

Driving and characterizing nucleation of urea and glycine polymorphs in water

Ziyue Zou^{a,†}, Eric R. Beyerle^{b,†}, Sun-Ting Tsai^c, and Pratyush Tiwary^{a,b,1}

^aDepartment of Chemistry and Biochemistry, University of Maryland, College Park, MD 20742; ^bInstitute for Physical Science and Technology, University of Maryland, College Park, MD 20742; ^cDepartment of Physics, University of Maryland, College Park, MD 20742

Crystal nucleation is relevant across the domains of fundamental and applied sciences. However, in many cases its mechanism remains unclear due to a lack of temporal or spatial resolution. To gain insights to the molecular details of nucleation, some form of molecular dynamics simulations is typically performed; these simulations, in turn, are limited by their ability to run long enough to sample the nucleation event thoroughly. To overcome the timescale limits in typical molecular dynamics simulations in a manner free of prior human bias, here we employ the machine learning augmented molecular dynamics framework “Reweighted Autoencoded Variational Bayes for enhanced sampling (RAVE)”. We study two molecular systems, urea and glycine in explicit all-atom water, due to their enrichment in polymorphic structures and common utility in commercial applications. From our simulations, we observe multiple back-and-forth liquid-solid transitions of different polymorphs and from these trajectories calculate the polymorph stability relative to the dissolved liquid state. We further observe that the obtained reaction coordinates and transitions are highly non-classical.

Molecular Simulations | Nucleation | Enhanced Sampling | Machine Learning

Introduction

Although traditionally viewed through the lens of classical nucleation theory (CNT) (1), which postulates the only important variable for describing nucleation and subsequent crystal growth is the size of the nascent crystal (2), more recent studies of nucleation for a number of physical systems under a variety of conditions via high-resolution experiments (3–6) and molecular dynamics (MD) simulations (2, 7–17) have shown that CNT alone is not sufficient for an accurate description of the mechanism of crystal formation, which often involves multiple, competing pathways (3, 18, 19). Thus it is of interest to discover other important variables besides or in addition to the size of the nucleus to accurately describe the nucleation process. One can imagine this situation is especially relevant for substances such as H₂O or organic molecules possessing multiple stable, crystalline structures or polymorphs, where both the lattice structure as well as crystal size are important in distinguishing the solid and the liquid states from each other.

From an experimental perspective, nucleation is a ubiquitous phenomenon, but the microscopic processes driving the nucleation event typically occur over time- and lengthscales that are, generally speaking, below the resolution of current physical apparatuses (3, 4). Thus, applying a combination of theory and simulation to describe the microscopic nature of nucleation is required to gain a detailed understanding of the molecular nature of the nucleation process and build a bridge to the time- and lengthscales accessible by computational and experimental approaches.

Even with supercomputers capable of simulating hundreds of microseconds per day (20), the timescale for sampling most nucleation processes, occurring on the timescale of seconds, still remains out-of-reach. Since nucleation is a rare event from the microscopic perspective, nucleation events can generally only be observed in unbiased simulations under conditions of heavy supersaturation or by applying an external driving force. Thus, finding excellent reaction coordinates (RCs) along which the system can be biased to accelerate the sampling is still very much of importance through a variety of enhanced sampling methods (21–26). In parallel to the development of more powerful computing machines and sampling algorithms, there has also been an explosion in the use of machine learning (ML) techniques to sift and interpret the results of long simulations (27–34). Neural networks and ML in general have been implemented to discover good RCs for describing nucleation in a more automatic manner. These neural network derived RCs can then be used as the biasing coordinates in different enhanced sampling methods. In fact, even traditionally RC-free methods, such as path sampling methods (e.g. forward flux sampling) have been shown to benefit from a judicious choice of low-dimensional projections along which to calculate the flux (35, 36).

In principle, the RC can be directly expressed as a function of simple variables such as high-dimensional atomic positions. But, for various computational reasons, generally the RC is expressed as a linear or non-linear function of lower-

Significance Statement

Although a common occurrence experimentally, in many instances, a detailed picture of crystal nucleation remains unresolved. Computer simulations are helpful for gaining molecular insights regarding nucleation, but typically unable to reach experimental timescales where nucleation events are readily observed. To overcome this limitation, we use a machine learning augmented molecular dynamics approach that automatically learns and biases along optimized reaction coordinates to enhance sampling of nucleation. We apply this procedure to aqueous solutions of urea and glycine, sampling multiple polymorphs for both and estimating their solvated free energies relative to their aqueous solvated states. Our reaction coordinates provide evidence of nonclassical nucleation mechanisms for both systems. This protocol should be to all-atom resolution studies of nucleation in different environments.

ZZ, EB, ST, PT designed research. ZZ, EB performed research. ZZ, EB, PT wrote the manuscript. There are no competing interests to declare.

[†] Ziyue Zou contributed equally to this work with Eric R. Beyerle.

¹ To whom correspondence should be addressed. E-mail: ptiwary@umd.edu

dimensional order parameters (OPs), which are able to distinguish metastable states. These OPs themselves are generally *a priori* designed functions of the atomic positions. Approximating the true RC with smoothly differentiable functions of the OPs has become one of the main aspects of studying crystal nucleation using enhanced sampling approaches such as well-tempered metadynamics (WTmetaD) (37), which is our interest in this work, but also more generally for other sampling approaches (13, 14, 34, 38–55).

In this manuscript, we utilize a neural network framework termed the state predictive information bottleneck (SPIB) (56, 57) to extract a two-dimensional RC for describing the nucleation of different polymorphs of urea and glycine from solution. SPIB is a variant of RAVE that discovers a set of low-dimensional RCs capable of faithfully predicting the metastable states given the input OPs, which can then be biased in enhanced sampling simulations (57). The input OPs originate from a library of candidate OPs and are system specific (56, 57). For the nucleation processes studied here, the OPs are derived from a set reporting on the global and local molecular orientations and packing of the nucleating species, e.g. coordination numbers (14, 58), radial distribution functions (59, 60) and Steinhardt bond OPs (61), *inter alia*. We describe SPIB in more detail in [State Predictive Information Bottleneck](#).

We find that biasing MD simulations of urea and glycine in water using the RCs discovered by SPIB allows for a sufficient enhancement of the sampling such that multiple back-and-forth liquid to solvated crystal polymorph transitions are observed over practical compute times on high-performance resources. Furthermore, through the use of a linear architecture during the encoding of the RCs, we are able to directly interpret the RCs discovered. Finally, by reweighting these biased simulations (62) the relative stability among the solvated crystalline phases sampled is ranked for both urea and glycine, and these rankings are compared with previous, related studies.

Discovering and Biasing Reaction Coordinates for Nucleation

Order parameters for nucleation. As mentioned in the Introduction, in this work as well as in work by others (13, 14, 38–51), it is common to build the reaction coordinate as a function of OPs that collectively distinguish among competing metastable states, which has led to the development of several rich OPs for the study of nucleation and phase transitions in general:

1. *Coordination numbers and associated moments:* As introduced above, CNT is based on the size of the nucleus, and the formation of such crystal nuclei is closely related to the formation of a denser phase. A continuous and differentiable coordination number of any molecule i can be defined as:

$$c(i) = \sum_j \frac{1 - (r_{ij}/r_c)^6}{1 - (r_{ij}/r_c)^{12}} \quad [1]$$

where r_c is a radial cutoff and r_{ij} denotes the distance between reference sites for two molecules i and j . When averaged over all molecules in the system, the coordination number can serve as an approximation to the true reaction coordinate in the study of gas-liquid transitions (14, 63, 64).

