10.566 ± 0.010	-10.445 ± 0.008	-10.751 ± 0.000	-10.566 ± 0.000	-10.445 ± 0.000

Table C.3: Polymorph free energy differences in OPLS-AA for each step in the PSCP path at 200K for both the 4 and 72 benzene super cell models

Step	Benzene I (4)	Benzene II (4)	Benzene III (4)	Benzene I (72)	Benzene II (72)	Benzene III (72)
Add Restraints	-1.874 ± 0.012	-1.928 ± 0.013	-1.893 ± 0.014	-2.077 ± 0.006	-2.205 ± 0.007	-2.281 ± 0.005
Remove Interactions	-9.007 ± 0.001	-8.698 ± 0.001	-8.636 ± 0.001	-8.949 ± 0.000	-8.594 ± 0.000	-8.514 ± 0.000
PV Correction	0.326 ± 0.000	0.322 ± 0.000	0.317 ± 0.000	0.000 ± 0.000	0.000 ± 0.000	0.000 ± 0.000
Total	-10.555 ± 0.007	-10.304 ± 0.022	-10.212 ± 0.008	-10.555 ± 0.007	-10.304 ± 0.022	-10.212 ± 0.008

Table C.4: Polymorph free energy differences in GROMOS54A7 for each step in the PSCP path at 200K for both the 4 and 72 benzene super cell models

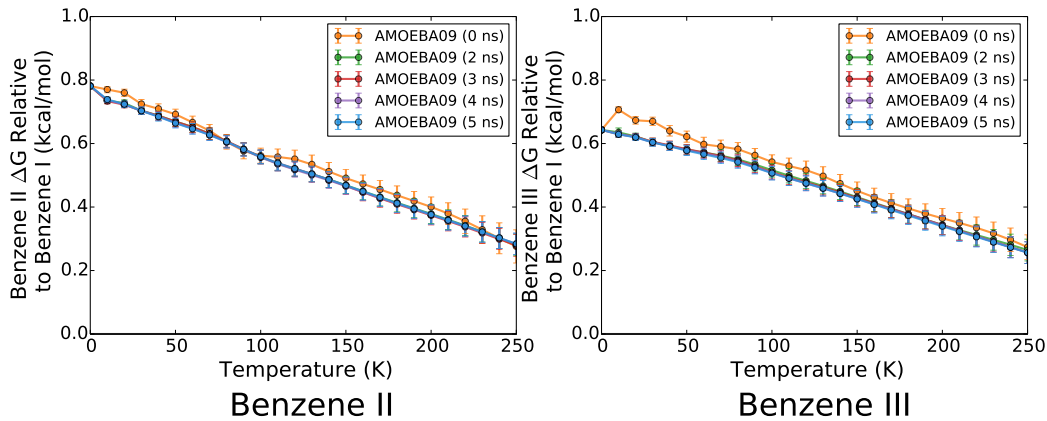


Figure C.6: The relative free energy as a function of temperature has a slight bias without any direct AMOEBA09 sampling, particularly for benzene III. This bias is remove with just a short 2ns of sampling directly in the AMOEBA09 potential.

Step	Benzene I (4)	Benzene II (4)	Benzene III (4)
Add Restraints	-1.797 ± 0.014	-1.867 ± 0.017	-1.834 ± 0.016
Remove Interactions	-9.424 ± 0.001	-9.171 ± 0.001	-9.084 ± 0.001
PV Correction	0.328 ± 0.000	0.324 ± 0.000	0.324 ± 0.000
Total	-10.893 ± 0.010	-10.714 ± 0.009	-10.594 ± 0.009

Table C.5: Polymorph free energy differences in Designed-GROMOS for each step in the PSCP path at 200K

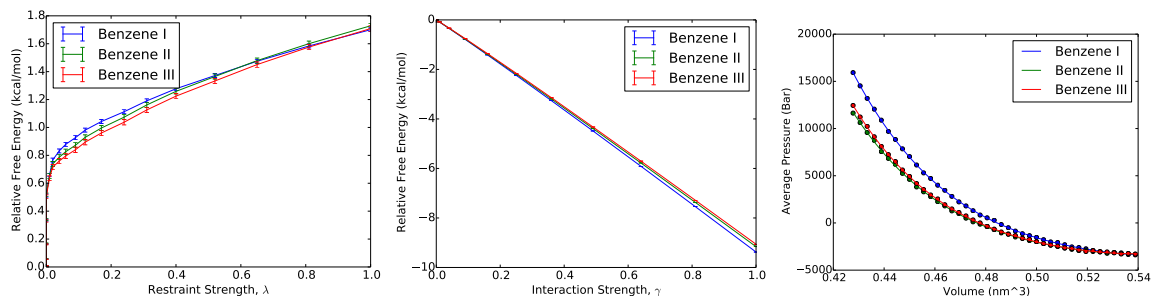


Figure C.7: Polymorph free energy profiles in OPLS-AA to A) restrain the benzene molecules to their crystalline position B) turn off intermolecular interactions and C) isotropically change the system volume.

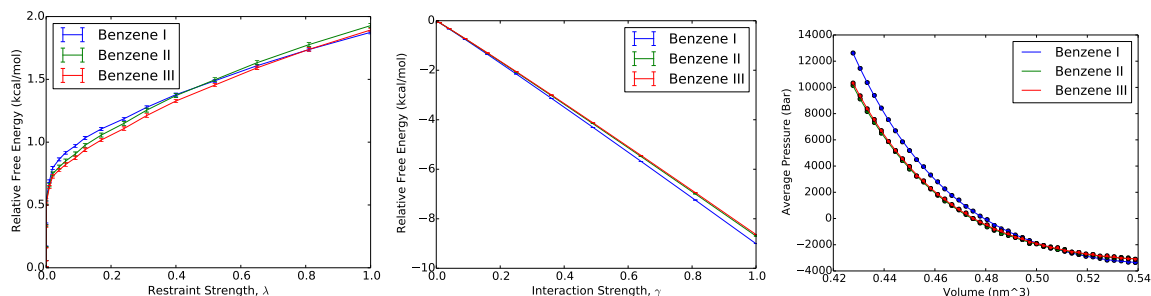


Figure C.8: Polymorph free energy profiles in GROMOS54A7 to A) restrain the benzene molecules to their crystalline position B) turn off intermolecular interactions and C) isotropically change the system volume.

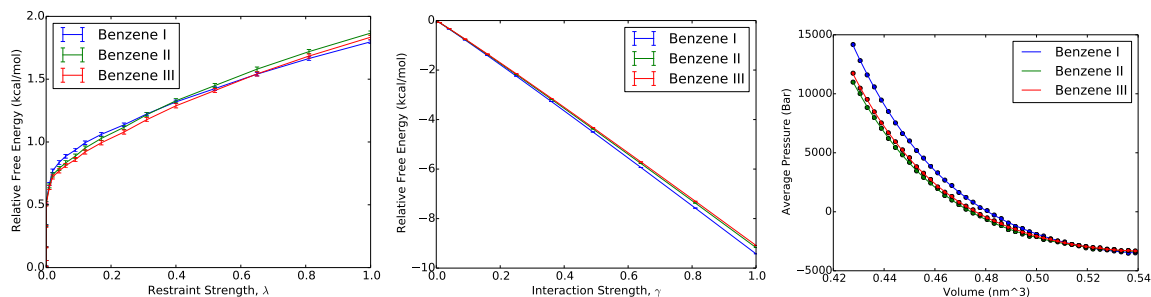


Figure C.9: Polymorph free energy profiles in Designed-GROMOS to A) restrain the benzene molecules to their crystalline position B) turn off intermolecular interactions and C) isotropically change the system volume.

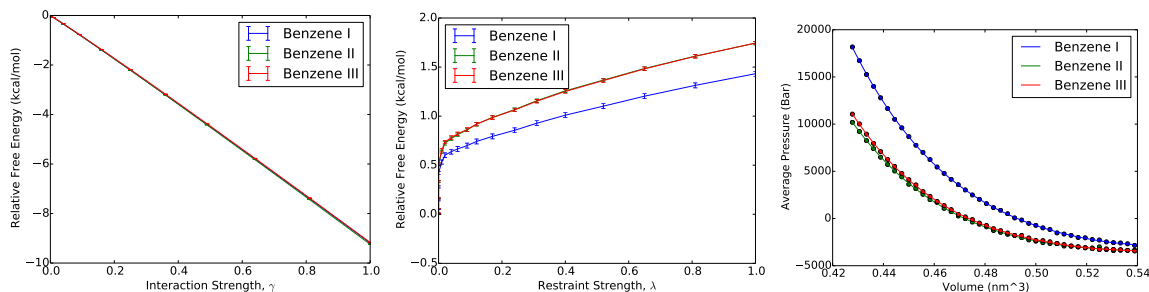


Figure C.10: Polymorph free energy profiles in Designed-AMOEBA to A) restrain the benzene molecules to their crystalline position B) turn off intermolecular interactions and C) isotropically change the system volume.

Step	Benzene I (4)	Benzene II (4)	Benzene III (4)
Add Restraints	-1.432 ± 0.027	-1.744 ± 0.016	-1.745 ± 0.018
Remove Interactions	-9.186 ± 0.001	-9.241 ± 0.001	-9.166 ± 0.001
PV Correction	0.320 ± 0.000	0.325 ± 0.000	0.322 ± 0.000
Total	-10.298 ± 0.010	-10.660 ± 0.009	-10.589 ± 0.009

Table C.6: Polymorph free energy differences in Designed-AMOEBA for each step in the PSCP path at 200K

Step	Benzene I (4)	Benzene II (4)	Benzene III (4)
Add Restraints	-0.471 ± 0.000	-0.472 ± 0.000	-0.504 ± 0.000
Remove Interactions	-9.931 ± 0.001	-9.704 ± 0.000	-9.611 ± 0.000
PV Correction	0.000 ± 0.000	0.000 ± 0.000	0.000 ± 0.000
Total	-10.402 ± 0.001	-10.176 ± 0.001	-10.115 ± 0.001

Table C.7: Polymorph free energy differences in OPLS-AA for each step in the PSCP path at 100K

Step	Benzene I (4)	Benzene II (4)	Benzene III (4)
Add Restraints	-0.544 ± 0.000	-0.583 ± 0.000	-0.621 ± 0.000
Remove Interactions	-9.601 ± 0.001	-9.284 ± 0.000	-9.193 ± 0.000
PV Correction	0.000 ± 0.000	0.000 ± 0.000	0.000 ± 0.000
Total	-10.145 ± 0.001	-9.867 ± 0.001	-9.849 ± 0.001

Table C.8: Polymorph free energy differences in GROMOS54A7 for each step in the PSCP path at 100K

Step	Benzene I (4)	Benzene II (4)	Benzene III (4)	Benzene I (72)	Benzene II (72)	Benzene III (72)
ΔG (without GROMOS54A7 sampling)	0.372 ± 0.010	0.425 ± 0.009	0.417 ± 0.009	0.367 ± 0.010	0.417 ± 0.006	0.416 ± 0.001
ΔG (with GROMOS54A7 sampling)	0.370 ± 0.010	0.423 ± 0.009	0.411 ± 0.009	0.366 ± 0.010	0.415 ± 0.006	0.413 ± 0.000

Table C.9: Free energy difference at 200K between GROMOS54A7 and Designed-GROMOS for each benzene polymorph in both the 4 and 72 benzene super cell models. The free energies are shown both with and without sampling in the target GROMOS54A7 potential

C.6 Polymorph Free Energy is Insensitive to Cutoff Distance

The cutoff distance used in the simulations in the main text is 0.7nm for both the Coulombic and Van der Waals interactions. A larger cutoff-distance provides a more accurate estimate of the total interaction energy, however a shorter cutoff-distance increases the speed of the simulation. The value of 0.7nm was chosen as an appropriate balance between the speed and accuracy of the simulations.

The sufficient accuracy of the potentials with a cutoff of 0.7nm was verified by examining the sensitivity of the system potential energy to the cutoff distance for the lattice energy minima in the OPLS-AA potential. The potential energy of the three polymorphs in the OPLS-AA potential as a function of cutoff distance is shown in figure C.11. The potential energy in each polymorph approaches an asymptotic value as the cutoff distance increases to 1.0 nm. At a cut-off of 0.7 nm the potential energies are all within $0.015 \text{ kcal}\cdot\text{mol}^{-1}$ of the value at 1.0nm. This is within the uncertainty of the polymorph free energy differences shown in the main text.

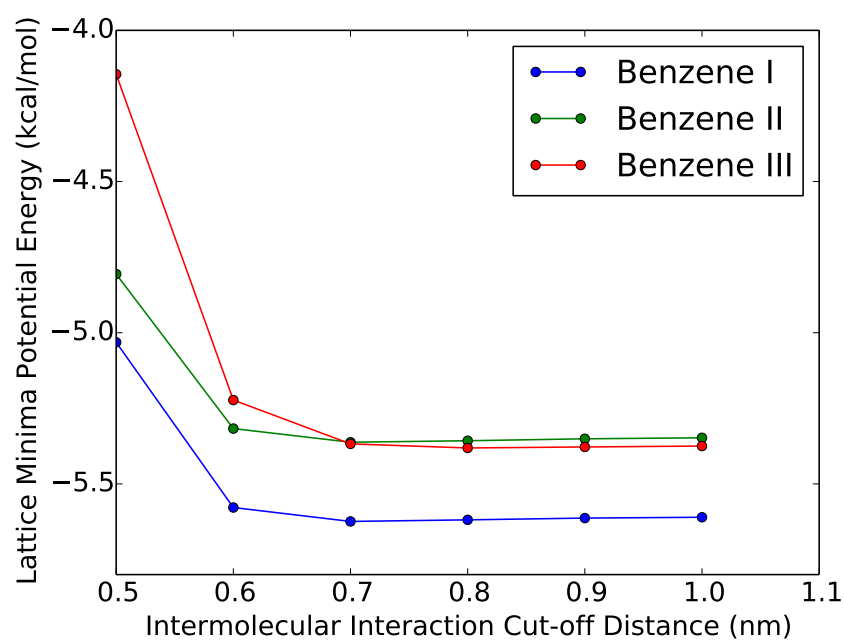


Figure C.11: The potential energies of three benzene lattice minima in the OPLS-AA potential are roughly insensitive to the interaction cut-off distance beyond 0.7 nm.

Appendix D

Appendix D

D.1 Offset Harmonic Oscillators

The methods NBB-indirect, three-stage MBAR(1) and three-stage MBAR(2) were used to estimate the free energy difference between a sampled harmonic oscillator and various unsampled harmonic oscillators through an intermediate sampled oscillator. The two cases examined in the main text are oscillators with equal offsets and different harmonic spring constants.

The bias in the free energy difference estimate for all Oscillators presented in figure D.1 are negligible compared with the uncertainty, as seen in Fig. D.2.

In the case presented above, the uncertainty in the free energy estimate with NBB-indirect is lower than the uncertainty in the free energy estimate with three-stage-MBAR(1). However, this result is only true in the limiting case that $N_1 = N_2$ where N_i is the number of samples drawn from oscillator i . In the above simulations, $N_1 = N_2 = 1000$. When $N_1 > N_2$ the uncertainty for NBB-indirect is larger than both three-stage MBAR (1) and three-stage MBAR (2) as shown below for $N_1 = 4000$ and $N_2 = 1000$.

In addition to perturbing the harmonic spring constants of the oscillators, the offset can also be changed. The sampling is biased for sampled oscillators with zero offset and unsampled oscillators with nonzero offset. An illustration of these oscillators is shown below:

This bias in the sampling leads to a bias in the resulting free energy difference estimate for all statistical estimators. The oscillators and the resulting uncertainties in the free energy difference estimates are shown below:

D.2 Mathematical differences between different reweighting schemes

The NBB estimator is written in eq 7 of König et al. [255] as the free energy between two states 0 and 1 as the free energy between two states 0 and 1 as:

$$\beta\Delta A_{0\rightarrow 1} = \ln \frac{\left\langle \frac{e^{\beta V_1^{bias}}}{1+e^{\beta(U_0-U_1+C)}} \right\rangle_{1,b} \left\langle e^{\beta V_0^{bias}} \right\rangle_{0,b}}{\left\langle \frac{e^{\beta V_0^{bias}}}{1+e^{-\beta(U_1-U_0+C)}} \right\rangle_{0,b} \left\langle e^{\beta V_1^{bias}} \right\rangle_{1,b}} + \beta C \quad (D.1)$$

Where $\beta\Delta A_{0\rightarrow 1} = \beta(A_1 - A_0)$, $V(x)^{bias}$ is the difference between the energy of the system including a biasing potential $U(x)^{bias}$ and the energy of the same system without the biasing potential $U(x)$, $V_i^{bias} = U_i^{biased} - U$. We also have $\beta C = \beta\Delta A_{0\rightarrow 1} + \ln \frac{N_1}{N_0}$, where N_1 and N_0 are samples collected from the states 1 and 0 simulated with biases, respectively, and we have inserted the definition of the function $f(x) = \frac{1}{1+e^x}$.