Inspired by this definition, in this manuscript we specifically consider two sets of populations of molecules, with coordination numbers greater than 8 and 11, denoted as N_{8+} , N_{11+} , respectively. We also consider the second moment of all coordination numbers as an OP, denoted μ_c^2 (14, 58, 63). However, such coordination number based descriptors are not sufficient to describe the reaction coordinate in events like nucleation from the melt and polymorphism in molecular systems because nucleation is accompanied with the appearance of a space group. Therefore, in addition to the OPs N_{8+} , N_{11+} and μ_c^2 derived from coordination number, we include other OPs into the dictionary, as introduced next.

2. *Steinhardt bond OPs:* These OPs, originally introduced in Ref. 61, have been used to distinguish Lennard-Jones solids (65–71) and have been extended to the study of ice polymorphs (72–76). They map a given local environment onto specific degrees of spherical harmonics (61), defined through :

$$q_{lm}(i) = \frac{\sum_j \sigma(r_{ij}) \mathbf{Y}_{lm}(\mathbf{r}_{ij})}{\sum_j \sigma(r_{ij})} \quad [2]$$

where \mathbf{Y}_{lm} is the l^{th} order of spherical harmonics and m ranges from $-l$ to l . As these bond OPs refer to distributions as superpositions of perfect crystalline structure, reliable results have been acquired in classifying face-centered-cubic (fcc), body-centered-cubic (bcc), hexagonal-close-packed (hcp) and liquid-like structures (65). Here, we add \bar{q}_4 , \bar{q}_6 from Ref. 61 corresponding to the average values of the fourth and sixth orders of spherical harmonics to our OP library. However, as pointed out by Dellago and Lechner (77) thermal fluctuations may destroy the ability in phase identification and, as such, the variants of original bond OPs have been developed for better differentiating ability at the expense of a higher computational cost; these are omitted here.

3. *Interfacial water:* The contribution of solvent molecules can hardly be neglected when studying the nucleation of solvated systems (15). Here we introduce a new OP measuring approximately the population of interfacial waters or more generally solvent molecules surrounding solid-like solute molecules by using tunable distance cutoffs. This variable is defined as

$$N_s = \sum_i \frac{1}{2} \xi(i) c_{solvent}(i) \quad [3]$$

where

$$\xi(i) = \tanh(c_{solute}(i) - c_o) + 1 \quad [4]$$

Here i denotes a sum over solute molecules and c_o is a tunable parameter for the determination of solid-like aggregates which is set to be 5.0 for all systems in this work (58). Effectively, the OP in Eq. 3 is a product of two smoothed Heaviside functions: $\xi(i)$ only counts solute clusters larger than the threshold c_o , and $c_{solvent}(i)$ is a coordination number between such selected solute atoms

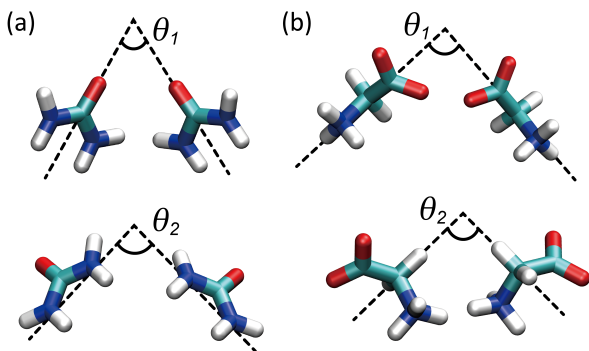


Fig. 1. Visualization of the feature vectors and corresponding intermolecular angles θ_1 and θ_2 for (a) urea and (b) glycine molecules. The feature vectors shown here are the vectors connecting C-O and N-N for urea, and C-N and $C_\alpha-H_\alpha$ for glycine. Different atom types are indicated with different colors, specifically carbons in cyan, oxygens in red, nitrogens in blue, and hydrogens in white. Snapshots are rendered using Visual Molecular Dynamics (VMD) (79).

and the solvent, as given in Eq. 1. Thus, N_s is large when two conditions are simultaneously satisfied: when the solute cluster is sufficiently large and also well-coordinated with solvent. If either one of these conditions is not satisfied, N_s is small. Thus, large values of N_s indicate cluster formation and low values indicate a liquid state or the presence of small, solvent-separated clusters.

4. Averaged intermolecular angles and associated moments:

The estimation of the size of solid-like aggregates and their solvation states is useful and intuitive for the study of molecular systems in different environments (14, 63, 64), but nucleation mechanisms can easily consist of non-classical behavior and polymorphs such that solely distance based OPs fail in categorizing them. In the spirit of Steinhardt bond OPs (61) and environment similarity OPs (78), we introduce a set of OPs which account for the orientation of molecular feature vectors of solutes to help better discriminate polymorphs. For any molecule i , we define:

$$\theta(i) = \frac{\sum_j \sigma(r_{ij}) \frac{1}{2} [(\pi - 2\tilde{\theta}) \tanh(5\tilde{\theta} - 9.25) + \pi]}{\sum_j \sigma(r_{ij})} \quad [5]$$

where $\tilde{\theta}$ is the angle formed between characteristic vector on molecule i and molecule j as illustrated in Fig. 1, and σ is a switching function of intermolecular distance. The hyperbolic tangent switching function is applied to remove the mirror image symmetry (see SI for detailed explanation). We then compute the mean of angular distribution functions, denoted as $\bar{\theta}_1$ and $\bar{\theta}_2$ where 1 and 2 denote the two intramolecular vectors specified for crystalline configuration identification (see Fig. 1 (a) and (b) for definitions of selected vectors for urea and glycine molecules, respectively). In addition to these OPs, we also consider in our OP dictionary the second moments $\mu_{\theta_1}^2$ and $\mu_{\theta_2}^2$ of the collection of θ values defined through Eq. 5.

5. *Pair entropy*: In Ref. 59 Piaggi and co-workers introduced an OP which approximates radial and orientational entropies from radial different distribution functions. This OP is formulated as:

$$S(r) = -2\pi\rho k_B \int_0^\infty [g(r) \ln g(r) - g(r) + 1] r^2 dr \quad [6]$$

where $g(r)$ is the radial distribution function, ρ is the density, and k_B is Boltzmann’s constant. This OP was introduced in order to account for the possible entropy-favored crystalline structures (60). An advanced version of this OP is developed in Ref. 60. This pair orientational entropy OP, denoted $S(r, \theta)$, depends on both the intermolecular distance and the relative angle θ between a given intramolecular vector equivalently defined for each solute molecule in the simulation:

$$S(r, \theta) = -\pi\rho k_B \int_0^\infty \int_0^\pi [g(r, \theta) \ln g(r, \theta) - g(r, \theta) + 1] \times r^2 \sin(\theta) dr d\theta \quad [7]$$

It was found (60) that such an approximation to the entropy is useful for distinguishing polymorphs in simulations of urea and naphthalene, and a variant of Eq. 7 is used in (55) to successfully discover and distinguish polymorphs of 1:1 mixtures of resorcinol and urea. Since Eq. 7 only accounts for a single angle, and not the three angles necessary to specify exactly the relative orientation between molecules, this expansion of the entropy, though useful, is necessarily approximate, as are the other OPs previously described in this section.

While individually these different OPs might have limitations, they can supplement each other and provide a palette through which the complex molecular processes governing nucleation can be described. It is by combining them through the State Predictive Information Bottleneck (SPIB) (56) approach that we construct low-dimensional projections that are much closer to the true RC than any of the OPs on their own. We next briefly summarize the SPIB approach, referring to Ref. 56 for further details.

State Predictive Information Bottleneck (SPIB). Given the evidence that CNT is insufficient to accurately describe the nucleation event for a number of chemical systems, for instance because knowledge of more than just the solute coordination number OP is required, we need an approach to combine the different possible OPs to develop a lower-dimensional reaction coordinate for describing nucleation in urea and glycine.