By plugging in the definition of C , gathering terms, and exponentiating, this simplifies to finding the $\Delta(A_1 - A_0)$ that satisfies the nonlinear equation:

$$N_0 \left\langle \frac{e^{\beta V_0^{bias}}}{1 + \frac{N_0}{N_1} e^{\beta(U_1-U_0-A_1+A_0)}} \right\rangle_{0,b} \left\langle e^{\beta V_1^{bias}} \right\rangle_{1,b} = N_1 \left\langle \frac{e^{\beta V_1^{bias}}}{1 + \frac{N_1}{N_0} e^{\beta(U_0-U_1-A_0+A_1)}} \right\rangle_{1,b} \left\langle e^{\beta V_0^{bias}} \right\rangle_{0,b} \quad (D.2)$$

One initial problem we note with this expression is that N_1 and N_0 cannot be the number of samples collected at the states 0 and 1 because these states do not have any samples collected at them; samples are instead collected at the biased version of states 0 and 1. The solution for previous applications of NBB was simply to set $N_0/N_1 = 1$. The equations remain true for any other choice of N_0/N_1 , though the convergence properties are different. This assumption of $N_0/N_1 = 1$ results in the nonlinear equation:

$$\left\langle \frac{e^{\beta V_0^{bias}}}{1 + e^{\beta(U_0-U_1-A_1+A_0)}} \right\rangle_{0,b} \left\langle e^{\beta V_1^{bias}} \right\rangle_{1,b} = \left\langle \frac{e^{\beta V_1^{bias}}}{1 + e^{\beta(U_1-U_0-A_0+A_1)}} \right\rangle_{1,b} \left\langle e^{\beta V_0^{bias}} \right\rangle_{0,b} \quad (D.3)$$

To better enable comparison to other methods, we relabel the four states of interest in the NBB estimator as:

- 0: State 0 in the above, from which no samples are collected
- 1: State 1 in the above, from which no samples are collected
- 2: The biased version of state 0 (0,b), from which N_2 samples are collected
- 3: The biased version of state 1 (1,b), from which N_3 samples are collected

Then $V_0^{bias}(x) = U_2(x) - U_0(x)$, and $V_1^{bias}(x) = U_3(x) - U_1(x)$. We are interested in the free energy difference from state 0 to state 1 (i.e. $\beta A_1 - \beta A_0$). We can simplify the math using the notation $\bar{u}_i(x) = e^{-\beta U_i(x)}$ and $\bar{f}_i(x) = e^{\beta A_i(x)}$, so that $\int \bar{f}_i \bar{u}_i(x) dx = 1$ is normalized.

N_i will represent the number of samples collected from each state of interest. The nonlinear equation then becomes:

$$\begin{aligned} \left\langle \frac{e^{\beta(U_2-U_0)}}{1 + e^{\beta(U_1-U_0-A_1+A_0)}} \right\rangle_2 \left\langle e^{\beta(U_3-U_1)} \right\rangle_3 &= \left\langle \frac{e^{\beta(U_3-U_1)}}{1 + e^{\beta(U_0-U_1-A_0+A_1)}} \right\rangle_3 \left\langle e^{\beta(U_2-U_0)} \right\rangle_2 \\ \left\langle \frac{\bar{u}_0}{1 + \frac{\bar{f}_0 \bar{u}_0}{\bar{f}_1 \bar{u}_1}} \right\rangle_2 \left\langle \frac{\bar{u}_1}{\bar{u}_3} \right\rangle_3 &= \left\langle \frac{\bar{u}_1}{1 + \frac{\bar{f}_1 \bar{u}_1}{\bar{f}_0 \bar{u}_0}} \right\rangle_3 \left\langle \frac{\bar{u}_0}{\bar{u}_2} \right\rangle_2 \\ \left\langle \frac{\bar{u}_0 \bar{u}_1 \bar{f}_1}{\bar{f}_0 \bar{u}_0 + \bar{f}_1 \bar{u}_1} \right\rangle_2 \left\langle \frac{\bar{u}_0}{\bar{u}_2} \right\rangle_2^{-1} &= \left\langle \frac{\bar{u}_1 \bar{u}_0 \bar{f}_0}{\bar{f}_0 \bar{u}_0 + \bar{f}_1 \bar{u}_1} \right\rangle_3 \left\langle \frac{\bar{u}_1}{\bar{u}_3} \right\rangle_3^{-1} \end{aligned}$$

We note that $\left\langle \frac{\bar{u}_1}{\bar{u}_3} \right\rangle_3$ and $\left\langle \frac{\bar{u}_0}{\bar{u}_2} \right\rangle_2$ are simply the Zwanzig estimators for the (exponential of) the free energy difference from states 0 to 2 and states 1 to 3. We write the ratio $\frac{k_0}{k_1} = \left\langle \frac{\bar{u}_1}{\bar{u}_2} \right\rangle_3 / \left\langle \frac{\bar{u}_0}{\bar{u}_2} \right\rangle_2$, obtaining:

$$\left\langle \frac{k_0 \bar{u}_1 \bar{f}_1 \frac{\bar{u}_0}{\bar{u}_2}}{\bar{f}_0 \bar{u}_0 + \bar{f}_1 \bar{u}_1} \right\rangle_2 = \left\langle \frac{k_1 \bar{u}_0 \bar{f}_0 \frac{\bar{u}_1}{\bar{u}_3}}{\bar{f}_0 \bar{u}_0 + \bar{f}_1 \bar{u}_1} \right\rangle_3 \quad (\text{D.4})$$

We will use the solution of this for f_1/f_0 as our operational definition of NBB for the purpose of comparisons.

The MBAR estimator between two unsampled states would be:

$$\frac{\bar{f}_1}{\bar{f}_0} = \frac{\left\langle \frac{\bar{u}_0}{f_2 \bar{u}_2 + f_3 \bar{u}_3} \right\rangle_{2+3}}{\left\langle \frac{\bar{u}_1}{f_2 \bar{u}_2 + f_3 \bar{u}_3} \right\rangle_{2+3}} \quad (\text{D.5})$$

Where f_2/f_3 can be solved for using MBAR:

$$1 = \left\langle \frac{\bar{f}_2 \bar{u}_2}{\bar{f}_2 \bar{u}_2 + \bar{f}_3 \bar{u}_3} \right\rangle_{2+3} \quad (\text{D.6})$$

$$\frac{\bar{f}_3}{\bar{f}_2} = \left\langle \frac{\bar{u}_2}{\frac{\bar{f}_2}{\bar{f}_3} \bar{u}_2 + \bar{u}_3} \right\rangle_{2+3} \quad (\text{D.7})$$

This can be shown to be mathematically equivalent to the Bennett Acceptance ratio for two states alone:

$$\left\langle \frac{\bar{f}_3 \bar{u}_3}{\bar{f}_2 \bar{u}_2 + \bar{f}_3 \bar{u}_3} \right\rangle_2 = \left\langle \frac{\bar{f}_2 \bar{u}_2}{\bar{f}_2 \bar{u}_2 + \bar{f}_3 \bar{u}_3} \right\rangle_3 \quad (\text{D.8})$$

Clearly, the MBAR estimator and the NBB estimator cannot be entirely equivalent, since MBAR includes the configurations from state 3 in the state 2 ensemble averaging as well as configurations from state 2 in the state 3 ensemble averaging, neither of which are used in the NBB algorithm.

So we instead first compare NBB to a simpler scheme, which includes only the same energy evaluations that NBB does, namely evaluations of U_0 in state 2, and U_1 in state 3. This is equivalent to the three-stage MBAR with one state discussed in the main paper.

To calculate $a = \frac{\bar{f}_1}{\bar{f}_0}$ using 3 state MBAR, we use two applications of Zwanzig from $3 \rightarrow 1$ and $0 \rightarrow 2$ plus MBAR (equivalently BAR) from 2 to 3. Since $\frac{k_0}{k_1} = \left\langle \frac{\bar{u}_1}{\bar{u}_3} \right\rangle_3 / \left\langle \frac{\bar{u}_0}{\bar{u}_2} \right\rangle_2 = \frac{\bar{f}_0 \bar{f}_3}{\bar{f}_2 \bar{f}_1}$, with the caveat these are free energies computed from the Zwanzig relationship, not exact free energies. We have

$$a k_0 / k_1 = \frac{\bar{f}_1}{\bar{f}_0} \frac{\bar{f}_0 \bar{f}_3}{\bar{f}_2 \bar{f}_1} = \frac{\bar{f}_3}{\bar{f}_2}$$

So the equation a (or $\frac{\bar{f}_1}{f_0}$) must satisfy in three-stage MBAR (with one state) is:

$$\begin{aligned} \left\langle \frac{ak_0\bar{u}_3}{k_1\bar{u}_2 + ak_0\bar{u}_3} \right\rangle_2 &= \left\langle \frac{k_1\bar{u}_2}{k_1\bar{u}_2 + ak_0\bar{u}_3} \right\rangle_3 \\ \left\langle \frac{k_0\bar{f}_1\bar{u}_3}{k_1\bar{f}_0\bar{u}_2 + k_0\bar{f}_1\bar{u}_3} \right\rangle_2 &= \left\langle \frac{k_1\bar{f}_0\bar{u}_2}{k_1\bar{f}_0\bar{u}_2 + k_0\bar{f}_1\bar{u}_3} \right\rangle_3 \end{aligned} \quad (\text{D.9})$$

while NBB is:

$$\left\langle \frac{k_0\bar{u}_1\bar{f}_1\frac{\bar{u}_0}{\bar{u}_2}}{\bar{f}_0\bar{u}_0 + \bar{f}_1\bar{u}_1} \right\rangle_2 = \left\langle \frac{k_1\bar{u}_0\bar{f}_0\frac{\bar{u}_1}{\bar{u}_3}}{\bar{f}_0\bar{u}_0 + \bar{f}_1\bar{u}_1} \right\rangle_3 \quad (\text{D.10})$$

which we rewrite as:

$$\left\langle \frac{k_0\frac{\bar{u}_0}{\bar{u}_2}}{1 + \frac{\bar{f}_0\bar{u}_0}{\bar{f}_1\bar{u}_1}} \right\rangle_2 = \left\langle \frac{k_1\frac{\bar{u}_1}{\bar{u}_3}}{1 + \frac{\bar{f}_1\bar{u}_1}{\bar{f}_0\bar{u}_0}} \right\rangle_3 \quad (\text{D.11})$$

And we rewrite three-stage MBAR as:

$$\left\langle \frac{1}{1 + \frac{k_1\bar{f}_0\bar{u}_2}{k_0\bar{f}_1\bar{u}_3}} \right\rangle_2 = \left\langle \frac{1}{1 + \frac{k_0\bar{f}_1\bar{u}_3}{k_1\bar{f}_0\bar{u}_2}} \right\rangle_3 \quad (\text{D.12})$$

Clearly if $k_0\frac{\bar{u}_0}{\bar{u}_2} = k_1\frac{\bar{u}_1}{\bar{u}_3} = 1$ then NBB and three-stage MBAR will be the same, though this is essentially the limiting case of just two state BAR.

We can rewrite the estimator for NBB as:

$$\left\langle \frac{k_0\bar{f}_0\bar{u}_2}{k_0\bar{f}_0\bar{u}_2 + k_1\bar{f}_1\bar{u}_3} \frac{k_0\bar{f}_0\bar{u}_2 + k_1\bar{f}_1\bar{u}_3}{\bar{f}_0\bar{u}_0 + \bar{f}_1\bar{u}_1} \right\rangle_2 = \left\langle \frac{k_1\bar{f}_1\bar{u}_3}{k_0\bar{f}_0\bar{u}_2 + k_1\bar{f}_1\bar{u}_3} \frac{k_0\bar{f}_0\bar{u}_2 + k_1\bar{f}_1\bar{u}_3}{\bar{f}_0\bar{u}_0 + \bar{f}_1\bar{u}_1} \right\rangle_3 \quad (\text{D.13})$$

If the term $\frac{k_0\bar{f}_0\bar{u}_2 + k_1\bar{f}_1\bar{u}_3}{\bar{f}_0\bar{u}_0 + \bar{f}_1\bar{u}_1} = 1$ for all configurations, then the methods would be equal, though clearly that is not always the case.

Expanding this expression for NBB out again, we obtain:

$$\begin{aligned}
 \frac{k_0 \bar{f}_0 \bar{u}_2 + k_1 \bar{f}_1 \bar{u}_3}{\bar{f}_0 \bar{u}_0 + \bar{f}_1 \bar{u}_1} &= \frac{\frac{\bar{f}_2}{\bar{f}_0} \bar{f}_0 \bar{u}_2 + \frac{\bar{f}_3}{\bar{f}_1} \bar{f}_1 \bar{u}_3}{\bar{f}_0 \bar{u}_0 + \bar{f}_1 \bar{u}_1} \\
 &= \frac{\frac{\bar{f}_2 \bar{u}_2}{\bar{f}_0 \bar{u}_0} \bar{f}_0 \bar{u}_0 + \frac{\bar{f}_3 \bar{u}_3}{\bar{f}_1 \bar{u}_1} \bar{f}_1 \bar{u}_1}{\bar{f}_0 \bar{u}_0 + \bar{f}_1 \bar{u}_1} \\
 &= \frac{\frac{\langle \frac{\bar{u}_0}{\bar{u}_2} \rangle_2}{\frac{\bar{u}_0}{\bar{u}_2}} \bar{f}_0 \bar{u}_0 + \frac{\langle \frac{\bar{u}_1}{\bar{u}_3} \rangle_3}{\frac{\bar{u}_1}{\bar{u}_3}} \bar{f}_1 \bar{u}_1}{\bar{f}_0 \bar{u}_0 + \bar{f}_1 \bar{u}_1} \tag{D.14}
 \end{aligned}$$

Inspecting this, we can see that on average, it will be relatively close to one in most cases. If u_0/u_2 and u_1/u_3 are constant (in other words, if $U_0 - U_2$ and $U_1 - U_3$ are constant), then the term will always be 1, but if it deviates, the term will not be 1, and introduces an extra source of statistical variation into the problem. Unless this variation is anti-correlated with the variance in the quantities averaged in the determination of BAR or MBAR themselves, it will mean the solutions will have additional statistical error, above and beyond BAR or MBAR.