Here we utilize a machine learning method known as the state predictive information bottleneck (SPIB) (56, 57, 81), which takes the form of a variational autoencoder (VAE) (82, 83). SPIB differs from the traditional VAE because instead of attempting to maximize the model’s ability to re-construct the input from the low-dimensional latent space, it instead maximizes the quality of reproducing the population of the system’s metastable states after a given lagtime τ in the future. That is, the SPIB finds the latent space that best reproduces

Table 1. Parameters for WTmetaD

System	ω (kJ/mol)	γ	σ_1 (RC unit)	σ_2 (RC unit)	T (K)
Urea	5.0	100	0.2	0.2	300
Glycine	5.0	100	1.15	1.15	300

the metastable state populations at time $t + \tau$ given the values of the input OPs at time t (56). This approach allows SPIB to simultaneously perform dimensionality reduction and accurate future state prediction by minimizing a loss function that is inspired by the variational information bottleneck formalism (56, 84, 85).

Conceptually, this approach is similar to performing a Markov state model and clustering analysis (86, 87) with a continuous basis set. In SPIB the number and neighborhood of each metastable state is adjusted on-the-fly during the training of the model. The final learned model is output once the neighborhoods spanned by each metastable state have converged to below a pre-defined threshold that is a tunable hyperparameter of the model.

Finally, since biased simulations are input to the SPIB, normally the WTmetaD weights would be used in the analysis to account for the sampling from a biased distribution (57). For urea and glycine, we find that the barriers to nucleation from the liquid state are so large (~ 100 kJ/mol or more) that using the metadynamics weights in the analysis effectively destroys any polymorph minima on the surface, yielding a single, global minimum corresponding to the liquid state (more details in the SI). Since SPIB finds the RC describing transitions between metastable minima, the metadynamics weights must be neglected in this case to find an effective set of RCs for nucleation.

Estimating free energies with Metadynamics. In order to learn the SPIB, or the approximate RC, as a combination of different OPs, we need access to a trajectory that has visited different possible conformations. Such a trajectory is generated here through the method Metadynamics, which helps the system escape free energy minima by periodically depositing Gaussian kernels as a function of the variable being biased. For the initial metadynamics runs we bias orientationally informative OPs such as intermolecular angles θ and orientational entropies S_θ . After performing SPIB analysis on these runs, we then bias

along the SPIB-learned RCs in future iterations. In particular, here WTmetaD is used where kernels are decreased in height for biasing regions being revisited, thus leading to better convergence of the free energy surface (37). The WTmetaD parameters are reported in Table 1. The height of Gaussian deposition, ω , is selected to be $2 k_B T$ at 300 K and the bias factor, γ , is set to be 100, in order to overcome high energy barrier associated with nucleation problems (88, 89). The width of the Gaussian, σ , is set to the thermal fluctuation of corresponding approximate reaction coordinate estimated from a short (~ 10 ns) unbiased MD run.

After classifying frames in the biased simulation into either the liquid state or a particular crystal polymorph, the free energy difference between two states a and b at temperature T can be calculated as follows:

$$\Delta G_{a \rightarrow b} = G_b - G_a = k_B T \log \frac{P_a}{P_b} \quad [8]$$

where k_B is Boltzmann’s constant and P_a, P_b denote the Boltzmann probabilities of states a, b obtained from reweighting the WTmetaD simulations (62).

Results and Discussion

We performed classical all-atom MD simulations for all systems using GROMACS version-2022.2 (90) patched with PLUMED 2.6.1 (91, 92) to perform the WTmetaD. Further simulation details are provided in **Materials and Methods**. Below we provide detailed results first for urea followed by glycine.

Urea.

Urea Nucleation. Urea is commonly used as fertilizer and as nitrogen feedstock for organic synthesis. Due to its industrial importance, urea has been studied extensively and several crystalline structures have been synthesized and reported. These include polymorphs named I (space group: $P42_1m$), III ($P2_12_12_1$), IV ($P2_12_12$) and V ($Pmcn$), where polymorph

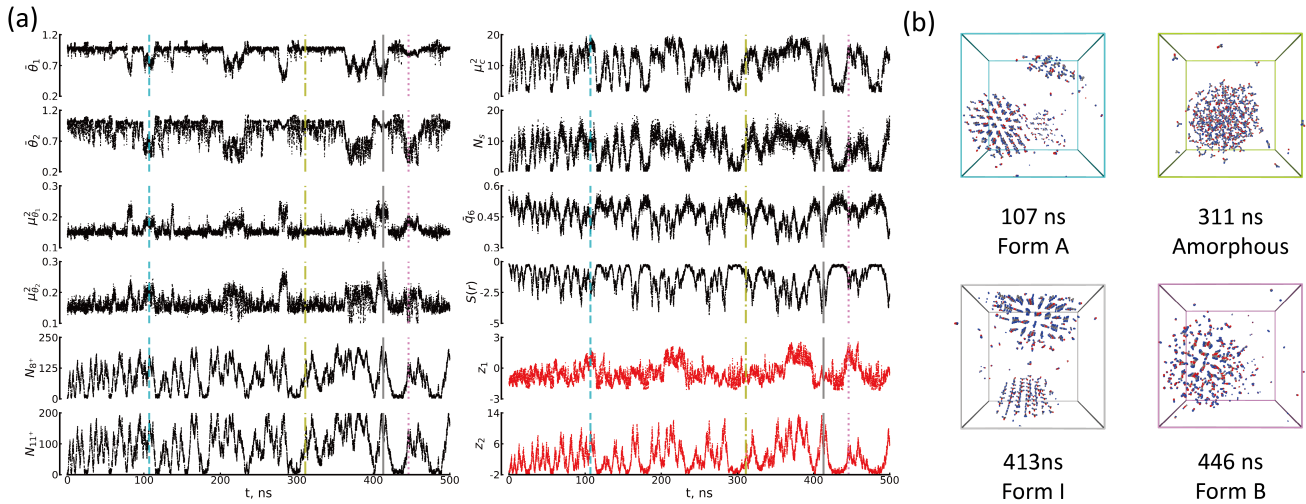


Fig. 2. Sampling different urea polymorphs with WTmetaD simulations biasing the SPIB-learned 2-dimensional reaction coordinate. (a) shows a representative time series of OPs (black) and reaction coordinates (red), clearly demonstrating multiple back-and-forth transitions between different phases. Vertical cyan dashed, green dash-dotted, grey solid, and pink dotted lines indicate representative transitions to form A, amorphous, form I, and form B from initial liquid phase. (b) shows snapshots of structures sampled during trajectory shown in (a), along with the time (vertical lines in (a)) at which they were observed. Ovito 3.3.5 (80) was the visualization tool used for snapshot generation.

I corresponds to the typical isoform stabilized at ambient conditions and the rest are high-pressure, high-temperature products (93–99). Moreover, an ever-enriched dictionary of polymorphism in urea has been established with computational methods using classical and quantum approaches (15, 60, 88, 89, 100–103). We adopt the same notations as in Ref. 60, where two more polymorphs of urea A (*Pnma*) and B (*P1*) are introduced. It remains unknown if these two crystal structures A and B are artifacts due to use of the generalized Amber force field (GAFF) (104) or simply experimentally yet to be observed. Despite the orientational difference in the lattice cell, hydrogen bonding pattern in both forms A and B (all in type III) has been shown to be different from experimental forms (type I and type II) (102, 103).