In the case of NBB-indirect, the only difference is that we set the energy function of state 1 to be same as the energy function of state 3. We start from the previous final result, replacing state 1 with state 3, such that $u_1 = u_3$ and $\bar{f}_1 = \bar{f}_3$. k_1/k_0 becomes simply $k_0 = \left\langle \frac{\bar{u}_0}{\bar{u}_2} \right\rangle_2$, and $k_1 = 1$.

$$\left\langle \frac{k_0 \bar{f}_0 \bar{u}_2}{k_0 \bar{f}_0 \bar{u}_2 + \bar{f}_3 \bar{u}_3} \frac{k_0 \bar{f}_0 \bar{u}_2 + \bar{f}_3 \bar{u}_3}{\bar{f}_0 \bar{u}_0 + \bar{f}_3 \bar{u}_3} \right\rangle_2 = \left\langle \frac{\bar{f}_3 \bar{u}_3}{k_0 \bar{f}_0 \bar{u}_2 + \bar{f}_3 \bar{u}_3} \frac{k_0 \bar{f}_0 \bar{u}_2 + \bar{f}_3 \bar{u}_3}{\bar{f}_0 \bar{u}_0 + \bar{f}_3 \bar{u}_3} \right\rangle_3 \tag{D.15}$$

Three-stage MBAR with one reweighted end state in this case becomes:

$$\left\langle \frac{k_0 \bar{f}_0 \bar{u}_2}{k_0 \bar{f}_0 \bar{u}_2 + \bar{f}_3 \bar{u}_3} \right\rangle_2 = \left\langle \frac{\bar{f}_3 \bar{u}_3}{k_0 \bar{f}_0 \bar{u}_2 + \bar{f}_3 \bar{u}_3} \right\rangle_3 \tag{D.16}$$

This would be equivalent if the factor equals 1, which in this case is:

$$\frac{\frac{\langle \frac{\bar{u}_0}{\bar{u}_2} \rangle_2}{\frac{\bar{u}_0}{\bar{u}_2}} \bar{f}_0 \bar{u}_0 + \bar{f}_3 \bar{u}_3}{\bar{f}_0 \bar{u}_0 + \bar{f}_3 \bar{u}_3} \tag{D.17}$$

With the exact same information, one could also compute MBAR from two states to a third unsampled state, using exactly the same energies, in which case it would satisfy:

$$\bar{f}_0 = \left\langle \frac{\bar{u}_0}{\bar{f}_2 \bar{u}_2 + \bar{f}_3 \bar{u}_3} \right\rangle_{2+3} \quad (\text{D.18})$$

$$\bar{f}_2 = \left\langle \frac{\bar{u}_2}{\bar{f}_2 \bar{u}_2 + \bar{f}_3 \bar{u}_3} \right\rangle_{2+3} \quad (\text{D.19})$$

or eliminating the redundant variable:

$$\bar{f}_2 / \bar{f}_3 = \left\langle \frac{\bar{u}_2}{\bar{f}_2 / \bar{f}_3 \bar{u}_2 + \bar{u}_3} \right\rangle_{2+3} \quad (\text{D.20})$$

$$\bar{f}_0 / \bar{f}_3 = \left\langle \frac{\bar{u}_0}{\bar{f}_2 / \bar{f}_3 \bar{u}_2 + \bar{u}_3} \right\rangle_{2+3} \quad (\text{D.21})$$

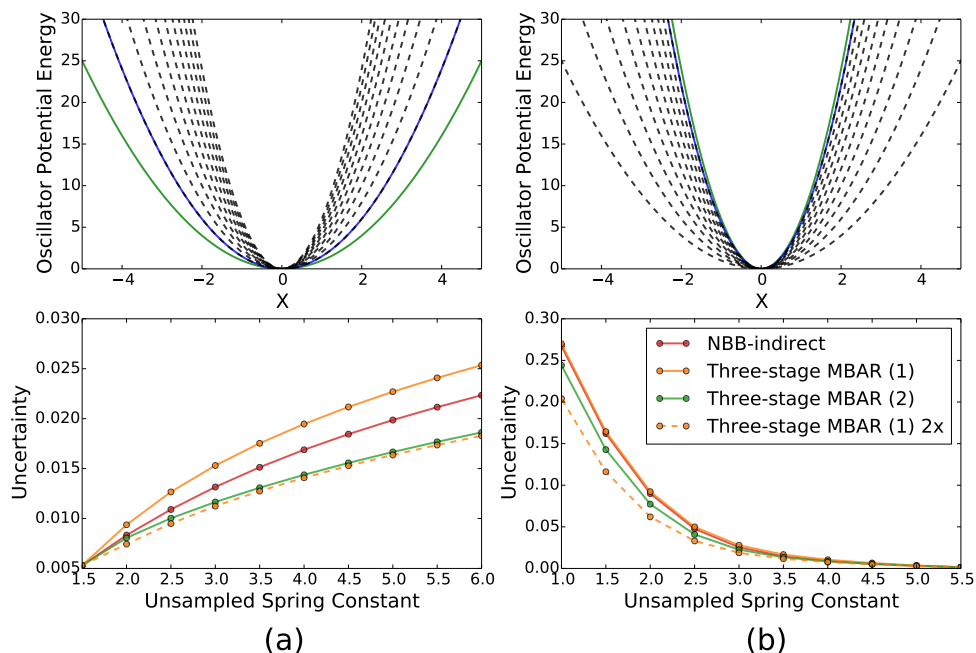


Figure D.1: The uncertainties in the free energy difference between the sampled and unsampled oscillators increase for all methods as the spring constant of the sampled and unsampled oscillators become farther apart. The uncertainties are lower with three-stage MBAR (2) than with NBB-indirect for all values of the spring constant of the unsampled oscillator in each case. For an equal amount of energy re-evaluations, the method with the lowest variance involves placing all the re-evaluations in the sampled oscillator with the highest overlap and reweighting with the Zwanzig equation. The set of harmonic oscillators are shown above in which the offset is identical and the spring constant of the oscillators changes. The green and blue solid lines represent the first and second sampled oscillator, and the dashed lines indicate the set of unsampled oscillators. The sampled oscillators are reweighted to the unsampled oscillators in two cases (a) the two oscillators with the smallest spring constant are sampled and (b) the two oscillators with the largest spring constant are sampled.

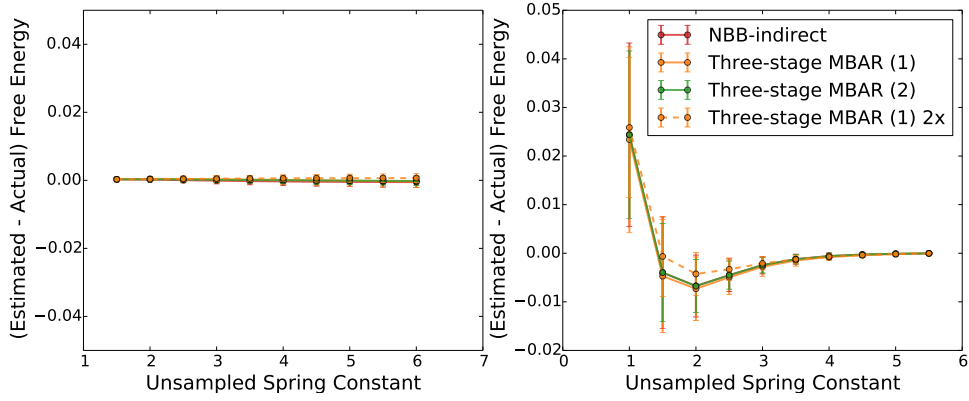


Figure D.2: The bias in the free energy estimate is small for all reweighting between oscillators of varying spring constant show in this work.

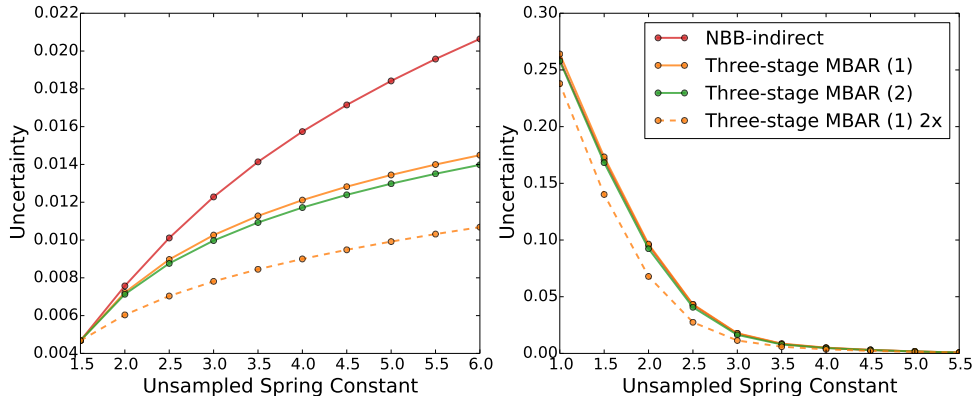


Figure D.3: When 4000 configurations are drawn from the first oscillator and 1000 configurations are drawn from the second, the variance in the free energy estimate is lower with three-stage MBAR(1) than with NBB when reweighting from wide to stiff oscillators. When reweighting from stiff to wide oscillators, the variance is essentially equal for all methods.

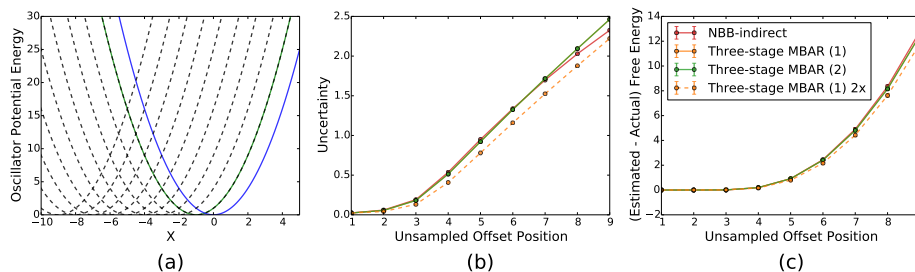


Figure D.4: The spring constant is identical for all harmonic oscillators, and the offset changes for the unsampled states leading to a biased sampling of the unsampled states configurations. The green and blue solid lines represent the first and second sampled oscillator, and the dashed lines indicate the set of unsampled oscillators.

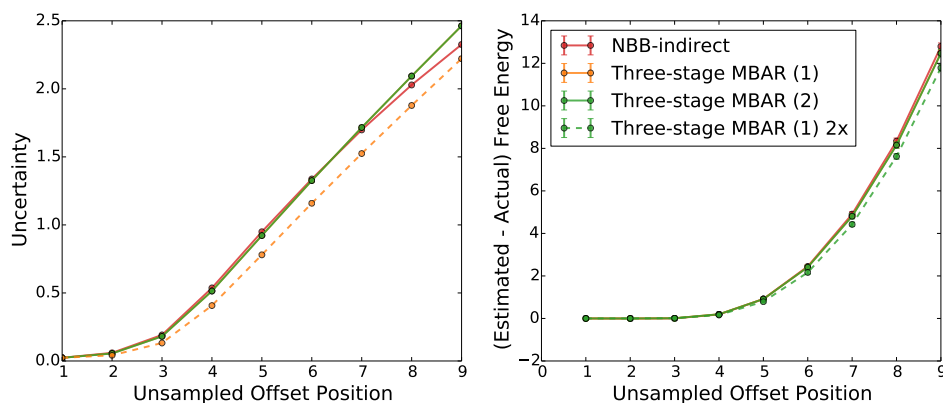


Figure D.5: The uncertainty in the free energy estimates for offset oscillators are relatively large for all methods. The variance in the estimate for unsampled offsets above 6 are lower for NBB-indirect than three-stage-MBAR (2), however the bias in the estimate for offsets above 6 are much larger than the uncertainties and all estimates with all methods are essentially meaningless.

Bibliography

- [1] Joel Bernstein. *Polymorphism in Molecular Crystals*. International Union of Crystallography, 14 edition, 2002.
- [2] Franco Cataldo. a Study on the Thermal Stability To 1000 C of Various Carbon Allotropes and Carbonaceous Matter Both Under Nitrogen. *Fullerenes, Nanotub. Carbon Nanostructures*, 10(4):293–311, 2002.
- [3] K. S. Novoselov, A. K. Geim, S. V. Morozov, D. Jiang, Y. Zhang, S. V. Dubonos, I. V. Grigorieva, and A. A. Firsov. Electric Field Effect in Atomically Thin Carbon Films. *Science*, 306(5696):666–669, 2004.
- [4] H. W. Kroto, J. R. Heath, S. C. O'Brien, R. F. Curl, and R. E. Smalley. C60 Buckminsterfullerene Smalley Nature 1985.pdf. *Nature*, 318:162–163, 1985.
- [5] R. Boehler. High-pressure experiments and the phase diagram of lower mantle and core materials. *Rev. Geophys.*, 38(2):221–245, 2000.
- [6] W. C. McCrone. *Polymorphism in Physics and Chemistry of the Organic Solid State*. Interscience, New York, 1965.
- [7] Jack D. Dunitz and Joel Bernstein. Disappearing Polymorphs. *Acc. Chem. Res.*, 28(4):193–200, 1995.
- [8] Angelo Gavezzotti. Are crystal structures predictable? *Acc. Chem. Res.*, 27(10):309–314, 1994.
- [9] G. R. Desiraju. Crystal Gazing: Structure Prediction and Polymorphism. *Science*, 278(5337):404–405, 1997.
- [10] Sally L Price. Computed crystal energy landscapes for understanding and predicting organic crystal structures and polymorphism. *Acc. Chem. Res.*, 42(1):117–126, jan 2009.
- [11] Graeme M. Day. Current approaches to predicting molecular organic crystal structures. *Crystallogr. Rev.*, 17(1):3–52, jan 2011.
- [12] A. G. Pandolfo and A. F. Hollenkamp. Carbon properties and their role in supercapacitors. *J. Power Sources*, 157(1):11–27, 2006.
- [13] A. A. Balandin. Thermal properties of graphene and nanostructured carbon materials. *Nat. Mater.*, 10(8):569, 2011.
- [14] J Bauer, S Spanton, R Henry, J Quick, W Dziki, W Porter, and J Morris. Ritonavir: an extraordinary example of conformational polymorphism. *Pharm. Res.*, 18(6):859–866, jun 2001.
- [15] Ivo B. Rietveld and Rene Ceolin. Rotigotine: Unexpected Polymorphism with Predictable Overall Monotropic Behavior. *J. Pharm. Sci.*, 104(12):4117–4122, 2015.

- [16] Peddy Vishweshwar, Jennifer A McMahon, Mark Oliveira, Matthew L Peterson, Michael J Zaworotko, V Uni, and East Fowler A V. The Predictably Elusive Form II of Aspirin. *J. Am. Chem. Soc.*, 127(48):16802–16803, 2005.
- [17] Philippe Espeau, Rene Ceolin, Josep Lluís Tamarit, Marc Antoine Perrin, Jean Pierre Gauchi, and Franck Leveiller. Polymorphism of paracetamol: Relative stabilities of the monoclinic and orthorhombic phases inferred from topological pressure-temperature and temperature-volume phase diagrams. *J. Pharm. Sci.*, 94(3):524–539, 2005.
- [18] M Joan Taylor, Sangeeta Tanna, and Tarsem Sahota. Conformational Polymorphism in Aripiprazole. *J. Pharm. Sci.*, 99(10):4215–4227, 2010.
- [19] Stephen Toon and William F Trager. Thermodynamic studies of Tolbutamide Polymorphs. *J. Pharm. Sci.*, 73(445):1673–1675, 1984.
- [20] Francesca P. A. Fabbiani and Colin R. Pulham. High-pressure studies of pharmaceutical compounds and energetic materials. *Chem. Soc. Rev.*, 35(10):932–942, 2006.
- [21] The Royal Society, Royal Society, and Physical Sciences. Deformation and fracture of b-HMX. *Proc. R. Soc. A*, 383(1785):399–407, 1982.
- [22] Antoine E. D. M. Van Der Heijden and Richard H B Bouma. Crystallization and characterization of RDX, HMX, and CL-20. *Cryst. Growth Des.*, 4(5):999–1007, 2004.
- [23] J Bernstein. *Polymorphism in Molecular Crystals*. Oxford University Press, New York, 2008.
- [24] David I. A. Millar, Helen E. Maynard-Casely, David R. Allan, Adam S. Cumming, Alistair R. Lennie, Alexandra J. Mackay, Iain D. H. Oswald, Chiu C. Tang, and Colin R. Pulham. Crystal engineering of energetic materials: Co-crystals of CL-20. *CrystEngComm*, 14(10):3742, 2012.
- [25] Erich F. Paulus, Frank J. J. Leusen, and Martin U. Schmidt. Crystal structures of quinacridones. *CrystEngComm*, 9(2):131–143, 2007.
- [26] S. Haas, A. Stassen, G. Schuck, K. Pernstich, D. Gundlach, B. Batlogg, U. Berens, and H.-J. Kirner. High charge-carrier mobility and low trap density in a rubrene derivative. *Phys. Rev. B*, 76(11):115203, sep 2007.
- [27] Gaurav Giri, Eric Verploegen, Stefan C. B. Mannsfeld, Sule Atahan-Evrenk, Do Hwan Kim, Sang Yoon Lee, Hector A. Becerril, Alán Aspuru-Guzik, Michael F. Toney, and Zhenan Bao. Tuning charge transport in solution-sheared organic semiconductors using lattice strain. *Nature*, 480(7378):504–508, 2011.
- [28] Yongbo Yuan, Gaurav Giri, Alexander L Ayzner, Arjan P Zoombelt, Stefan C B Mannsfeld, Jihua Chen, Dennis Nordlund, Michael F Toney, Jinsong Huang, and Zhenan Bao. Ultra-high mobility transparent organic thin film transistors grown by an off-centre spin-coating method. *Nat. Commun.*, 5:3005, 2014.