Chronologically, we start with a preliminary round of well-tempered metadynamics to facilitate state-to-state transitions along trial reaction coordinates. For this, we biased the two averages of the two feature angles (Fig. 1 (a)), $\bar{\theta}_1$ and $\bar{\theta}_2$ as they have been considered as OPs in distinguishing crystal structures of urea from each other (46). Several transitions to and from the solid states of interest are observed in the WTmetaD simulations (data not shown) biasing this two-dimensional coordinate (4 independent runs of 400 ns each, totalling $\sim 1.5 \mu\text{s}$) and trained as input data for SPIB along with other candidate OPs.

Urea Polymorphism. From this preliminary WTmetaD simulation, we performed SPIB to learn a 2-d latent representation for the reaction coordinate (denoted as z_1 and z_2) as linear combinations of OPs discussed in **Materials and Methods**. We then perform a second round of WTmetaD biasing along this SPIB learned 2-d RC consisting of 4 independent runs of 500 ns each. The time series for different OPs and the 2-d RC components are shown in Fig. 2(a) in black and red, respectively. This unequivocally shows the presence of numerous back and forth state-to-state transitions (see SI for detailed sampling efficiency comparison), a hallmark of good enhanced sampling (105). In particular, three crystal structures, namely polymorphs I, A, and B, are visited during simulations. The corresponding snapshots for these from one are shown in Fig. 2(b) along with that of an amorphous bulk phase. The evidence for the formation of these polymorphs can be found in the time series of intermolecular angle OPs (top left two rows in Fig. 2(a)), as either one or both of these OPs drop, indicating formation of an ordered phase based on their definitions. Comparing to the known polymorphs introduced above, structures from experimental products III, IV, V (97) and zero-temperature prediction C (103) are missing. It is possible that the former high-pressure structures are less stable as the existence of mediating solvent molecules, and the latter isoform, could potentially be unstabilized by addition of entropic effects under finite temperature. As these solid states have not been found in previous studies using MD simulations (88, 100) (except form IV being observed in nucleation from the melt (59)), we therefore believe there is remaining scope here for exploration of the configuration space in future studies from other perspectives such as force-field refinement.

In Fig. 3(a), the reweighted free energy surface along the two-dimensional latent representation z_1 and z_2 is plotted (62). The states of interest, one liquid and three solid phases, have been labeled on this surface. However, no minima are observed corresponding to regions marked as polymorph states

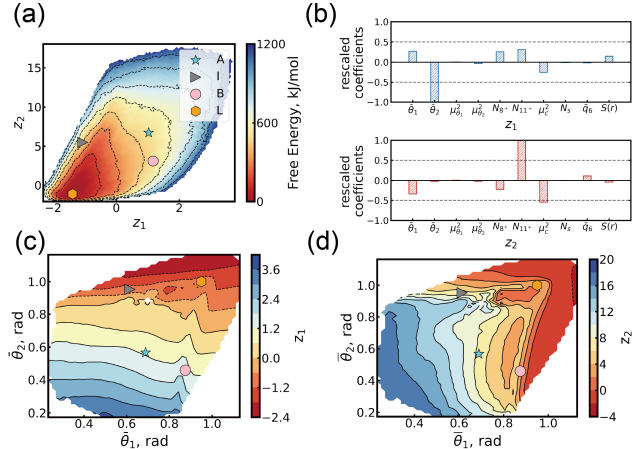


Fig. 3. (a) Reweighted free energy surface in kJ/mol along SPIB-learned reaction coordinates with landmarks for states of interest (solid states A, I, B, and liquid state L). (b) Rescaled coefficients for each OP in construction of SPIB-learned RCs. The coefficients are rescaled with respect to the corresponding fluctuations of components in WTmetaD simulations (see **Materials and Methods** for details). The importance of the urea molecules relative orientation can be seen from the first reaction coordinate while the second coordinate shows the relevance of clustering. (c) z_1 -projection onto $(\bar{\theta}_1, \bar{\theta}_2)$ space; and d) z_2 -projection onto $(\bar{\theta}_1, \bar{\theta}_2)$ space. The markers asterisks (cyan), triangles (grey), circles (pink) and hexagons (orange) indicate crystalline structures A, I, B and liquid phase L, respectively.

of urea, which could be caused by multiple factors. One of the main reasons suggested is the effect of finite size, which can be corrected analytically using CNT or increasing the size of the simulation box (88, 89). However, in this work, we emphasize the construction of ML RCs that explore the configuration space reported in previous studies rather than conveying detailed, more quantitative, investigations on thermodynamic and kinetic properties, which are left for future studies.

In Fig. 3(b), we provide the coefficients for different OPs, indicating how much they contribute to the RC (see **Materials and Methods** for details on calculating these coefficients). From these bar plots, it can be seen that several OPs are weighted heavily in the RC, comprising attributes of intermolecular angle and coordination number. Referring to the definition of these two categories of OPs, the coordination number is a simple but useful OP in the identification of phases, suggestive of the CNT formalism where only size of the nucleating cluster is sufficient. However Fig. 3(b) clearly shows that the coordination number in such complicated molecular systems is not enough for capturing the slow degrees of freedom for crystallization. In addition, our RC indicates that the orientation of feature vectors needs to be taken into account as these vectors are significant when classifying polymorphic urea (102).

To support the above rationale, we project the SPIB-learned RCs z_1 and z_2 used to bias onto orientational informative OP space $(\bar{\theta}_1, \bar{\theta}_2)$ in Fig. 3(c) and (d) respectively. Here it can be seen that the driving forces behind the two representations are clearly different. z_1 specifically enforces the transitions from initial liquid state (orange hexagon) to states with lower $\bar{\theta}_2$ values (which are form A (cyan star) and form B (pink circle)) and it shows strong correlation to $\bar{\theta}_2$ OP itself. On the other hand, z_2 pushes the system away from liquid state to all potential solid phases.

Comparison with previous studies. For the purpose of accurate classification of crystal structures for further analyses, the local crystallinity (SMAC) OPs are used (see SI for full expressions) (15, 88, 89, 100). In this way, polymorphs can be screened precisely by properly positioning the centers of various switching functions. The size of the nucleus can then be computed using graph theoretic algorithms (106). After this classification, the free energy corresponding to each phase can be evaluated, giving the free energy difference between two arbitrary states a and b at temperature T as described in **Estimating free energies with Metadynamics**. Fig. 4 shows the free energy differences of different forms of solvated, crystalline urea with respect to the initial solvated liquid phase obtained after reweighting the biased WTmetaD simulations (62).

Under the consideration of the existing finite-size effect, we draw only qualitative analyses on these visited metastable states in terms of the relative stability of each crystalline structure in aqueous solution. As can be seen from this figure (Fig. 4), our simulations suggest that form A is the most stable polymorph in the given water model, followed by form I and finally form B. A similar relationship between forms A and I in aqueous solution has been established and studied in detail in previous work: whereas a small nucleus favors crystal structure A, interconversions between A and I occur when cluster grows larger. In particular, Ref. 88 suggests a barrierless A-to-I transition occurs as clusters larger than ~ 50 molecules when applying corrections to leverage finite-size effect, and Ref. 101 claims such a transition occur at a cluster size of ~ 530 molecules associated to a $\sim 125 k_B T$ energy barrier when performing a seeding method by inserting pre-built crystals into large simulation boxes.

To the best of our knowledge, metastable state B has never been synthesized in aqueous solution of urea, and its very high free energy relative to the other forms may explain why such a structure is difficult to sample. We can explain why we were able to sample this metastable polymorph in our simulations and also in Ref. 60. Structurally, form B has a completely different orientation among other polymorphs in which its dipole moment in the direction of C-O vector is no longer parallel or anti-parallel to its neighbors. This is likely a result of entropic contributions (60). By having entropy-related OPs in our dictionary of OPs that comprise the RC, we are able to accelerate the substantial entropic degrees of freedom together with other more energy dominated modes, and therefore both enthalpic- and entropic-favored crystalline structures of urea are formed with enhanced sampling method biasing along machine learned RC. We then demonstrate the ability of the SPIB approach on the nucleation of a more complex molecule, glycine, in aqueous solution.

Glycine.

Glycine Nucleation. Glycine is the simplest amino acid, with its R-group consisting of a single hydrogen atom (107). Glycine has multiple physiological functions: it serves as a precursor to more complex biomolecules such as heme, purines, creatine, and more complex amino acids; glycine is also an inhibitory neurotransmitter (108) and is involved in cytoprotection and immune system function (109). Glycine forms three well-known polymorphs from solution called the α (space group: $P2_1/n$), β ($P2_1$), and γ ($P3_1$) polymorphs; thermodynamically,

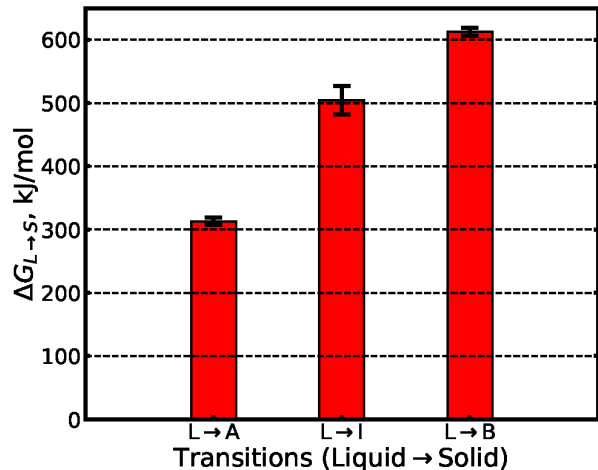


Fig. 4. Free energy difference in kJ/mol between different solvated polymorphs with respect to liquid state as calculated from WTmetaD simulations. This indicates that form A is more stable than form I and form B is the least stable among the three solid phases under solvated condition, comparatively. The error bar is computed over four independent product runs.

γ -glycine is the the most stable, but α -glycine tends to be the form that crystallizes from solution at neutral pH and ambient atmospheric pressure (110–113). Several other high-pressure polymorphs have been discovered (113, 114), but they are not typically observed at atmospheric pressures or temperatures. For performing the WTmetaD (37) simulations of glycine, we use a similar OP library as urea, with one major exception: the pair entropy $S(r)$ is swapped for the orientational pair entropy $S(r, \theta)$ (60). The angles used for calculating $S(r, \theta)$ are those given in Fig. 1(b). For glycine, the first intramolecular vector is defined as that from the N-terminal nitrogen to the C-terminal carbon atom; the second intramolecular vector starts at the alpha-carbon C_α , and ends at the R-group hydrogen atom H_α .

In contrast to urea, we find it necessary to include $S(r, \theta)$ in the OP library because, without it, reversible nucleation events are not observed (data not shown). We postulate that this increased importance in the case of glycine is due simply to $S(r, \theta)$ being a two-body term and, hence, a first-order correction to the excess entropy of a fluid (115, 116) compared to the entropy of the equivalent ideal gas. Since glycine is more massive and zwitterionic, as well as possessing a stronger dipole moment, it is less ideal in the liquid phase compared to urea and, therefore, it is to be expected that it will possess more excess entropy compared to urea, making it necessary to place $S(r, \theta)$ and not just $S(r)$ in the OP library.

For the preliminary round of WTmetaD, we bias along $S(r, \theta_1)$ and $S(r, \theta_2)$ since these OPs were found to give the best sampling from the OP library. This 400ns preliminary simulation and its SPIB analysis is described in detail in the SI. As with urea, the trajectories of the OPs from this preliminary simulation are used as input data for training SPIB. Also in the SI are projections of the metadynamics bias onto the $(\bar{\theta}_1, \bar{\theta}_2)$ space, where it can be seen that using the SPIB RCs greatly enhances the amount of bias deposited, and hence the amount of sampling, of regions of $(\bar{\theta}_1, \bar{\theta}_2)$ space where the glycine polymorphs of interest reside. This result provides support for using the SPIB RCs as the biasing coordinates

over the hand-picked orientational entropies.

Glycine Polymorphism. Here we are interested in the α -, β -, and γ -glycine polymorphs observed at ambient temperature and pressure (113). A necessary but not sufficient method (117) to differentiate the three polymorphs is through the collective orientation of the C-N vectors (Fig. 1(b)). For α -glycine, C-N vectors alternate in layers of two in a ‘positive’ orientation (C-N) followed by a ‘negative’ orientation (N-C) repeating. For β -glycine, positive and negative layers alternate one at a time. In γ -glycine, all layers of the crystal possess the same orientation, either all positive or all negative. Using only the C-N axis orientations to classify polymorphs, we find a relative stability rank order of form- $\gamma > \text{form-}\beta \geq \text{form-}\alpha$, which recovers form- γ as the most stable form in agreement with experiments of bulk glycine (111, 113).

While this classification approach correctly accounts for the large-scale orientation of the crystal, it ignores the local orientation, such as hydrogen bond patterns and relative angles and distances between layers of the crystal (113, 118). This coarse procedure for identifying polymorphs is selected over the more rigorous one given in Ref. 119, which has not been tested for use in nucleation from solution or long timescale simulations, where there can be significant perturbations to the crystal structure. Based on these considerations, we have chosen to characterize the polymorphs using the relative orientation of monomers along the axis parallel to the C-N vector only. Ignoring local structure, this protocol will introduce some contamination of the polymorph classification. Furthermore, the crystal structures of glycine are much more similar to each other compared to urea, making the SMAC protocol used for urea ineffective for glycine, with details for why this is the case given in the SI. Since this study is, to our knowledge, the first to perform enhanced sampling simulations of glycine to study nucleation without seeding in water alone at or below the saturation concentration, our reported numbers at the very least have qualitative significance in being able to correctly rank order glycine polymorph stability in pure aqueous solution.

Fig. 5a shows the trajectory of the OPs used in the SPIB analysis for glycine, and, in the final panel, the trajectories of the two SPIB RCs; all trajectories come from the final set of simulations biased along the discovered linear SPIB RCs. These trajectories clearly show many transitions between the solvated, crystalline states and the liquid state of glycine; a comparison of the sampling efficiency of this set of trajectories biased along the SPIB RCs and the hand-picked orientational entropy CVs is given in the SI for both glycine and urea.

Fig. 5 shows the α -glycine, β -glycine and γ -glycine polymorphs extracted from the SPIB-biased trajectories. For all three snapshots, the appropriate orientations for each polymorph along the monomer C-N axis should be apparent, although, given the size of the system (72 glycine molecules), finite volume and surface area effects are likely to be significant due to the large surface area to volume ratio of the clusters. We also expect solvent-induced perturbations compared to the respective crystal structures (120, 121) from the melt.

An analysis of the SPIB RCs used to bias the nucleation simulation of glycine is given in Fig. 6. In Fig. 6(a), the reweighted free energy surface along the 2-d linear SPIB latent space used as the biasing variables in metadynamics is shown with the locations of the putative α -, β -, and γ -glycine poly-

morphs labeled with colored stars; for reference, the liquid (isotropic) state is also labeled with a magenta star. Fig. 6(a) shows that, while none of the glycine polymorphs are located in free energy minima they are neither at free energy maxima, which indicates that the SPIB RCs are doing a reasonable job sampling the crystal polymorphs.

Fig. 6(b) shows the rescaled contribution of each input variable to the final two-dimensional SPIB model. For both RCs, the orientational entropies are weighted heavily in-line with the physical intuition explained in [Glycine nucleation](#) that higher-order corrections to the entropy of the zwitterionic glycine will be more important than an uncharged molecule such as urea. Also highly weighted are the two angular coordinates, $\bar{\theta}_1$ and $\bar{\theta}_2$, which justifies performing the polymorph labeling and analysis in that subspace.