- [29] Loah A. Stevens, Katelyn P. Goetz, Alexandr Fonari, Ying Shu, Rachel M. Williamson, Jean-Luc Brédas, Veaceslav Coropceanu, Oana D. Jurchescu, and Gavin E. Collis. Temperature-Mediated Polymorphism in Molecular Crystals: The Impact on Crystal Packing and Charge Transport. *Chem. Mater.*, 27(1):112–118, 2015.
- [30] A. Chithambararaj, N. Rajeswari Yogamalar, and A. Chandra Bose. Hydrothermally Synthesized h-MoO₃ and α -MoO₃ Nanocrystals: New Findings on Crystal-Structure-Dependent Charge Transport. *Cryst. Growth Des.*, 16(4):1984–1995, 2016.
- [31] Sarah L Price. From crystal structure prediction to polymorph prediction: interpreting the crystal energy landscape. *Phys. Chem. Chem. Phys.*, 10(15):1996–2009, apr 2008.
- [32] M. M. Thiery and J. M. Leger. High pressure solid phases of benzene. I. Raman and x-ray studies of C₆H₆ at 294 K up to 25 GPa. *J. Chem. Phys.*, 89(7):4255, 1988.
- [33] T. Beyer, G. M. Day, and S. L. Price. The prediction, morphology, and mechanical properties of the polymorphs of paracetamol. *J. Am. Chem. Soc.*, 123(13):5086–5094, 2001.
- [34] Manish Kumar Mishra, Upadrasta Ramamurty, and Gautam R Desiraju. Solid Solution Hardening of Molecular Crystals: Tautomeric Polymorphs of Omeprazole. *J. Am. Chem. Soc.*, 137(5):1794–1797, 2014.
- [35] Baiju P Krishnan and Kana M Sureshan. A Molecular Level Study of Metamorphosis and Strengthening of Gel via Spontaneous Polymorphic Transition. *ChemPhysChem*, 17(19):3062–3067, 2016.
- [36] Gamidi Rama Krishna, Ramesh Devarapalli, Garima Lal, and C. Malla Reddy. Mechanically Flexible Organic Crystals Achieved by Introducing Weak Interactions in Structure: Supramolecular Shape Synthons. *J. Am. Chem. Soc.*, 138(41):13561–13567, 2016.
- [37] V. Ramar, K. Saravanan, S. R. Gajjela, S. Hariharan, and P. Balaya. The effect of synthesis parameters on the lithium storage performance of LiMnPO₄/C. *J. Power Sources*, 306:552–558, 2016.
- [38] Hyunjoong Chung and Ying Diao. Polymorphism as an Emerging Design Strategy for High Performance Organic Electronics. *J. Mater. Chem. C*, 2016.
- [39] Chun Ju Hou, Jorge Botana, Xu Zhang, Xianlong Wang, and Mao sheng Miao. Pressure-induced structural and valence transition in AgO. *Phys. Chem. Chem. Phys.*, 2016.
- [40] Konstantin I. Hadjiivanov and Dimitar G. Klissurski. Surface chemistry of titania (anatase) and titania-supported catalysts. *Chem. Soc. Rev.*, 25(1):61, 1996.
- [41] Barr Halevi, Eric J. Peterson, Aaron Roy, Andrew Delariva, Ese Jeroro, Feng Gao, Yong Wang, John M. Vohs, Boris Kiefer, Edward Kunkes, Michael Hävecker, Malte Behrens, Robert Schlögl, and Abhaya K. Datye. Catalytic reactivity of face centered cubic PdZn for the steam reforming of methanol. *J. Catal.*, 291:44–54, 2012.

- [42] Aron Walsh and David O. Scanlon. Polymorphism of indium oxide: Materials physics of orthorhombic In₂O₃. *Phys. Rev. B*, 88(16):161201, oct 2013.
- [43] Lisa N Hutfluss and Pavle V Radovanovic. Controlling the Mechanism of Phase Transformation of Colloidal In₂O₃ Nanocrystals. *J. Am. Chem. Soc.*, 137(3):1101–1108, jan 2015.
- [44] P Sulovsky and T Stanek. The influence of the alite polymorphism on the strength of the Portland cement. *Cem. Concr. Res.*, 32(7):1169–1175, 2002.
- [45] Mino R. Caira, Khoulood A. Alkhamis, and Rana M. Obaidat. Preparation and Crystal Characterization of a Polymorph, a Monohydrate, and an Ethyl Acetate Solvate of the Antifungal Fluconazole. *J. Pharm. Sci.*, 93(3):601–611, 2004.
- [46] Lian Yu, Jun Huang, and Karen J. Jones. Measuring free-energy difference between crystal polymorphs through eutectic melting. *J. Phys. Chem. B*, 109(42):19915–19922, 2005.
- [47] G. C. Viscomi, M. Campana, M. Barbanti, F. Grepioni, M. Polito, D. Confortini, G. Rosini, P. Righi, V. Cannata, and D. Braga. Crystal forms of rifaximin and their effect on pharmaceutical properties. *CrystEngComm*, 10(8):1074, 2008.
- [48] Donald A. McAfee, Jonathan Hadgraft, and Majella E. Lane. Rotigotine: The first new chemical entity for transdermal drug delivery. *Eur. J. Pharm. Biopharm.*, 88(3):586–593, 2014.
- [49] Dominic Cheuk, Dikshitkumar Khamar, Patrick McArdle, and Åke C. Rasmuson. Solid Forms, Crystal Habits, and Solubility of Danthron. *J. Chem. Eng. Data*, 60(7):2110–2118, 2015.
- [50] Sanjay R Chemburkar, John Bauer, Kris Deming, Harry Spiwek, Ketan Patel, John Morris, Rodger Henry, Stephen Spanton, Walter Dziki, William Porter, John Quick, Phil Bauer, John Donaubaue, B A Narayanan, Mauro Soldani, Dave Riley, and Kathryn McFarland. Dealing with the Impact of Ritonavir Polymorphs on the Late Stages of Bulk Drug Process Development. *Org. Process Res. Dev.*, 4(5):413–417, 2000.
- [51] C Rustichelli, G Gamberini, V Ferioli, M.C Gamberini, R Ficarra, and S Tommasini. Solid-state study of polymorphic drugs: carbamazepine. *J. Pharm. Biomed. Anal.*, 23(1):41–54, aug 2000.
- [52] Andrew D Bond, Roland Boese, and Gautam R Desiraju. On the polymorphism of aspirin: crystalline aspirin as intergrowths of two “polymorphic” domains. *Angew. Chem. Int. Ed. Engl.*, 46(4):618–622, jan 2007.
- [53] Ashley T Hulme, Sarah L Price, and Derek A Tocher. A New Polymorph of 5-Fluorouracil Found Following Computational Crystal Structure Predictions. *J. Am. Chem. Soc.*, 127(4):1116–1117, 2005.
- [54] Said Hamad, Changman Moon, C Richard A Catlow, Ashley T Hulme, and Sarah L Price. Kinetic insights into the role of the solvent in the polymorphism of 5-fluorouracil from molecular dynamics simulations. *J. Phys. Chem. B*, 110(7):3323–3329, feb 2006.

- [55] O I Corrigan, K Sabra, and E M Holohan. Physicochemical properties of spray-dried drugs: phenobarbitone and hydroflumethiazide. *Drug Dev. Ind. Pharm.*, 9(1&2):1–20, 1983.
- [56] Dharmendra Singhal and William Curatolo. Drug polymorphism and dosage form design: A practical perspective. *Adv. Drug Deliv. Rev.*, 56(3):335–347, 2004.
- [57] Dejan-Krešimir Bučar, Robert W. Lancaster, and Joel Bernstein. Disappearing Polymorphs Revisited. *Angew. Chemie Int. Ed.*, 54(24):6972–6993, 2015.
- [58] Nir Giladi, Babak Boroojerdi, Amos D. Korczyn, David J. Burn, Carl E. Clarke, and Anthony H V Schapira. Rotigotine transdermal patch in early Parkinson’s disease: A randomized, double-blind, controlled study versus placebo and ropinirole. *Mov. Disord.*, 22(16):2398–2404, 2007.
- [59] Armin Breitenbach. Transdermal delivery system for the administration of rotigotine, 2003.
- [60] William A Rakoczy and Deanne M Mazzochi. The case of the disappearing polymorph: Inherent anticipation’ and the impact of SmithKline Beecham Corp. v Apotex Corp. (Paxil®) on patent validity and infringement by inevitable conversion. *J. Generic Med.*, 3(2):131–139, 2006.
- [61] Walter Cabri, Paolo Ghetti, Giovanni Pozzi, and Marco Alpegiani. Polymorphisms and patent, market, and legal battles: Cefdinir case study. *Org. Process Res. Dev.*, 11(1):64–72, 2007.
- [62] Darren J. Lipomi, Benjamin C K Tee, Michael Vosgueritchian, and Zhenan Bao. Stretchable organic solar cells. *Adv. Mater.*, 23(15):1771–1775, 2011.
- [63] Kazunori Kuribara, He Wang, Naoya Uchiyama, Kenjiro Fukuda, Tomoyuki Yokota, Ute Zschieschang, Chernojaye, Daniel Fischer, Hagen Klauk, Tatsuya Yamamoto, Kazuo Takimiya, Masaaki Ikeda, Hirokazu Kuwabara, Tsuyoshi Sekitani, Yueh-Lin Loo, and Takao Someya. Organic transistors with high thermal stability for medical applications. *Nat. Commun.*, 3:723, 2012.
- [64] George Malliaras and Richard Friend. An Organic Electronics Primer. *Phys. Today*, 58(5):53–58, 2005.
- [65] Aldo Brillante, Ivano Bilotti, Raffaele Guido Della Valle, Elisabetta Venuti, and Alberto Girlando. Probing polymorphs of organic semiconductors by lattice phonon Raman microscopy. *CrystEngComm*, 10(8):937, 2008.
- [66] James Patterson, Andrew Bary, and Thomas Rades. Physical stability and solubility of the thermotropic mesophase of fenoprofen calcium as pure drug and in a tablet formulation. *Int. J. Pharm.*, 247(1-2):147–157, 2002.
- [67] Evgeniy Shalaev, Ke Wu, Sheri Shamblin, Joseph F Krzyzaniak, and Marc Descamps. Crystalline mesophases : Structure , mobility , and pharmaceutical properties . *Adv. Drug Deliv. Rev.*, 100:194–211, 2016.
- [68] Baiju P Krishnan and Kana M Sureshan. A spontaneous single-crystal-to-single-crystal polymorphic transition involving major packing changes. *J. Am. Chem. Soc.*, 137(4):1692–1696, jan 2015.

- [69] Chelsea M. Hess, Angela R. Rudolph, and Philip J. Reid. Imaging the Effects of Annealing on the Polymorphic Phases of Poly(vinylidene fluoride). *J. Phys. Chem. B*, 119(10):4127–4132, 2015.
- [70] Horng Long Cheng and Jr Wei Lin. Controlling polymorphic transformations of pentacene crystal through solvent treatments: An experimental and theoretical study. *Cryst. Growth Des.*, 10(10):4501–4508, 2010.
- [71] Paul D. Schmitt, Emma L. DeWalt, Ximeng Y. Dow, and Garth J. Simpson. Rapid Discrimination of Polymorphic Crystal Forms by Nonlinear Optical Stokes Ellipsometric Microscopy. *Anal. Chem.*, 88(11):5760–5768, 2016.
- [72] Lucia Carlucci and Angelo Gavezzotti. Molecular recognition and crystal energy landscapes: An X-ray and computational study of caffeine and other methylxanthines. *Chem. - A Eur. J.*, 11(1):271–279, 2005.
- [73] WZ Ostwald. Studies on formation and transformation of solid materials. *Zeitschrift für Phys. Chemie*, 22:289–330, 1897.
- [74] Jie Chen and Bernhardt L Trout. Computational Study of Solvent Effects on the Molecular Self-Assembly of Tetrolic Acid in Solution and Implications for the Polymorph Formed from Crystallization. *J. Phys. Chem. B*, 112(26):7794–7802, 2008.
- [75] Sherry L. Morissette, Örn Almarsson, Matthew L. Peterson, Julius F. Remenar, Michael J. Read, Anthony V. Lemmo, Steve Ellis, Michael J. Cima, and Colin R. Gardner. High-throughput crystallization: Polymorphs, salts, co-crystals and solvates of pharmaceutical solids. *Adv. Drug Deliv. Rev.*, 56(3):275–300, 2004.
- [76] Francois Cansell, Denise Fabre, and Jean-Pierre Petit. Phase transitions and chemical transformations of benzene up to 550C and 30 GPa. *J. Chem. Phys.*, 99(10):7300, 1993.
- [77] E. V. Boldyreva, T. P. Shakhshneider, H. Ahsbahs, H. Sowa, and H. Uchtmann. Effect of high pressure on the polymorphs of paracetamol. *J. Therm. Anal. Calorim.*, 68(2):437–452, 2002.
- [78] Elena V. Boldyreva. High-pressure diffraction studies of molecular organic solids. A personal view. *Acta Crystallogr. Sect. A Found. Crystallogr.*, 64(1):218–231, 2008.
- [79] Roman Gajda and Andrzej Katrusiak. Pressure-Promoted CH...O Hydrogen Bonds in Formamide Aggregates. *Cryst. Growth Des.*, 11(11):4768–4774, 2011.
- [80] Damian Paliwoda, Kamil F. Dziubek, and Andrzej Katrusiak. Imidazole hidden polar phase. *Cryst. Growth Des.*, 12(9):4302–4305, 2012.
- [81] Yury V. Seryotkin, Tatiana N. Drebuschak, and Elena V. Boldyreva. A high-pressure polymorph of chlorpropamide formed on hydrostatic compression of the α -form in saturated ethanol solution. *Acta Crystallogr. Sect. B Struct. Sci. Cryst. Eng. Mater.*, 69(1):77–85, 2013.
- [82] Makoto Otsuka, Mika Onoe, and Yoshihisa Matsuda. Physicochemical Stability of Phenobarbital Polymorphs at Various Levels of Humidity and Temperature, 1993.