The SPIB RCs are further interpreted by projecting them onto the $(\bar{\theta}_1, \bar{\theta}_2)$ subspace in Figs. 6(c), (d). The first SPIB RC, z_1 , corresponds to nucleation from the liquid to the three common polymorphs labeled on the surface. Examining the contours more closely, it specifically describes transformation from the liquid and γ -glycine states, which are the two most populated states observed in our simulations, to the α - and β -glycine polymorphs. There is also a strong correlation of $\bar{\theta}_1$ with z_1 , as expected from Fig. 6(b). The second SPIB RC, z_2 , corresponds to transitions from the liquid state to γ -glycine via the α and β polymorphs. Interestingly, this mechanism of transition from the liquid state to the most stable crystal polymorph (form γ) via the less stable crystal polymorphs (forms α and β , Figure 7) is consistent with the Ostwald step rule (122) and is evidence for two-step nucleation in aqueous glycine, which already has been observed experimentally (123). This RC also correlates strongly with transitions from low-to-high values of $S(r, \theta_1)$ (data shown in the SI), again as expected from Fig. 6(b).

In summary, we find the linear RCs discovered by SPIB are able to 1) separate the common glycine polymorphs in the two-dimensional RC space and 2) yield a reasonable and interpretable set of reaction coordinates for describing nucleation of glycine from aqueous solution.

Comparison with previous studies. While the SMAC OP is used to classify urea polymorphs, we find that the $g(r, \theta)$ distributions for the polymorphs of glycine to be too similar for SMAC to be effective in distinguishing polymorphs and instead using the classification method outlined in [Glycine polymorphism](#). We also propose a second method based on Voronoi tessellation of the $(\bar{\theta}_1, \bar{\theta}_2)$ space around the ‘landmark’ polymorph and liquid structures shown in Figure 6 in the SI.

To our knowledge, no explicit free energy differences between polymorphs for crystallization from aqueous solution for glycine have been reported in the literature. Solubility experiments in pure water and water-organic solvent-antisolvent mixtures of varying composition have been reported in Ref. 124, where the authors show that, at each solvent composition (ranging from 80% water/20% methanol to 0% water/100% methanol) the solubility order is $\beta > \alpha > \gamma$. Using the argument that the solubility of the crystal in solvent decreases as the free energy barrier increases (since the solubility is roughly the ratio of the concentration of the solvated liquid form to the solvated, crystal form) (88, 125), this solubility order implies a relative stability order with reference to the solvated liquid state of $\gamma > \alpha > \beta$, which is in rough agreement with the

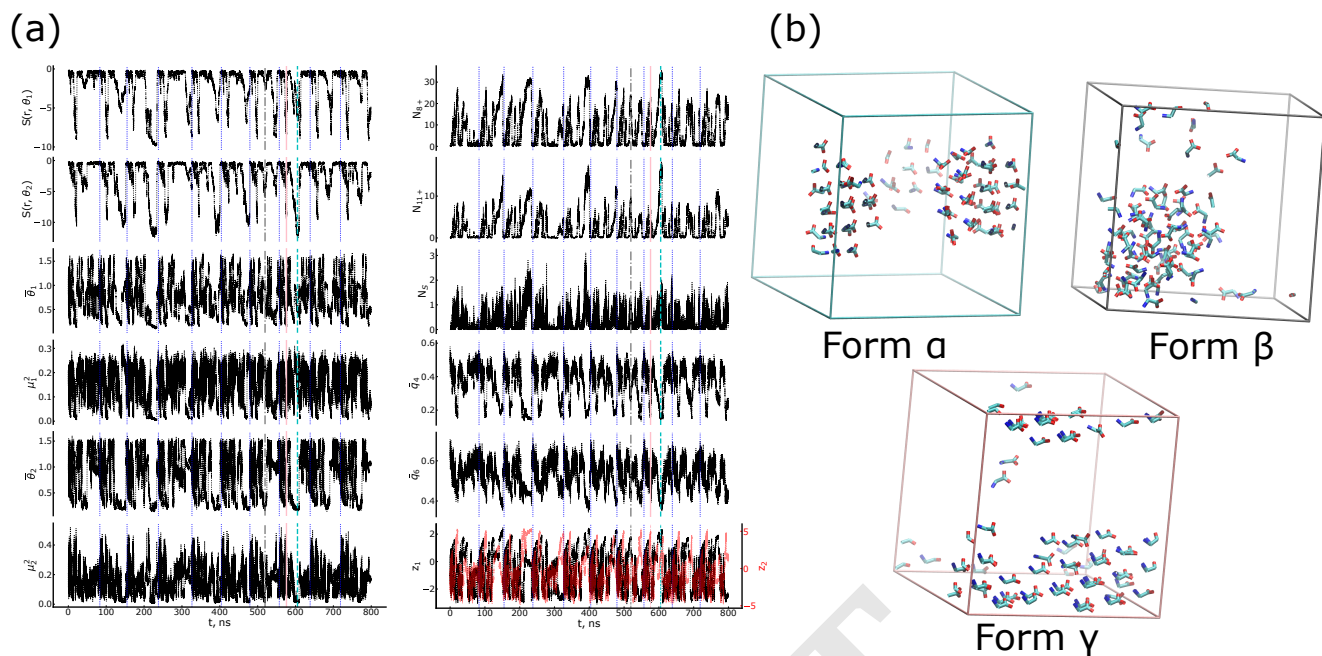


Fig. 5. Sampling glycine polymorphs with WTmetaD simulations biasing the SPIB-learned 2-dimensional reaction coordinate. (a) shows the time series of the OPs (black) and SPIB RCs (black, z_1 : red, z_2) from the ten replicate simulations of glycine, with each trajectory demarcated by the dotted blue lines. Vertical cyan dashed, gray dashed-dot, and pink solid lines indicate representative transitions to form α , form β , and form γ from the liquid phase. (b) shows snapshots of the polymorph structures sampled during trajectories shown in (a); color of the box framing each snapshot corresponds to the color of the vertical time slice in (a). Hydrogen atoms on glycine and water molecules are removed for clarity; snapshots are rendered using VMD.

relative stability we find, namely $\gamma > \beta \geq \alpha$.

We additionally note that the only study performed to rank the relative stability of glycine polymorphs is Ref. 125, where the authors find that the γ -form of glycine is the least kinetically accessible and the β -form the most kinetically accessible, at the nanoscale, using calculations based on seeded cluster MD simulations and CNT. Since we use a vastly different setup (classical MD coupled with WTmetaD and machine learning versus seeded MD and CNT), we are not concerned about disagreements between polymorph accessibility in our study and the one performed in Ref. 125. However, we concede that our glycine clusters are well below the critical nucleus size predicted by CNT and thus may be plagued by finite size effects (126, 127), especially considering that many times the γ polymorph is observed to form in the simulations, the crystal spans the simulation box along one dimension, leading to an infinite crystal when periodic boundary conditions are taken into account (see Fig. 5); this effect could be a primary reason we observe that polymorph to be sampled the most frequently. Additional study is needed to determine the extent to which finite size effects affect the relative stabilities of glycine polymorphs reported here.

We conclude this section by noting that while glycine is a common molecule for experimental study of nucleation and crystal polymorphism, it is also a notoriously difficult system to obtain consistent crystallization results experimentally (111, 113). For example, although most experiments report that the α polymorph is the primary crystal product extracted from solution (112, 123, 128–130), the γ form is observed to form spontaneously and nearly exclusively from solution given the appropriate conditions (113, 130–134). Furthermore, the timescales of most experimental studies of glycine nucleation

tend to be orders of magnitude longer (hours to days) than the MD simulations studied here, making comparisons to experimental results difficult (e.g. the kinetic product seen in those experiments could be the thermodynamics product observed in this study). However, regardless of any quantitative comparison to experiments, we are encouraged by the results presented, as we are able to sample the three well-known experimental polymorphs of glycine using the SPIB approach.