- [83] Olímpia Maria Martins Santos, Maria Esther Dias Reis, Jennifer Tavares Jacon, Mônica Esselin De Sousa Lino, Juliana Savioli Simões, and Antonio Carlos Doriguetto. Polymorphism: An evaluation of the potential risk to the quality of drug products from the Farmácia Popular Rede Própria. *Brazilian J. Pharm. Sci.*, 50(1):1–24, 2014.
- [84] G Patrick Stahly. Diversity in Single- and Multiple-Component Crystals. The Search for and Prevalence of Polymorphs and Cocrystals. *Cryst. Growth Des.*, 7(6):1007–1026, 2007.
- [85] Andrew D. Bond. Polymorphism in molecular crystals. *Curr. Opin. Solid State Mater. Sci.*, 13(3-4):91–97, jun 2009.
- [86] Peng Yi and Gregory C Rutledge. Molecular origins of homogeneous crystal nucleation. *Annu. Rev. Chem. Biomol. Eng.*, 3:157–82, jan 2012.
- [87] Amanuel Gebrekrestos, Maya Sharma, Giridhar Madras, and Suryasarathi Bose. New Physical Insights into Shear History Dependent Polymorphism in Poly(vinylidene fluoride). *Cryst. Growth Des.*, 16(5):2937–2944, 2016.
- [88] Aurora J Cruz-Cabeza, Susan M Reutzel-Edens, and Joel Bernstein. Facts and fictions about polymorphism. *Chem. Soc. Rev.*, 44(23):8619–8635, 2015.
- [89] Jean-Baptiste Arlin, Louise S. Price, Sarah L. Price, and Alastair J. Florence. A strategy for producing predicted polymorphs: catemeric carbamazepine form V. *Chem. Commun.*, 47(25):7074, 2011.
- [90] Yong Chen, Long Li, Jia Yao, Yu-Yu Ma, Jia-Mei Chen, and Tong-Bu Lu. Improving the Solubility and Bioavailability of Apixaban via ApixabanOxalic Acid Cocrystal. *Cryst. Growth Des.*, 16(5):2923–2930, 2016.
- [91] Sean P. Delaney and Timothy M. Korter. Terahertz Spectroscopy and Computational Investigation of the Flufenamic Acid/Nicotinamide Cocrystal. *J. Phys. Chem. A*, 119(13):3269–3276, 2015.
- [92] Tomislav Stolar, Stipe Lukin, Josip Požar, Mirta Rubčić, Graeme M. Day, Ivana Biljan, Dubravka Šišak Jung, Gordan Horvat, Krunoslav Užarević, Ernest Meštrović, and Ivan Halasz. Solid-State Chemistry and Polymorphism of the Nucleobase Adenine. *Cryst. Growth Des.*, 16(6):3262–3270, 2016.
- [93] Adam L. Grzesiak, Meidong Lang, Kibum Kim, and Adam J. Matzger. Comparison of the Four Anhydrous Polymorphs of Carbamazepine and the Crystal Structure of Form I. *J. Pharm. Sci.*, 92(11):2260–2271, 2003.
- [94] M. Sacchetti. Thermodynamic analysis of DSC data for acetaminophen polymorphs. *J. Therm. Anal. Calorim.*, 63(2):345–350, 2001.
- [95] Suryanarayan Cherukuvada, Ranjit Thakuria, and Ashwini Nangia. Pyrazinamide polymorphs: Relative stability and vibrational spectroscopy. *Cryst. Growth Des.*, 10(9):3931–3941, 2010.
- [96] Gen Hasegawa, Takao Komasa, Rui Bando, Yasuo Yoshihashi, Etsuo Yonemochi, Kotaro Fujii, Hidehiro Uekusa, and Katsuhide Terada. Reevaluation of solubility of tolbutamide and polymorphic transformation from Form I to unknown crystal form. *Int. J. Pharm.*, 369(1-2):12–18, 2009.

- [97] Chandra Vemavarapu, Matthew J Mollan, and Thomas E Needham. Crystal doping aided by rapid expansion of supercritical solutions. *AAPS PharmSciTech*, 3(4):17–31, 2002.
- [98] Masaki Yoshino, Kazuhiro Takahashi, Youhei Okuda, Takashi Yoshizawa, Nobuaki Fukushima, and Moto-suke Naoki. Contribution of Hydrogen Bonds to Equilibrium α/β Transition of Resorcinol. *J. Phys. Chem. A*, 103(15):2775–2783, 1999.
- [99] Attilio Cesaro and Giorgio Starec. Thermodynamic properties of caffeine crystal forms. *J. Phys. Chem.*, 84(11):1345–1346, 1980.
- [100] Antonio Torrisi, Charlotte K. Leech, Kenneth Shankland, William I. F. David, Richard M. Ibberson, Jordi Benet-Buchholz, Roland Boese, Maurice Leslie, C. Richard A. Catlow, and Sarah L. Price. Solid phases of cyclopentane: combined experimental and simulation study. *J. Phys. Chem. B*, 112(12):3746–58, 2008.
- [101] Elena Badea, Ignazio Blanco, and Giuseppe Della Gatta. Fusion and solid-to-solid transitions of a homologous series of alkane- α,ω -dinitriles. *J. Chem. Thermodyn.*, 39(10):1392–1398, 2007.
- [102] Anthony Maher, Åke C. Rasmuson, Denise M. Croker, and Benjamin K. Hodnett. Solubility of the metastable polymorph of piracetam (Form II) in a range of solvents. *J. Chem. Eng. Data*, 57(12):3525–3531, 2012.
- [103] Xavier Alcobe, Eugenia Estop, Kenneth D M Harris, Abil E. Aliev, Juan Rodriguez-Carvajal, and Jordi Rius. Temperature-Dependent Structural Properties of p-Diiodobenzene- Neutron Diffraction and High-Resolution Solid State C13 NMA Investigations Rius JSSC 1994.pdf. *J. Solid State Chem.*, 110(1):20–27, 1994.
- [104] Saikat Roy, N. Rajesh Goud, and Adam J. Matzger. Polymorphism in phenobarbital: discovery of a new polymorph and crystal structure of elusive form V. *Chem. Commun.*, 52(23):4389–4392, 2016.
- [105] Sarah A. Barnett, Ashley T. Hulme, Nizar Issa, Thomas C. Lewis, Louise S. Price, Derek A. Tocher, and Sarah (Sally) L. Price. The observed and energetically feasible crystal structures of 5-substituted uracils. *New J. Chem.*, 32(10):1761, 2008.
- [106] Doris E. Braun, Thomas Gelbrich, Klaus Wurst, and Ulrich J. Griesser. Computational and Experimental Characterization of Five Crystal Forms of Thymine: Packing Polymorphism, Polytypism/Disorder and Stoichiometric 0.8-Hydrate. *Cryst. Growth Des.*, 16(6):3480–3496, 2016.
- [107] Solid Form Solutions Ltd, 2015.
- [108] Amazon EC2, 2016.
- [109] Jos P. M. Lommerse, W. D. Sam Motherwell, Herman L. Ammon, Jack D. Dunitz, Angelo Gavezzotti, Detlef W. M. Hofmann, Frank J. J. Leusen, Wijnand T. M. Mooij, Sarah L. Price, Bernd Schweizer, Martin U. Schmidt, Bouke P. van Eijck, Paul Verwer, and Donald E. Williams. A test of crystal structure prediction of small organic molecules. *Acta Crystallogr. Sect. B Struct. Sci.*, 56(4):697–714, aug 2000.

- [110] W. D. Sam Motherwell, Herman L. Ammon, Jack D. Dunitz, Alexander Dzyabchenko, Peter Erk, Angelo Gavezzotti, Detlef W. M. Hofmann, Frank J. J. Leusen, Jos P. M. Lommerse, Wijnand T. M. Mooij, Sarah L. Price, Harold Scheraga, Bernd Schweizer, Martin U. Schmidt, Bouke P. van Eijck, Paul Verwer, and Donald E. Williams. Crystal structure prediction of small organic molecules: a second blind test. *Acta Crystallogr. Sect. B Struct. Sci.*, 58(4):647–661, jul 2002.
- [111] G M Day, W D S Motherwell, H L Ammon, S X M Boerrigter, R G Della Valle, E Venuti, A Dzyabchenko, J D Dunitz, B Schweizer, B P van Eijck, P Erk, J C Facelli, V E Bazterra, M B Ferraro, D W M Hofmann, F J J Leusen, C Liang, C C Pantelides, P G Karamertzanis, S L Price, T C Lewis, H Nowell, A Torrisi, H A Scheraga, Y A Arnautova, M U Schmidt, and P Verwer. A third blind test of crystal structure prediction. *Acta Crystallogr. B.*, 61(Pt 5):511–527, oct 2005.
- [112] Graeme M Day, Timothy G Cooper, Aurora J Cruz-Cabeza, Katarzyna E Hejczyk, Herman L Ammon, Stephan X M Boerrigter, Jeffrey S Tan, Raffaele G Della Valle, Elisabetta Venuti, Jovan Jose, Shridhar R Gadre, Gautam R Desiraju, Tejender S Thakur, Bouke P van Eijck, Julio C Facelli, Victor E Bazterra, Marta B Ferraro, Detlef W M Hofmann, Marcus A Neumann, Frank J J Leusen, John Kendrick, Sarah L Price, Alston J Misquitta, Panagiotis G Karamertzanis, Gareth W A Welch, Harold A Scheraga, Yelena A Arnautova, Martin U Schmidt, Jacco van de Streek, Alexandra K Wolf, and Bernd Schweizer. Significant progress in predicting the crystal structures of small organic molecules—a report on the fourth blind test. *Acta Crystallogr. B.*, 65:107–25, apr 2009.
- [113] David A Bardwell, Claire S Adjiman, Yelena A Arnautova, Ekaterina Bartashevich, Stephan X M Boerrigter, Doris E Braun, Aurora J Cruz-Cabeza, Graeme M Day, Raffaele G Della Valle, Gautam R Desiraju, Bouke P van Eijck, Julio C Facelli, Marta B Ferraro, Damian Grillo, Matthew Habgood, Detlef W M Hofmann, Fridolin Hofmann, K V Jovan Jose, Panagiotis G Karamertzanis, Andrei V Kazantsev, John Kendrick, Liudmila N Kuleshova, Frank J J Leusen, Andrey V Maleev, Alston J Misquitta, Sharmarke Mohamed, Richard J Needs, Marcus A Neumann, Denis Nikylov, Anita M Orendt, Rumpa Pal, Constantinos C Pantelides, Chris J Pickard, Louise S Price, Sarah L Price, Harold A Scheraga, Jacco van de Streek, Tejender S Thakur, Siddharth Tiwari, Elisabetta Venuti, and Ilia K Zhitkov. Towards crystal structure prediction of complex organic compounds—a report on the fifth blind test. *Acta Crystallogr. B.*, 67:535–51, dec 2011.
- [114] Anthony M. Reilly, Richard I. Cooper, Claire S. Adjiman, Saswata Bhattacharya, A. Daniel Boese, Jan Gerit Brandenburg, Peter J. Bygrave, Rita Bylsma, Josh E. Campbell, Roberto Car, David H. Case, Renu Chadha, Jason C. Cole, Katherine Cosburn, Herma M. Cuppen, Farren Curtis, Graeme M. Day, Robert A. DiStasio, Alexander Dzyabchenko, Bouke P. Van Eijck, Dennis M. Elking, Joost A. Van Den Ende, Julio C. Facelli, Marta B. Ferraro, Laszlo Fusti-Molnar, Christina Anna Gatsiou, Thomas S. Gee, Rene De Gelder, Luca M. Ghiringhelli, Hitoshi Goto, Stefan Grimme, Rui Guo, Detlef W M Hofmann, Johannes Hoja, Rebecca K. Hylton, Luca Iuzzolino, Wojciech Jankiewicz, Daniel T. De Jong, John Kendrick, Niek J J De Klerk, Hsin Yu

- Ko, Liudmila N. Kuleshova, Xiayue Li, Sanjaya Lohani, Frank J J Leusen, Albert M. Lund, Jian Lv, Yanming Ma, Noa Marom, Artem E. Masunov, Patrick McCabe, David P. McMahon, Hugo Meekes, Michael P. Metz, Alston J. Misquitta, Sharmarke Mohamed, Bartomeu Monserrat, Richard J. Needs, Marcus A. Neumann, Jonas Nyman, Shigeaki Obata, Harald Oberhofer, Artem R. Oganov, Anita M. Orendt, Gabriel I. Pagola, Constantinos C. Pantelides, Chris J. Pickard, Rafal Podeszwa, Louise S. Price, Sarah L. Price, Angeles Pulido, Murray G. Read, Karsten Reuter, Elia Schneider, Christoph Schober, Gregory P. Shields, Pawanpreet Singh, Isaac J. Sugden, Krzysztof Szalewicz, Christopher R. Taylor, Alexandre Tkatchenko, Mark E. Tuckerman, Francesca Vacarro, Manolis Vasileiadis, Alvaro Vazquez-Mayagoitia, Leslie Vogt, Yanchao Wang, Rona E. Watson, Gilles A. De Wijs, Jack Yang, Qiang Zhu, and Colin R. Groom. Report on the sixth blind test of organic crystal structure prediction methods. *Acta Crystallogr. Sect. B Struct. Sci. Cryst. Eng. Mater.*, 72(4):439–459, 2016.
- [115] Salvatore Cardamone, Timothy J Hughes, and Paul L A Popelier. Multipolar electrostatics. *Phys. Chem. Chem. Phys.*, 16(22):10367–87, 2014.
- [116] Carole Ouvrard and Sarah L. Price. Toward Crystal Structure Prediction for Conformationally Flexible Molecules: The Headaches Illustrated by Aspirin. *Cryst. Growth Des.*, 4(6):1119–1127, nov 2004.
- [117] Vijay K. Srirambhatla, Rui Guo, Sarah L. Price, and Alastair J. Florence. Isomorphous template induced crystallisation: a robust method for the targeted crystallisation of computationally predicted metastable polymorphs. *Chem. Commun.*, 52(46):7384–7386, 2016.
- [118] Marcus A Neumann. Tailor-Made Force Fields for Crystal-Structure Prediction. *J. Phys. Chem. B*, 112(32):9810–9829, 2008.
- [119] Eric C. Dybeck, Natalie P. Schieber, and Michael R. Shirts. Effects of a More Accurate Polarizable Hamiltonian on Polymorph Free Energies Computed Efficiently by Reweighting Point-Charge Potentials. *J. Chem. Theory Comput.*, 12(8):3491–3505, 2016.
- [120] Marcello Colapietro, Aldo Domenicano, Gustavo Portalone, György Schultz, and István Hargittai. Molecular Structure of p -Diaminobenzene in the Gaseous Phase and in the Crystal. *J. Phys. Chem. A*, 91(15):1728–1737, 1987.
- [121] A Warshel, E Huler, D Rabinovich, and Z Shakked. Examination of intramolecular potential surfaces of flexible conjugated molecules by calculation of crystal structures. equilibrium geometries of chalcones and diphenyloctatetraene in the crystal and gaseous state. *J. Mol. Struct.*, 23:175–191, 1974.
- [122] Arne Almenningen, Otto Bastiansen, Liv Fernholt, Bjorg Cyvin, Sven Cyvin, and Svein Samdal. Structure and Barrier of Internal Rotation of Biphenyl Derivatives in the Gaseous State. *J. Mol. Struct.*, 128:59–76, 1985.
- [123] J. L. Baudour and M. Sanquer. Structural phase transition in polyphenyls. VIII. The modulated structure of phase III of biphenyl (T=20 K) from neutron diffraction data. *Acta Crystallogr. Sect. B*, B39(1):75–84, 1983.

- [124] Thomas Steiner. research papers Frequency of Z H values in organic and organo- metallic crystal structures research papers. *Struct. Sci.*, B56:673–676, 2000.
- [125] Panagiotis G. Karamertzanis and Constantinos C. Pantelides. Ab initio crystal structure prediction - I. Rigid molecules. *J. Comput. Chem.*, 26:304–324, 2005.
- [126] P. G. Karamertzanis and C. C. Pantelides. Ab initio crystal structure prediction. II. Flexible molecules. *Mol. Phys.*, 105(2-3):273–291, 2007.
- [127] Marcus a Neumann and Marc-Antoine Perrin. Energy ranking of molecular crystals using density functional theory calculations and an empirical van der waals correction. *J. Phys. Chem. B*, 109(32):15531–15541, 2005.
- [128] Victor E. Bazterra, Marta B. Ferraro, and Julio C. Facelli. Modified genetic algorithm to model crystal structures. I. Benzene, naphthalene and anthracene. *J. Chem. Phys.*, 116(14):5984–5991, 2002.
- [129] Victor E. Bazterra, Marta B. Ferraro, and Julio C. Facelli. Modified genetic algorithm to model crystal structures. II. Determination of a polymorphic structure of benzene using enthalpy minimization. *J. Chem. Phys.*, 116(14):5992–5995, 2002.
- [130] Victor E. Bazterra, Marta B. Ferraro, and Julio C. Facelli. Modified genetic algorithm to model crystal structures. III. Determination of Crystal Structures Allowing Simultaneous Molecular Geometry Relaxation. *Int. J. Quantum Chem.*, 96(4):312–320, 2004.
- [131] Colin W. Glass, Artem R. Oganov, and Nikolaus Hansen. USPEX Evolutionary crystal structure prediction. *Comput. Phys. Commun.*, 175(11-12):713–720, 2006.
- [132] Artem R Oganov and Colin W Glass. Crystal structure prediction using ab initio evolutionary techniques: principles and applications. *J. Chem. Phys.*, 124(2006):244704, 2006.
- [133] D. W. M. Hofmann and T. Lengauer. A Discrete Algorithm for Crystal Structure Prediction of Organic Molecules. *Acta Crystallogr. Sect. A*, A53:225–235, 1997.
- [134] Detlef W.M. Hofmann and Thomas Lengauer. Prediction of crystal structures of organic molecules. *J. Mol. Struct.*, 474(98):13–23, 1999.
- [135] Wijnand T M Mooij, Bouke P Van Eijck, and Jan Kroon. Ab Initio Crystal Structure Predictions for Flexible Hydrogen-Bonded Molecules. *J. Am. Chem. Soc.*, 122(14):3500–3505, 2000.
- [136] Angelo Gavezzotti. OPiX, 2003.
- [137] Jay W. Ponder. TINKER Molecular Modeling Package, 2015.
- [138] Sander Pronk, Szilárd Páll, Roland Schulz, Per Larsson, Pär Bjelkmar, Rossen Apostolov, Michael R Shirts, Jeremy C Smith, Peter M Kasson, David van der Spoel, Berk Hess, and Erik Lindahl. GROMACS 4.5: a high-throughput and highly parallel open source molecular simulation toolkit. *Bioinformatics*, 29(7):845–54, apr 2013.