Conclusion

Although nucleation is a frequent event on the experimental and biological timescales, with currently accessible compute resources nucleation is a rare event computationally, even when utilizing enhanced sampling techniques (12, 14, 38, 46, 63, 100, 135). Thus, even for relatively simple systems such as the aqueous urea and glycine systems studied here, sampling polymorph configurations and, more dauntingly, determining relative thermodynamic stabilities of these polymorphs, is difficult computationally. In this work we have demonstrated a more-or-less automatable protocol to obtain back-and-forth sampling of different polymorphs in aqueous solution. Specifically, we have shown that good sampling of polymorphs is possible when classical molecular dynamics simulations of urea and glycine are biased along optimized reaction coordinates found using the machine learning approach State Predictive Information Bottleneck (SPIB) (56) which belongs to the RAVE family of methods (27, 136). A linear, two-dimensional approximation to the reaction coordinate extracted from SPIB (56) can be used to 1) effectively bias nucleation simulations of aqueous urea and glycine to enhance

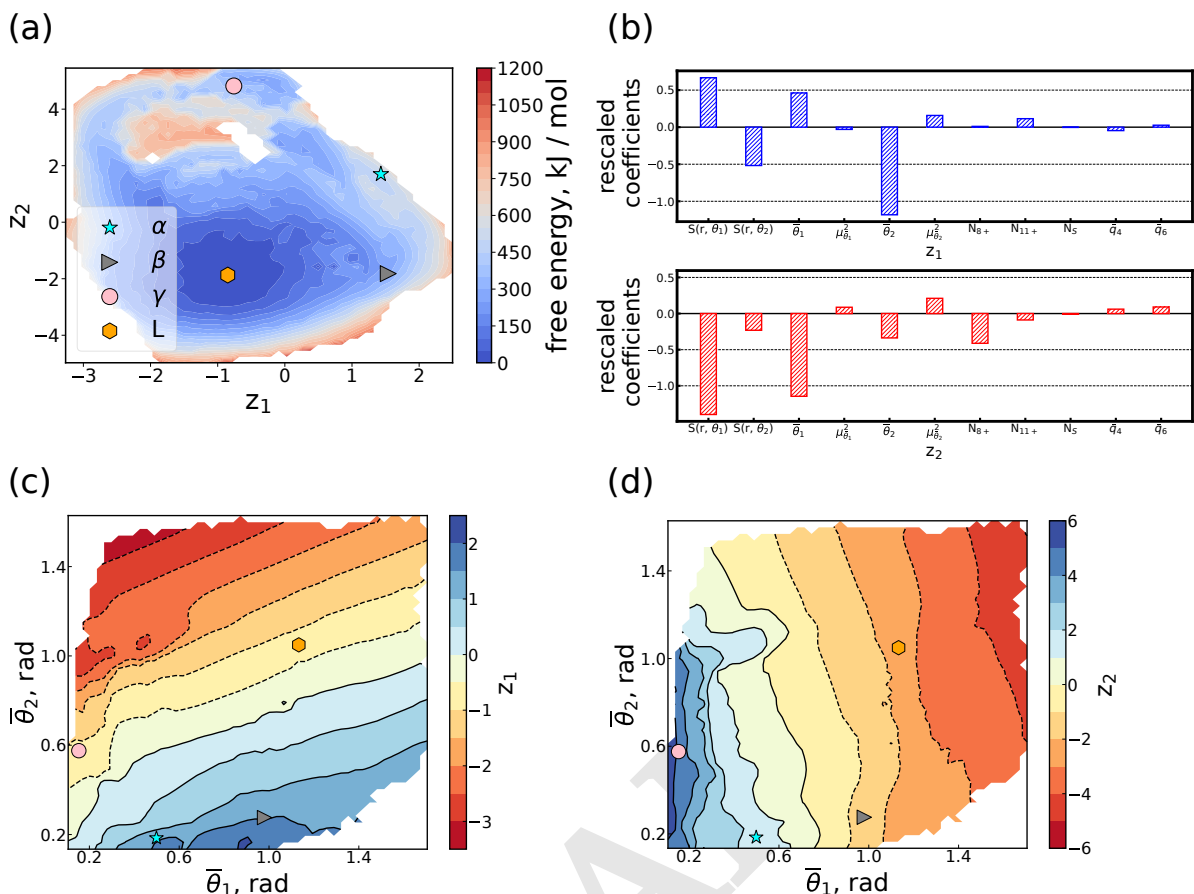


Fig. 6. (a) Reweighted free energy surface in kJ/mol in the two-dimensional space of the linear SPIB RCs, which were used as the metadynamics biasing variables. Markers show the locations of the putative glycine polymorphs sampled during the simulation: liquid, orange hexagon; α , cyan star; β , grey sideways triangle; and γ , pink circle. (b) SPIB coefficients from the linear model, reweighted by the standard deviation of the input variable. For the first SPIB RC (top panel, blue bars), the most important input variables is the trajectory of the average second intermolecular angle $\bar{\theta}_2$, and the orientational entropy $S(r, \theta_1)$. For the second SPIB RC (bottom panel, red bars), the two most important input variables are $\bar{\theta}_1$ and its orientational entropy $S(r, \theta_1)$. (c) The first SPIB RC z_1 projected into the $(\bar{\theta}_1, \bar{\theta}_2)$ subspace. Colored symbols show the locations of the putative glycine polymorphs, with the coloring identical to that in panel (a). (d) is the same as (c), except the projection is of the second SPIB RC, z_2 .

the sampling of polymorph nucleation, and 2) interpret the SPIB RCs to determine which OPs are the most important drivers of nucleation.

When interpreting the RCs derived from the SPIB analysis, we find that for urea, the intermolecular angle and cluster size OPs are the most important contributors to the RC; for glycine, the intermolecular angles and orientational entropy of these angles are the most important to determining the reaction coordinate. These observations reflect that size informative OPs suggested by CNT are insufficient in describing slow modes towards nucleation of urea and glycine molecules; in addition orientational OPs are needed. Each polymorph is identified using different metrics and then ranked appropriately as per their Boltzmann weights at 300 K (62). The relative stability of the solvated crystal polymorphs compared to the solvated liquid state found by us shows form-A > form-I > form-B for urea and form- γ > form- β \geq form- α for glycine. Since we performed the SPIB analysis on the unweighted energy surface, we cannot yet offer any kinetic insights regarding the nucleation of either organic molecule studied here. In addition we did not consider the question of finite size effects, which would especially impact kinetic properties. This will be done in future work using approaches such as Ref. 63.

Secondly, our simulations would have been more representative of supersaturation levels of experiments with the use of constant chemical potential ensemble, which is an active area of research for studying nucleation (135).

To summarize, this work makes several significant contributions to the field of nucleation research and enhanced sampling in general. We have applied deep learning based RC construction methods to arguably the most complex set of systems to date and used them to successfully find reaction coordinates that accelerate the nucleation process of multiple polymorphs for both urea and glycine aqueous solutions. Furthermore, the simulations of glycine reported here are the first, to our knowledge, to utilize metadynamics to induce nucleation of the amino acid in water. Based on these results, we believe using machine learning methods to construct approximate reaction coordinates can provide novel solutions and information regarding non-classical effects in the nucleation of molecules from aqueous media.

Associated Content

Supporting Information. The Supporting Information is available free of charge at xxx. It contains detailed further numerical analyses for different systems.