- [139] Richard Byrd, Peihuang Lu, Jorge Nocedal, and Ciyu Zhu. A limited. *J. Sci. Comput.*, 16(5):1190–1208, 1995.
- [140] Graeme M. Day, James Chisholm, Ning Shan, W. D Sam Motherwell, and William Jones. An assessment of lattice energy minimization for the prediction of molecular organic crystal structures. *Cryst. Growth Des.*, 4(6):1327–1340, 2004.
- [141] Panagiotis G Karamertzanis, Paolo Raiteri, Michele Parrinello, Maurice Leslie, and Sarah L Price. The thermal stability of lattice-energy minima of 5-fluorouracil: metadynamics as an aid to polymorph prediction. *J. Phys. Chem. B*, 112(14):4298–4308, apr 2008.
- [142] Jonas Nyman and Graeme M. Day. Static and lattice vibrational energy differences between polymorphs. *CrystEngComm*, 17(28):5154–5165, 2015.
- [143] Carsten Müller and Daniel Spångberg. Calculation of the stability of nonperiodic solids using classical force fields and the method of increments: N2O as an example. *J. Comput. Chem.*, 36(18):1420–1427, 2015.
- [144] Paolo Raiteri, Roman Martonák, and Michele Parrinello. Exploring polymorphism: the case of benzene. *Angew. Chem. Int. Ed. Engl.*, 44(24):3769–3773, jun 2005.
- [145] Tai Boon Tan, Andrew J Schultz, and David a Kofke. Efficient calculation of temperature dependence of solid-phase free energies by overlap sampling coupled with harmonically targeted perturbation. *J. Chem. Phys.*, 133(13):134104, oct 2010.
- [146] Tai Boon Tan, Andrew J Schultz, and David a Kofke. Suitability of umbrella- and overlap-sampling methods for calculation of solid-phase free energies by molecular simulation. *J. Chem. Phys.*, 132(21):214103, jun 2010.
- [147] Tang-Qing Yu and Mark E. Tuckerman. Temperature-Accelerated Method for Exploring Polymorphism in Molecular Crystals Based on Free Energy. *Phys. Rev. Lett.*, 107(1):015701, jun 2011.
- [148] Tang-Qing Yu, Pei-Yang Chen, Ming Chen, Amit Samanta, Eric Vanden-Eijnden, and Mark Tuckerman. Order-parameter-aided temperature-accelerated sampling for the exploration of crystal polymorphism and solid-liquid phase transitions. *J. Chem. Phys.*, 140(21):214109, jun 2014.
- [149] Ming Chen, Tang-Qing Yu, and Mark E. Tuckerman. Locating landmarks on high-dimensional free energy surfaces. *Proc. Natl. Acad. Sci.*, 112(11):3235–3240, 2015.
- [150] A. Bruce, N. Wilding, and G. Ackland. Free Energy of Crystalline Solids: A Lattice-Switch Monte Carlo Method. *Phys. Rev. Lett.*, 79(16):3002–3005, oct 1997.
- [151] Ad Bruce, An Jackson, Gj Ackland, and Nb Wilding. Lattice-switch monte carlo method. *Phys. Rev. E. Stat. Phys. Plasmas. Fluids. Relat. Interdiscip. Topics*, 61(1):906–919, jan 2000.
- [152] Dorothea Wilms, Nigel B. Wilding, and Kurt Binder. Transitions between imperfectly ordered crystalline structures: A phase switch Monte Carlo study. *Phys. Rev. E*, 85(5):056703, may 2012.

- [153] Michael J. Schnieders, Jonas Baltrusaitis, Yue Shi, Gaurav Chattree, Lianqing Zheng, Wei Yang, and Pengyu Ren. The structure, thermodynamics, and solubility of organic crystals from simulation with a polarizable force field. *J. Chem. Theory Comput.*, 8(5):1721–1736, 2012.
- [154] Jooyeon Park, Ian Nessler, Brian McClain, Dainius Macikenas, Jonas Baltrusaitis, and Michael J Schnieders. Absolute Organic Crystal Thermodynamics: Growth of the Asymmetric Unit into a Crystal via Alchemy. *J. Chem. Theory Comput.*, 10(7):2781–2791, 2014.
- [155] Wolfgang Beckmann, Roland Boisteile, and Kiyotaka Sate. Solubility of the A , B , and C Polymorphs of Stearic Acid in Decane , Methanol , and Butanone. *J. Chem. Eng. Data*, 29(2):211–214, 1984.
- [156] J. Willard Gibbs. *Elementary Principles in Statistical Mechanics*. Charles Scribner’s Sons, New York, 1902.
- [157] Richard Chace Tolman. *The Principles of Statistical Mechanics*. Dover Publications, 1938.
- [158] Chia-en Chang, Wei Chen, and Michael K Gilson. Evaluating the Accuracy of the Quasiharmonic Approximation. *J. Chem. Theory Comput.*, 1(5):1017–1028, 2005.
- [159] A. Gavezzotti and G Filippini. Polymorphic Forms of Organic Crystals at Room Conditions: Thermodynamic and Structural Implications. *J. Am. Chem. Soc.*, 117(49):12299–12305, 1995.
- [160] R Ramírez, N Neuerburg, M-V Fernández-Serra, and C P Herrero. Quasi-harmonic approximation of thermodynamic properties of ice Ih, II, and III. *J. Chem. Phys.*, 137(2012):044502, 2012.
- [161] A. Otero-de-la Roza and Victor Luana. Equations of state and thermodynamics of solids using empirical corrections in the quasiharmonic approximation. *Phys. Rev. B*, 84:184103, 2011.
- [162] James C. Phillips, Rosemary Braun, Wei Wang, James Gumbart, Emad Tajkhorshid, Elizabeth Villa, Christophe Chipot, Robert D. Skeel, Laxmikant Kalé, and Klaus Schulten. Scalable molecular dynamics with NAMD. *J. Comput. Chem.*, 26(16):1781–1802, 2005.
- [163] B R Brooks, C L Brooks Iii, A D Mackerell, L Nilsson, R J Petrella, B Roux, Y Won, G Archontis, C Bartels, S Boresch, A Caflisch, L Caves, Q Cui, A R Dinner, and M Feig. CHARMM : The Biomolecular Simulation Program. *J. Comput. Chem.*, 30(10):1545–1614, 2009.
- [164] David A. Case, Thomas E. Cheatham, Tom Darden, Holger Gohlke, Ray Luo, Kenneth M. Merz, Alexey Onufriev, Carlos Simmerling, Bing Wang, and Robert J. Woods. The Amber biomolecular simulation programs. *J. Comput. Chem.*, 26(16):1668–1688, 2005.
- [165] S Plimpton. Fast parallel algorithms for short-range molecular dynamics, 1995.
- [166] Tom Darden, Darrin York, and Lee Pedersen. Particle mesh Ewald: An Nlog(N) method for Ewald sums in large systems. *J. Chem. Phys.*, 98(12):10089, 1993.
- [167] P. P. Ewald. Die Berechnung optischer und elektrostatischer Gitterpotentiale. *Ann. Phys.*, 369(3):253–287, 1921.

-
- [168] Ulrich Essmann, Lalith Perera, Max L Berkowitz, Tom Darden, Hsing Lee, and Lee G Pedersen. A smooth particle mesh Ewald method. *J Chem Phys*, 103(1995):8577–8593, 1995.
- [169] Hans C. Andersen. Molecular dynamics simulations at constant pressure and/or temperature. *J. Chem. Phys.*, 72(4):2384, 1980.
- [170] H. J. C. Berendsen, J. P. M. Postma, W. F. van Gunsteren, a. DiNola, and J. R. Haak. Molecular dynamics with coupling to an external bath. *J. Chem. Phys.*, 81(8):3684, 1984.
- [171] Giovanni Bussi, Davide Donadio, and Michele Parrinello. Canonical sampling through velocity rescaling. *J. Chem. Phys.*, 126(1), 2007.
- [172] William G Hoover. Canonical dynamics: Equilibrium phase-space distributions. *Phys. Rev. A*, 31(3):1695–1697, 1985.
- [173] N. Goga, A. J. Rzepiela, A. H. de Vries, S. J. Marrink, and H. J. C. Berendsen. Efficient Algorithms for Langevin and DPD Dynamics. *J. Chem. Theory Comput.*, 8(10):3637–3649, 2012.
- [174] M. Parrinello and A. Rahman. Polymorphic Transitions in Single Crystals: a New Molecular Dynamics Method. *J. Appl. Phys.*, 52(12):7182–7190, 1981.
- [175] Mark E Tuckerman, José Alejandro, Roberto López-Rendón, Andrea L Jochim, and Glenn J Martyna. A Liouville-operator derived measure-preserving integrator for molecular dynamics simulations in the isothermal-isobaric ensemble. *J. Phys. A: Math. Gen.*, 39(19):5629–5651, 2006.
- [176] Jerome P Nilmeier, Gavin E Crooks, David D L Minh, and John D Chodera. Nonequilibrium candidate Monte Carlo is an efficient tool for equilibrium simulation. *Proc. Natl. Acad. Sci. U. S. A.*, 108(45):E1009–1018, nov 2011.
- [177] A. C. Bailey and B. Yates. Anisotropic thermal expansion of pyrolytic graphite at low temperatures. *J. Appl. Phys.*, 41(13):5088–5091, 1970.
- [178] H. Iwanaga, A. Kunishige, and S. Takeuchi. Anisotropic thermal expansion in wurtzite-type crystals. *J. Mater. Sci.*, 35(10):2451–2454, 2000.
- [179] Michael R Shirts. Simple Quantitative Tests to Validate Sampling from Thermodynamic Ensembles. *J. Chem. Theory Comput.*, 9(2):909–926, 2012.
- [180] Erik E Santiso and Bernhardt L Trout. A general set of order parameters for molecular crystals. *J. Chem. Phys.*, 134(6):064109, feb 2011.
- [181] Nathan Duff and Baron Peters. Polymorph specific RMSD local order parameters for molecular crystals and nuclei: α -, β - and γ -glycine. *J. Chem. Phys.*, 135(13), 2011.

- [182] Alexandar T Tzanov, Michel a Cuendet, and Mark E Tuckerman. How accurately do current force fields predict experimental peptide conformations? An adiabatic free energy dynamics study. *J. Phys. Chem. B*, 118(24):6539–52, jun 2014.
- [183] David M Eike, Joan F Brennecke, and Edward J Maginn. Toward a robust and general molecular simulation method for computing solid-liquid coexistence. *J. Chem. Phys.*, 122(1):14115, jan 2005.
- [184] David M Eike and Edward J Maginn. Atomistic simulation of solid-liquid coexistence for molecular systems: application to triazole and benzene. *J. Chem. Phys.*, 124(16):164503, apr 2006.
- [185] Saivenkataraman Jayaraman and Edward J Maginn. Computing the melting point and thermodynamic stability of the orthorhombic and monoclinic crystalline polymorphs of the ionic liquid 1-n-butyl-3-methylimidazolium chloride. *J. Chem. Phys.*, 127(21):214504, dec 2007.
- [186] Martin B. Sweatman. Comparison of absolute free energy calculation methods for fluids and solids. *Mol. Phys.*, 113(9-10):1206–1216, 2015.
- [187] Song-Ho Chong and Sihyun Ham. Thermodynamic-Ensemble Independence of Solvation Free Energy. *J. Chem. Theory Comput.*, 11(2):378–380, jan 2015.
- [188] Michael R Shirts and John D Chodera. Statistically optimal analysis of samples from multiple equilibrium states. *J. Chem. Phys.*, 129(12):124105, 2008.
- [189] Zhiqiang Tan. On a Likelihood Approach for Monte Carlo Integration. *J. Am. Stat. Assoc.*, 99(468):1027–1036, 2004.
- [190] Charles H Bennett. Efficient Estimation of Free Energy Differences from Monte Carlo Data. *J. Comput. Phys.*, 22(2):245–268, 1976.
- [191] R Zwanzig. High - Temperature Equation of State by a Perturbation Method. I. Nonpolar Gases. *J. Chem. Phys.*, 22(8):1420–1426, 1954.
- [192] Jamshed Anwar and Dirk Zahn. Uncovering molecular processes in crystal nucleation and growth by using molecular simulation. *Angew. Chem. Int. Ed. Engl.*, 50(9):1996–2013, feb 2011.
- [193] Federico Giberti, Matteo Salvalaglio, and Michele Parrinello. Metadynamics studies of crystal nucleation. *IUCrJ*, 2(2):256–266, 2015.
- [194] Johannes Seibel, Manfred Parschau, and Karl-Heinz Ernst. From Homochiral Clusters to Racemate Crystals: Viable Nuclei in 2D Chiral Crystallization. *J. Am. Chem. Soc.*, 137(25):7970–7973, 2015.
- [195] Manas Shah, Erik E Santiso, and Bernhardt L Trout. Computer simulations of homogeneous nucleation of benzene from the melt. *J. Phys. Chem. B*, 115(35):10400–12, sep 2011.

- [196] C Heidelberger, N. K. Chaudhuri, P Danneberg, D Mooren, L Greisbach, R Duschinsky, R. J. Schnitzer, E. Plevin, and J. Scheiner. Fluorinated Pyrimidines, a new class of tumour-inhibitory compounds. *Nature*, 179:663–669, 1957.
- [197] V Bertolasi and M Sacerdoti. The Crystal and Molecular Structure of. *Acta Crystallogr. Sect. B*, B29:2549–2556, 1978.
- [198] Qiang Zhu, Alexander G Shtukenberg, Damien J Carter, Tang-Qing Yu, Jingxiang Yang, Ming Chen, Paolo Raiteri, Artem R Oganov, Boaz Pokroy, Iryna Polishchuk, Peter J Bygrave, Graeme M Day, Andrew L Rohl, Mark E Tuckerman, and Bart Kahr. Resorcinol Crystallization from the Melt: A New Ambient Phase and New Riddles. *J. Am. Chem. Soc.*, 138(14):4881–4889, 2016.
- [199] Kyungho Park, James M B Evans, Allan S Myerson, and James M B. Determination of Solubility of Polymorphs Using Differential Scanning Calorimetry Determination of Solubility of Polymorphs Using Differential Scanning Calorimetry. *Cryst. Growth Des.*, 3(6):991–995, 2003.
- [200] Doris E. Braun, Karol P. Nartowski, Yaroslav Z. Khimyak, Kenneth R. Morris, Stephen R. Byrn, and Ulrich J. Griesser. Structural Properties, Order-Disorder Phenomena, and Phase Stability of Orotic Acid Crystal Forms. *Mol. Pharm.*, 13(3):1012–1029, 2016.
- [201] W.T.M Mooij, B.P. van Eijck, S.L. Price, Paul Verwer, and J. Kr Loon. Crystal Structure Predictions for Acetic Acid. *J Comp. Chem.*, 19(4):459–474, 1998.
- [202] Sarah L Price. Predicting crystal structures of organic compounds. *Chem. Soc. Rev.*, 43(7):2098–111, 2014.
- [203] Royston C B Copley, Sarah A. Barnett, Panagiotis G. Karamertzanis, Kenneth D M Harris, Benson M. Kar-iuki, Mingcan Xu, E. Anne Nickels, Robert W. Lancaster, and Sarah L. Price. Predictable disorder versus polymorphism in the rationalization of structural diversity: A multidisciplinary study of eniluracil. *Cryst. Growth Des.*, 8(9):3474–3481, 2008.
- [204] Panagiotis G Karamertzanis, Andrei V Kazantsev, Nizar Issa, W A Gareth, Claire S Adjiman, Constantinos C Pantelides, and Sarah L Price. Can the Formation of Pharmaceutical Cocrystals Be Computationally Predicted I. Comparison of Lattice Energy. *Cryst. Growth Des.*, 9(1):442–453, 2009.
- [205] Matthew Habgood. Form II caffeine: A case study for confirming and predicting disorder in organic crystals. *Cryst. Growth Des.*, 11(8):3600–3608, 2011.
- [206] W L Jorgensen, W L Jorgensen, D S Maxwell, D S Maxwell, J Tirado-Rives, and J Tirado-Rives. Development and Testing of the OLPS All-Atom Force Field on Conformational Energetics and Properties of Organic Liquids. *J. Am. Chem. Soc.*, 118(15):11225–11236, 1996.
- [207] William L. Jorgensen and Nora A. McDonald. Development of an all-atom force field for heterocycles. Properties of liquid pyridine and diazenes. *J. Mol. Struct. THEOCHEM*, 424(1-2):145–155, 1998.

-
- [208] George A Kaminski, Richard A Friesner, Julian Tirado-rives, and William L Jorgensen. Comparison with Accurate Quantum Chemical Calculations on Peptides . *J. Phys. Chem. B*, 105(28):6474–6487, 2001.
- [209] Sarmimala Hore, Robert Dinnebier, Wen Wen, Johnathan Hanson, and Joachim Maier. Structure of plastic crystalline succinonitrile: High-resolution in situ powder diffraction. *Zeitschrift fur Anorg. und Allg. Chemie*, 635(1):88–93, 2009.
- [210] Himanshu Paliwal and Michael R Shirts. Using Multistate Reweighting to Rapidly and Efficiently Explore Molecular Simulation Parameters Space for Nonbonded Interactions. *J. Chem. Theory Comput.*, 9(11):4700–4717, 2013.
- [211] H. Müller-Krumbhaar and K. Binder. Dynamic properties of the Monte Carlo method in statistical mechanics. *J. Stat. Phys.*, 8(1):1–24, 1973.
- [212] William C. Swope, Hans C. Andersen, Peter H. Berens, and Kent R. Wilson. A computer simulation method for the calculation of equilibrium constants for the formation of physical clusters of molecules: Application to small water clusters. *J. Chem. Phys.*, 76(1):637–649, 1982.
- [213] John D Chodera, William C Swope, Jed W Pitera, Chaok Seok, and Ken A Dill. Use of the Weighted Histogram Analysis Method for the Analysis of Simulated and Parallel Tempering Simulations. *J. Chem. Theory Comput.*, 3(1):26–41, 2007.
- [214] John D Chodera. A simple method for automated equilibration detection in molecular simulations. *J. Chem. Theory Comput.*, 12(4):1799–1805, 2016.
- [215] Graeme M. Day, W D S Motherwell, and W. Jones. Beyond the Isotropic Atom Model in Crystal Structure Prediction Atomic Multipoles vs Point Charges. *Cryst. Growth Des.*, 5(3):1023–1033, 2005.
- [216] Theresa Beyer and Sarah L. Price. The errors in lattice energy minimisation studies: sensitivity to experimental variations in the molecular structure of paracetamol. *CrystEngComm*, 2:183–190, 2000.
- [217] Thomas C. Lewis, Derek A. Tocher, and Sarah L. Price. An experimental and theoretical search for polymorphs of barbituric acid: The challenges of even limited conformational flexibility. *Cryst. Growth Des.*, 4(5):979–981, 2004.
- [218] Sarah L. Price. The computational prediction of pharmaceutical crystal structures and polymorphism. *Adv. Drug Deliv. Rev.*, 56(3):301–319, 2004.
- [219] G M Day, W D S Motherwell, and W Jones. A strategy for predicting the crystal structures of flexible molecules: the polymorphism of phenobarbital. *Phys. Chem. Chem. Phys.*, 9(14):1693–1704, 2007.
- [220] Gregory J. O. Beran and Kaushik Nanda. Predicting Organic Crystal Lattice Energies with Chemical Accuracy. *J. Phys. Chem. Lett.*, 1(24):3480–3487, dec 2010.

- [221] Garnet Kin-Lic Chan Jun Yang, Weifeng Hu, Denis Usvyat, Devin Matthews, Martin Schutz. Ab Initio Determination of the Crystalline Benzene Lattice Energy to sub-kilojoule/mol Accuracy. *Science*, 345(6197):640–643, 2014.
- [222] Dian Jiao, Pavel A Golubkov, Thomas A Darden, and Pengyu Ren. Calculation of protein ligand binding free energy by using a polarizable potential. *Proc. Natl. Acad. Sci.*, 105(17):6290–6295, 2008.
- [223] Alan Grossfield, Pengyu Ren, and Jay W. Ponder. Ion Solvation Thermodynamics from Simulation with a Polarizable Force Field. *J. Am. Chem. Soc.*, 125(50):15671–15682, 2003.
- [224] Zhixiong Lin and Wilfred F. van Gunsteren. Effects of Polarizable Solvent Models upon the Relative Stability of an α -Helical and a β -Hairpin Structure of an Alanine Decapeptide. *J. Chem. Theory Comput.*, 11(5):1983–1986, 2015.
- [225] Filip Moučka, Ivo Nezbeda, and William R. Smith. Chemical Potentials, Activity Coefficients, and Solubility in Aqueous NaCl Solutions: Prediction by Polarizable Force Fields. *J. Chem. Theory Comput.*, 11(4):1756–1764, 2015.
- [226] Hao Jiang, Zoltan Mester, Othonas A. Moultos, Ioannis G. Economou, and Athanassios Z. Panagiotopoulos. Thermodynamic and Transport Properties of H₂O + NaCl from Polarizable Force Fields. *J. Chem. Theory Comput.*, 11(8):3802–3810, 2015.
- [227] Xiangda Peng, Yuebin Zhang, Huiying Chu, Yan Li, Dinglin Zhang, Liaoran Cao, and Guohui Li. Accurate Evaluation of Ion Conductivity of the Gramicidin A Channel Using a Polarizable Force Field without Any Corrections. *J. Chem. Theory Comput.*, 12(6):2973–2982, 2016.
- [228] X Zheng, C Wu, J W Ponder, and G R Marshall. Molecular Dynamics of β -Hairpin Models of Epigenetic Recognition Motifs. *J. Am. Chem. Soc.*, 134(38):15970–15978, 2012.
- [229] D E Williams and R R Weller. Lone-pair electronic effects on the calculated ab initio SCF-MO electric potential and the crystal structures of azabenzenes. *J. Am. Chem. Soc.*, 105(10):4143–4148, 1983.
- [230] Ds Coombes, Sl Price, Dj Willock, and M Leslie. Role of electrostatic interactions in determining the crystal structures of polar organic molecules. A distributed multipole study. *J. Phys. Chem.*, 3654(96):7352–7360, 1996.
- [231] Wijnand T. M. Mooij, Bouke P. van Eijck, and Jan Kroon. Transferable ab Initio Intermolecular Potentials. 2. Validation and Application to Crystal Structure Prediction. *J. Phys. Chem. A*, 103(48):9883–9890, 1999.
- [232] Timothy G. Cooper, Katarzyna E. Hejczyk, William Jones, and Graeme M. Day. Molecular polarization effects on the relative energies of the real and putative crystal structures of valine. *J. Chem. Theory Comput.*, 4(10):1795–1805, 2008.

- [233] Sarah L. Price, Maurice Leslie, Gareth W. A. Welch, Matthew Habgood, Louise S. Price, Panagiotis G. Karamertzanis, and Graeme M Day. Modelling organic crystal structures using distributed multipole and polarizability-based model intermolecular potentials. *Phys. Chem. Chem. Phys.*, 12(30):8466–8477, 2010.
- [234] Anthony M. Reilly and Alexandre Tkatchenko. Role of dispersion interactions in the polymorphism and entropic stabilization of the aspirin crystal. *Phys. Rev. Lett.*, 113(5):1–5, 2014.
- [235] Carlos E. S. Bernardes and Abhinav Joseph. Evaluation of the OPLS-AA Force Field for the Study of Structural and Energetic Aspects of Molecular Organic Crystals. *J. Phys. Chem. A*, 119(12):3023–3034, 2015.
- [236] Alberto Milani, Chiara Castiglioni, and Stefano Radice. Joint Experimental and Computational Investigation of the Structural and Spectroscopic Properties of Poly(vinylidene fluoride) Polymorphs. *J. Phys. Chem. B*, 119(14):4888–4897, 2015.
- [237] Anthony M. Reilly, Richard I. Cooper, Claire S. Adjiman, Saswate Bhattacharya, A. Daniel Boese, Jan Gerit Brandenburg, Peter J. Bygrave, Rita Bylsma, Josh E. Campbell, Roberto Car, David H. Case, Renu Chadha, Jason C. Cole, Katherine Cosburn, Herma M. Cuppen, Farren Curtis, Graeme M. Day, Robert A. DiStasio Jr, Alexander Dzyabchenko, Bouke P. van Eijck, Dennis M. Elking, Joost A. van den Ende, Julio C. Facelli, Marta B. Ferraro, Laszlo Fusti-Molnar, Chritina-Anna Gatsiou, Thomas S. Gee, Rene de Gelder, Luca M. Ghiringhelli, Hitoshi Goto, Stefan Grimme, Rui Guo, Detlef W.M. Hofmann, Johannes Hoja, Rebecca K. Hylton, Luca Iuzzolino, Wojciech Jankiewicz, Daniel T. de Jong, John Kendrick, Niek J.J. de Klerk, Hsin-Yu Ko, Liudmila N. Kuleshova, Xiayue Li, Sanjaya Lohani, Frank J.J. Leusen, Albert M. Lund, Jian Lv, Yanming Ma, Noa Marom, Artem E. Masunov, Patrick McCabe, David P. McMahon, Hugo Meekes, Michael P. Metz, Alston J. Misquitta, Sharmarke Mohamed, Bartomeu Monserrat, Richard J. Needs, Marcus A. Neumann, Jonas Nyman, Shigeaki Obata, Harald Oberhofer, Artem R. Oganov, Anita M. Orendt, Gabriel I. Pagola, Constantinos C. Pantelides, Chris J. Pickard, Rafal Podeszwa, Louise S. Price, Sarah L. Price, Angeles Pulido, Murray G. Read, Karsten Reuter, Elia Schneider, Christoph Schober, Gregory P. Shields, Pawanpreet Singh, Isaac J. Sugden, Krzysztof Szaleqicz, Christopher R. Taylor, Alexandre Tkatchenko, Mark E. Tuckerman, Francesca Vacarro, Manolis Vasileiadis, Alvaro Vazquez-Mayagoitia, Leslie Vogt, Yanchao Wang, Rona E. Watson, Gilles A. de Wijs, Jack Yang, Qiang Zhu, and Colin R. Groom. Report on the sixth blind test of organic crystal-structure prediction methods. (February):439–459, 2016.
- [238] A Warshel and M Levitt. Theoretical studies of enzymic reactions: dielectric, electrostatic and steric stabilization of the carbonium ion in the reaction of lysozyme. *J. Mol. Biol.*, 103(2):227–249, 1976.
- [239] John G. Kirkwood. Statistical Mechanics of Fluid Mixtures. *J. Chem. Phys.*, 3(5):300–313, 1935.
- [240] N Vaidehi, T A Wesolowski, and A Warshel. Quantum-Mechanical Calculations of Solvation Free Energies. A Combined ab initio Pseudopotential Free-Energy Perturbation Approach. *J. Chem. Phys.*, 97(1992):4264–4271, 1992.

- [241] V Luzhkov and A Warshel. Microscopic models for quantum mechanical calculations of chemical processes in solutions: LD/AMPAC and SCAAS/AMPAC calculations of solvation energies. *J. Comput. Chem.*, 13(2):199–213, 1992.
- [242] Jiali Gao and Xinfu Xia. A Priori Evaluation of Aqueous Polarization Effects Through Monte Carlo QM-MM Simulations. *Science*, 258(5082):631–635, 1992.
- [243] Mats H M Olsson, Gongyi Hong, and Arieh Warshel. Frozen density functional free energy simulations of redox proteins: computational studies of the reduction potential of plastocyanin and rusticyanin. *J. Am. Chem. Soc.*, 125(17):5025–5039, 2003.
- [244] Edina Rosta, Marco Klähn, and Arieh Warshel. Towards accurate ab initio QM/MM calculations of free-energy profiles of enzymatic reactions. *J. Phys. Chem. B*, 110(6):2934–2941, 2006.
- [245] J Bentzien, R P Muller, J Florian, and Arieh Warshel. Hybrid ab initio quantum mechanics molecular mechanics calculations of free energy surfaces for enzymatic reactions: The nucleophilic attack in subtilisin. *J. Phys. Chem. B*, 102(12):2293–2301, 1998.
- [246] Nikolay V. Plotnikov, Shina C L Kamerlin, and Arieh Warshel. Paradynamics: An effective and reliable model for Ab initio QM/MM free-energy calculations and related tasks. *J. Phys. Chem. B*, 115(24):7950–7962, 2011.
- [247] Frank R Beierlein, Julien Michel, and Jonathan W Essex. A simple QM/MM approach for capturing polarization effects in protein-ligand binding free energy calculations. *J. Phys. Chem. B*, 115(17):4911–4926, may 2011.
- [248] Stephen John Fox, Chris Pittock, Christofer S Tautermann, Thomas Fox, Clara Christ, N. O. J. Malcom, Jon W Essex, and Chris-Kriton Skylaris. Free Energies of Binding from Large-Scale First Principles Quantum Mechanical Calculations : Application to Ligand Hydration Energies. *J. Phys. Chem. B*, 117(32):9478–9485, 2013.
- [249] Samuel Genheden, Ana I. Cabedo Martinez, Michael P. Criddle, and Jonathan W. Essex. Extensive all-atom Monte Carlo sampling and QM/MM corrections in the SAMPL4 hydration free energy challenge. *J. Comput.-Aided. Mol. Des.*, 28(3):187–200, 2014.
- [250] Christopher Cave-Ayland, Chris-Kriton Skylaris, and Jonathan W Essex. Direct Validation of the Single Step Classical to Quantum Free Energy Perturbation. *J. Phys. Chem. B*, 119(3):1017–1025, oct 2015.
- [251] Christopher J. Woods, Frederick R. Manby, and Adrian J. Mulholland. An efficient method for the calculation of quantum mechanics/molecular mechanics free energies. *J. Chem. Phys.*, 128(1):014109, 2008.
- [252] Chris Sampson, Thomas Fox, Christofer S Tautermann, Christopher J. Woods, and Chris-Kriton Skylaris. A “Stepping Stone” Approach for Obtaining Quantum Free Energies of Binding. *J. Phys. Chem. B*, 119(23):7030–7040, 2015.

- [253] Eric R. Pinnick, Camilo E. Calderon, Andrew J. Rusnak, and Feng Wang. Achieving fast convergence of ab initio free energy perturbation calculations with the adaptive force-matching method. *Theor. Chem. Acc.*, 131(3):1146, feb 2012.
- [254] Himanshu Paliwal and Michael R Shirts. Multistate reweighting and configuration mapping together accelerate the efficiency of thermodynamic calculations as a function of molecular geometry by orders of magnitude. *J. Chem. Phys.*, 138(15):154108, apr 2013.
- [255] Gerhard König and Stefan Boresch. Non-Boltzmann Sampling and Bennett’s Acceptance Ratio Method : How to Profit from Bending the Rules. *J. Comput. Chem.*, 32(6):1082–1090, 2011.
- [256] Gerhard König, Phillip S. Hudson, Stefan Boresch, and H. Lee Woodcock. Multiscale free energy simulations: An efficient method for connecting classical MD simulations to QM or QM/MM free energies using non-Boltzmann Bennett reweighting schemes. *J. Chem. Theory Comput.*, 10(4):1406–1419, 2014.
- [257] Gerhard König, Frank C. Pickard IV, Ye Mei, and Bernard R. Brooks. Predicting hydration free energies with a hybrid QM/MM approach: An evaluation of implicit and explicit solvation models in SAMPL4. *J. Comput.-Aided. Mol. Des.*, 28(3):245–257, 2014.
- [258] Phillip S Hudson, Justin K White, Fiona L Kearns, Milan Hodoscek, Stefan Boresch, and H Lee Woodcock. Efficiently computing pathway free energies: New approaches based on chain-of-replica and Non-Boltzmann Bennett reweighting schemes. *Biochim. Biophys. Acta*, 1850(5):944–53, 2015.
- [259] N. Lu and D. A. Kofke. Accuracy of free-energy perturbation calculations in molecular simulation. II. Heuristics. *J. Chem. Phys.*, 115(15):6866–6875, 2001.
- [260] Omololu Akin-Ojo, Yang Song, and Feng Wang. Developing ab initio quality force fields from condensed phase quantum-mechanics/molecular-mechanics calculations through the adaptive force matching method. *J. Chem. Phys.*, 129(6):064108, 2008.
- [261] Jimmy Heimdal and Ulf Ryde. Convergence of QM/MM free-energy perturbations based on molecular-mechanics or semiempirical simulations. *Phys. Chem. Chem. Phys.*, 14(36):12592, 2012.
- [262] Gerhard König and Bernard R. Brooks. Correcting for the free energy costs of bond or angle constraints in molecular dynamics simulations. *Biochim. Biophys. Acta - Gen. Subj.*, 1850(5):932–943, 2015.
- [263] Nathan Schmid, Andreas P Eichenberger, Alexandra Choutko, Sereina Riniker, Moritz Winger, Alan E Mark, and Wilfred F Van Gunsteren. Definition and testing of the GROMOS force-field versions 54A7 and 54B7. *Eur. Biophys. J.*, 40(7):843–856, 2011.
- [264] Pengyu Ren and J W Ponder. Polarizable Atomic Multipole Water Model for Molecular Mechanics Simulation. *J. Phys. Chem. B*, 107(24):5933–5947, 2003.

- [265] Jay W Ponder, Chuanjie Wu, Pengyu Ren, Vijay S Pande, John D Chodera, Michael J Schnieders, Imran Haque, David L Mobley, Daniel S Lambrecht, Robert A DiStasio, Martin Head-Gordon, Gary N I Clark, Margaret E Johnson, and Teresa Head-Gordon. Current status of the AMOEBA polarizable force field. *J. Phys. Chem. B*, 114(8):2549–64, mar 2010.
- [266] Pengyu Ren, Chuanjie Wu, and Jay W. Ponder. Polarizable atomic multipole-based molecular mechanics for organic molecules. *J. Chem. Theory Comput.*, 7(10):3143–3161, 2011.
- [267] Kaushik D Nanda and Gregory J O Beran. Prediction of organic molecular crystal geometries from MP2-level fragment quantum mechanical/molecular mechanical calculations. *J. Chem. Phys.*, 137(17):174106, nov 2012.
- [268] Sarah Price. Why don't we find more polymorphs? *Acta Crystallogr. Sect. B*, 69:313–328, 2013.
- [269] A T Anghel, G M Day, and S L Price. A study of the known and hypothetical crystal structures of pyridine: why are there four molecules in the asymmetric unit cell? *CrystEngComm*, 4(62):348–355, 2002.
- [270] T.M. Winjnad, Bouke P. van Eijck, Sarah L Price, Paul Verwer, and Jan Kroon. Crystal Structure Predictions for Acetic Acid. *J. Comput. Chem.*, 19(4):459–474, 1997.
- [271] Marc Descamps, Natalia T. Correia, Patrick Derollez, Florence Danede, and Frédéric Capet. Plastic and glassy crystal states of caffeine. *J. Phys. Chem. B*, 109(33):16092–16098, 2005.
- [272] A. Srinivasan, F. Bermejo, A. de Andrés, J. Dawidowski, J. Zúñiga, and A. Criado. Evidence for a supercooled plastic-crystal phase in solid ethanol. *Phys. Rev. B*, 53(13):8172–8175, 1996.
- [273] Andreas Hermann. High-pressure phase transitions in rubidium and caesium hydroxides. *Phys. Chem. Chem. Phys.*, 2016.
- [274] Gautam R. Desiraju. Reflections on the hydrogen bond in crystal engineering. *Cryst. Growth Des.*, 11(4):896–898, 2011.
- [275] Peter A Kollman. Free-Energy Calculations - Applications to Chemical and Biochemical Phenomena. *Chem. Rev.*, 93(7):2395–2417, 1993.
- [276] David L. Mobley, Elise Dumont, John D. Chodera, and Ken A. Dill. Comparison of charge models for fixed-charge force fields: small-molecule hydration free energies in explicit solvent. *J. Phys. Chem. B*, 111(9):2242–54, 2007.
- [277] David L Mobley, Christopher I Bayly, Matthew D Cooper, Michael R Shirts, and Ken A. Dill. Small molecule hydration free energies in explicit solvent: An extensive test of fixed-charge atomistic simulations. *J. Chem. Theory Comput.*, 5(2):350–358, feb 2009.
- [278] J Marelius, T Hansson, and J Aqvist. Calculation of Ligand Binding Free Energies from Molecular Dynamics Simulations. *Int. J. Quantum Chem.*, 69(1):77–88, 1998.

- [279] William L Jorgensen. The many roles of computation in drug discovery. *Science*, 303(5665):1813–1818, 2004.
- [280] Emilio Gallicchio and Ronald M. Levy. Recent theoretical and computational advances for modeling proteinligand binding affinities. *Adv. Protein Chem. Struct. Biol.*, 85:27–80, 2011.
- [281] Jeff Wereszczynski and J. Andrew McCammon. Statistical mechanics and molecular dynamics in evaluating thermodynamic properties of biomolecular recognition. *Q. Rev. Biophys.*, 45(01):1–25, 2012.
- [282] Wilfred F. van Gunsteren, Dirk Bakowies, Riccardo Baron, Indira Chandrasekhar, Markus Christen, Xavier Daura, Peter Gee, Daan P. Geerke, Alice Glättli, Philippe H. Hünenberger, Mika A. Kastenholz, Chris Oostenbrink, Merijn Schenk, Daniel Trzesniak, Nico F. A. van der Vegt, and Haibo B. Yu. Biomolecular Modeling: Goals, Problems, Perspectives. *Angew. Chemie Int. Ed.*, 45(25):4064–4092, 2006.
- [283] José Nelson Onuchic and Peter G Wolynes. Theory of protein folding. *Curr. Opin. Struct. Biol.*, 14(1):70–75, 2004.
- [284] Martin Karplus and J Andrew Mccammon. Molecular dynamics simulations of biomolecules. *Nat. Struct. Biol.*, 9(9):646 – 652, 2002.
- [285] M. Karplus and J. Kuriyan. Molecular dynamics and protein function. *Proc. Natl. Acad. Sci. U. S. A.*, 102(19):6679–6685, 2005.
- [286] Anthony Nicholls, David L. Mobley, J. Peter Guthrie, John D. Chodera, Christopher I. Bayly, Matthew D. Cooper, and Vijay S. Pande. Predicting small-molecule solvation free energies: An informal blind test for computational chemistry. *J. Med. Chem.*, 51(4):769–779, 2008.
- [287] J Peter Guthrie. A Blind Challenge for Computational Solvation Free Energies : Introduction and Overview. *J. Phys. Chem. B*, 113(14):4501–4507, 2009.
- [288] Matthew T. Geballe, A. Geoffrey Skillman, Anthony Nicholls, J. Peter Guthrie, and Peter J. Taylor. The SAMPL2 blind prediction challenge: Introduction and overview. *J. Comput.-Aided. Mol. Des.*, 24(4):259–279, 2010.
- [289] Hari S Muddana, C Daniel Varnado, Christopher W Bielawski, Adam R Urbach, Lyle Isaacs, Matthew T Geballe, and Michael K Gilson. Blind prediction of host-guest binding affinities: a new SAMPL3 challenge. *J. Comput.-Aided. Mol. Des.*, 26(5):475–87, may 2012.
- [290] Pavel V. Klimovich and David L. Mobley. Predicting hydration free energies using all-atom molecular dynamics simulations and multiple starting conformations. *J. Comput.-Aided. Mol. Des.*, 24(4):307–316, 2010.
- [291] Gerhard König and Bernard R Brooks. Predicting binding affinities of host-guest systems in the SAMPL3 blind challenge: the performance of relative free energy calculations. *J. Comput.-Aided. Mol. Des.*, 26(5):543–50, may 2012.

- [292] Morgan Lawrenz, Jeff Wereszczynski, Juan Manuel Ortiz-Sánchez, Sara E Nichols, and J Andrew McCammon. Thermodynamic integration to predict host-guest binding affinities. *J. Comput.-Aided. Mol. Des.*, 26(5):569–76, may 2012.
- [293] David L. Mobley, Shuai Liu, David S. Cerutti, William C. Swope, and Julia E. Rice. Alchemical prediction of hydration free energies for SAMPL. *J. Comput.-Aided. Mol. Des.*, 26(5):551–562, 2012.
- [294] Emilio Gallicchio, Nanjie Deng, Peng He, Lauren Wickstrom, Alexander L. Perryman, Daniel N. Santiago, Stefano Forli, Arthur J. Olson, and Ronald M. Levy. Virtual screening of integrase inhibitors by large scale binding free energy calculations: The SAMPL4 challenge. *J. Comput.-Aided. Mol. Des.*, 28(4):475–490, 2014.
- [295] David L. Mobley, Karisa L. Wymer, Nathan M. Lim, and J. Peter Guthrie. Blind prediction of solvation free energies from the SAMPL4 challenge. *J. Comput.-Aided. Mol. Des.*, 28(3):135–150, 2014.
- [296] Paulius Mikulskis, Daniela Cioloboc, Milica Andrejić, Sakshi Khare, Joakim Brorsson, Samuel Genheden, Ricardo A. Mata, Pär Söderhjelm, and Ulf Ryde. Free-energy perturbation and quantum mechanical study of SAMPL4 octa-acid host-guest binding energies. *J. Comput.-Aided. Mol. Des.*, 28(4):375–400, 2014.
- [297] Niels Hansen and Wilfred F. Van Gunsteren. Practical aspects of free-energy calculations: A review. *J. Chem. Theory Comput.*, 10(7):2632–2647, 2014.
- [298] Franziska Fischer, Adrian Heidrich, Sebastian Greiser, Sigrid Benemann, Klaus Rademann, and Franziska Emmerling. Polymorphism of mechanochemically synthesized Cocrystals: a case study. *Cryst. Growth Des.*, 16(3):1701–1707, 2016.
- [299] Dritan Hasa, Elvio Carlino, and William Jones. Polymer-Assisted Grinding, a Versatile Method for Polymorph Control of Cocrystallization. *Cryst. Growth Des.*, 16(3):1772–1779, 2016.
- [300] Roberto Centore, Sandra Fusco, Fabio Capone, and Mauro Causà. Competition between Polar and Centrosymmetric Packings in Molecular Crystals: Analysis of Actual and Virtual Structures. *Cryst. Growth Des.*, 16(4):2260–2265, 2016.
- [301] Ingrid Hallsteinsen, Magnus Nord, Torsatein Bolstad, Per-erik Vullum, Jos Emiel Boschker, Paolo Longo, Ryota Takahashi, Randi Holmestad, Mikk Lippmaa, and Thomas Tybell. The effect of polar (111)-oriented SrTiO₃ on initial perovskite growth. (111), 2016.
- [302] David Allan and Stewart Clark. Comparison of the high-pressure and low-temperature structures of ethanol and acetic acid. *Phys. Rev. B*, 60(9):6328–6334, 1999.
- [303] Shankar Kumar, John M Rosenberg, Djamal Bouzida, Robert H Swendsen, and Peter A Kollman. THE weighted histogram analysis method for free-energy calculations on biomolecules. I. The method. *J. Comput. Chem.*, 13(8):1011–1021, 1992.

- [304] Himanshu Paliwal and Michael R Shirts. A Benchmark Test Set for Alchemical Free Energy Transformations and Its Use to Quantify Error in Common Free Energy Methods. *J. Comput. Chem.*, 7(12):4115–4134, 2011.
- [305] Marc Souaille and Benoît Roux. Extension to the weighted histogram analysis method : combining umbrella sampling with free energy calculations. *Comput. Phys. Commun.*, 135(1):40–57, 2001.
- [306] Xiangyu Jia, Meiting Wang, Yihan Shao, Gerhard Ko, Bernard R Brooks, John Z H Zhang, and Ye Mei. Calculations of Solvation Free Energy through Energy Reweighting from Molecular Mechanics to Quantum Mechanics. *J. Chem. Theory Comput.*, 2016.
- [307] Glenn M. Torrie and John P. Valleau. Monte Carlo free energy estimates using non-Boltzmann sampling: Application to the sub-critical Lennard-Jones fluid. *Chem. Phys. Lett.*, 28(4):578–581, 1974.
- [308] J Valleau. Nonphysical Sampling Distributions in Monte Carlo Free-Energy Estimation: Umbrella Sampling. *J. Comput. Phys.*, 23(2):187–199, 1977.
- [309] Donghong Min, Lianqing Zheng, William Harris, Mengen Chen, Chao Lv, and Wei Yang. Practically efficient QM/MM alchemical free energy simulations: The orthogonal space random walk strategy. *J. Chem. Theory Comput.*, 6(8):2253–2266, 2010.
- [310] Leslie Kish. *Survey Sampling*. Wiley, New York, 1965.
- [311] Bernard R. Brooks, Robert E. Bruccoleri, Barry D. Olafson, David J. States, S. Swaminathan, and Martin Karplus. CHARMM: A program for macromolecular energy, minimization, and dynamics calculations. *J. Comput. Chem.*, 4(2):187–217, 1983.
- [312] Ad MacKerell and D Bashford. All-atom empirical potential for molecular modeling and dynamics studies of proteins. *J. Phys. Chem. B*, 102(18):3586–3616, 1998.
- [313] Yihan Shao, Laszlo Fusti Molnar, Yousung Jung, Jörg Kussmann, Christian Ochsenfeld, Shawn T Brown, Andrew T B Gilbert, Lyudmila V Slipchenko, Sergey V Levchenko, Darragh P O'Neill, Robert A DiStasio, Rohini C Lochan, Tao Wang, Gregory J O Beran, Nicholas A Besley, John M Herbert, Ching Yeh Lin, Troy Van Voorhis, Siu Hung Chien, Alex Sodt, Ryan P Steele, Vitaly a Rassolov, Paul E Maslen, Prakashan P Korambath, Ross D Adamson, Brian Austin, Jon Baker, Edward F C Byrd, Holger Dachsel, Robert J Doerksen, Andreas Dreuw, Barry D Dunietz, Anthony D Dutoi, Thomas R Furlani, Steven R Gwaltney, Andreas Heyden, So Hirata, Chao-Ping Hsu, Gary Kedziora, Rustam Z Khallilulin, Phil Klunzinger, Aaron M Lee, Michael S Lee, Wanzhen Liang, Itay Lotan, Nikhil Nair, Baron Peters, Emil I Proynov, Piotr A Pieniazek, Young Min Rhee, Jim Ritchie, Edina Rosta, C David Sherrill, Andrew C Simmonett, Joseph E Subotnik, H Lee Woodcock, Weimin Zhang, Alexis T Bell, Arup K Chakraborty, Daniel M Chipman, Frerich J Keil, Arie Warshel, Warren J Hehre, Henry F Schaefer, Jing Kong, Anna I Krylov, Peter M W Gill, and Martin Head-Gordon. Advances in methods and algorithms in a modern quantum chemistry program package. *Phys. Chem. Chem. Phys.*, 8(27):3172–3191, 2006.

- [314] Lee Woodcock, Milan Hodoscek, Andrew Gilbert, Peter Gill, Henry Schaefer, and Bernard R. Brooks. Interfacing Q-Chem and CHARMM to Perform QM/MM Reaction Path Calculations. *J. Comput. Chem.*, 28(9):1485–1502, 2007.
- [315] H Lee Woodcock, Benjamin T Miller, Milan Hodoscek, Asim Okur, Joseph D Larkin, Jay W Ponder, and Bernard R Brooks. MSCALE : A General Utility for Multiscale Modeling. *J. Chem. Theory Comput.*, 7(4):1208–1219, 2011.
- [316] Thomas A. Halgren. The representation of van der Waals (vdW) interactions in molecular mechanics force fields: potential form, combination rules, and vdW parameters. *J. Am. Chem. Soc.*, 114(20):7827–7843, 1992.