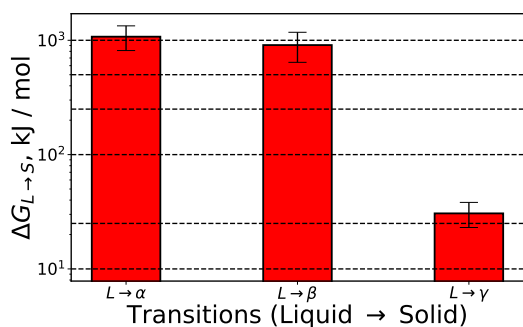


Fig. 7. Free energy difference in kJ/mol between different solvated polymorphs with respect to liquid state as calculated from WTmetaD simulations. This result indicates that form γ is more stable than form β , which is slightly (though not significantly) more stable than form α when solvated with the solvated liquid state as a common reference. Error bars are one standard error of the mean over the ten replicate simulations.

Notes. The code needed to reproduce the models used in this work is available at <https://github.com/tiwarylab/Driving-and-characterizing-nucleation-of-urea-and-glycine-polymorphs-in-water>. The input files necessary to reproduce the simulations done in this work are available on PLUMED NEST at <https://www.plumed-nest.org/eggs/22/039/>.

Materials and Methods

Simulations setup. A 5.0 nm \times 5.0 nm \times 5.0 nm simulation box containing 300 urea and 3085 TIP3P (137) water molecules is built to mimic the work of Salvalaglio *et al.* (88). The partial charges of urea were adopted from the Amber03 (138) database. For the aqueous simulations of glycine, we follow the simulation protocol established in Ref. 139. A simulation box of size 3.6 nm \times 3.6 nm \times 3.6 nm is constructed with 72 glycine molecules and 1200 SPC/E (140) water molecules with generalized Amber force field (GAFF) (104), and the charges on glycine are those derived using the CNDO methodology in Ref. 141. This combination of the GAFF and CNDO charges for glycine is utilized because it has been found to accurately reproduce the physical properties of both aqueous and crystalline glycine as well as inducing α -glycine crystallization from solution (128). The glycine simulations correspond to a saturated glycine solution (whose molality is experimentally measured to be 3.33 mol/kg) (139). However, since we do not explicitly calculate the solubility of glycine with the GAFF parameterization using the SPC/E water model, the solubility and the saturation level of glycine is unknown during the simulations.

For both the systems, all MD and metadynamics simulations were performed in the constant molecule number, pressure, temperature (NPT) ensemble with an integration time step of 2 fs. Systems were coupled with a thermostat of velocity rescaling scheme (142) at 300 K. The pressure was controlled using the Parrinello-Rahman barostat (143) at 1 bar. The relaxation times to thermostat and barostat were 0.1 ps and 1 ps, respectively, for both system investigated. Particle-mesh Ewald method (144) was used for the computation of long range electrostatic interaction and hydrogen bonds were constrained with the LINCS algorithm (145). The cutoffs of electrostatic and Van der Waals interaction in real space were selected to be 1.2 nm for glycine and 1.0 nm for

urea.

Rescaling Coefficients of SPIB-learned Representations. As the coefficients learned by SPIB are in dimensions of the inverse of the corresponding OPs and unstandardized as input, proper rescaling of these coefficients is needed for the purpose of evaluating the importance of each OP to the RC. Following the “betasq” protocol from Ref. 146, these rescaled coefficients are calculated as the product of the coefficients learned by SPIB and the fluctuation of the OP individually from the input trajectories (WTmetaD simulations biasing θ_1 and θ_2 OPs).

Acknowledgments

The authors thank Dr. Pablo M. Piaggi for providing source codes for the entropy OPs. We also thank Prof. John D. Weeks, Prof. Matteo Salvalaglio, Dr. Yihang Wang, Dr. Ruiyu Wang, Zachary Smith, Luke Evans, Dedi Wang, and Renjie Zhao for discussions and Dr. Ruiyu Wang, Luke Evans for proofreading the manuscript. This research was entirely supported by the U.S. Department of Energy, Office of Science, Basic Energy Sciences, CPIMS Program, under Award DE-SC0021009. We also thank Deepthought2, MARCC and XSEDE (147) (projects CHE180007P and CHE180027P) for computational resources used in this work.

1. B Peters, *Reaction Rate Theory and Rare Events Simulations*. (Elsevier, Amsterdam), (2017).
2. S Karthika, T Radhakrishnan, P Kalaichelvi, A review of classical and nonclassical nucleation theories. *Cryst. Growth & Des.* **16**, 6663–6681 (2016).
3. MH Nielsen, S Aloni, JJ De Yoreo, In situ tem imaging of caco3 nucleation reveals coexistence of direct and indirect pathways. *Science* **345**, 1158–1162 (2014).
4. T Nakamuro, M Sakakibara, H Nada, K Harano, E Nakamura, Capturing the moment of emergence of crystal nucleus from disorder. *J. Am. Chem. Soc.* **143**, 1763–1767 (2021) PMID: 33475359.
5. JJD Yoreo, et al., Crystallization by particle attachment in synthetic, biogenic, and geologic environments. *Science* **349**, aaa6760 (2015).
6. B Abécassis, et al., Persistent nucleation and size dependent attachment kinetics produce monodisperse pbs nanocrystals. *Chem. Sci.* **13**, 4977–4983 (2022).
7. PM Piaggi, J Weis, AZ Panagiotopoulos, PG Debenedetti, R Car, Homogeneous ice nucleation in an ab initio machine-learning model of water. *Proc. Natl. Acad. Sci.* **119**, e2207294119 (2022).
8. RA LaCour, TC Moore, SC Glotzer, Tuning stoichiometry to promote formation of binary colloidal superlattices. *Phys. Rev. Lett.* **128**, 188001 (2022).
9. AA Bertolazzo, D Dhabal, V Molinero, Polymorph selection in zeolite synthesis occurs after nucleation. *The J. Phys. Chem. Lett.* **13**, 977–981 (2022) PMID: 35060725.
10. LC Jacobson, W Hujo, V Molinero, Amorphous precursors in the nucleation of clathrate hydrates. *J. Am. Chem. Soc.* **132**, 11806–11811 (2010) PMID: 20669949.
11. AG Shtrikman, et al., Melt crystallization for paracetamol polymorphism. *Cryst. Growth & Des.* **19**, 4070–4080 (2019).
12. F Giberti, GA Tribello, M Parrinello, Transient polymorphism in nacl. *J. chemical theory computation* **9**, 2526–2530 (2013).
13. AR Finney, M Salvalaglio, Multiple pathways in nacl homogeneous crystal nucleation. *Faraday Discuss.* (2022).
14. ST Tsai, Z Smith, P Tiwary, Reaction coordinates and rate constants for liquid droplet nucleation: Quantifying the interplay between driving force and memory. *The J. chemical physics* **151**, 154106 (2019).
15. M Salvalaglio, T Vetter, F Giberti, M Mazzotti, M Parrinello, Uncovering molecular details of urea crystal growth in the presence of additives. *J. Am. Chem. Soc.* **134**, 17221–17233 (2012).
16. H Niu, YI Yang, M Parrinello, Temperature dependence of homogeneous nucleation in ice. *Phys. Rev. Lett.* **122**, 245501 (2019).
17. G Gobbo, MA Bellucci, GA Tribello, G Ciccotti, BL Trout, Nucleation of molecular crystals driven by relative information entropy. *J. Chem. Theory Comput.* **14**, 959–972 (2018) PMID: 29272581.
18. J De Yoreo, More than one pathway. *Nat. Mater.* **12**, 284–285 (year?).
19. J De Yoreo, A Perspective on Multistep Pathways of Nucleation. (ACS Publications), pp. 1–17 (2020).
20. DE Shaw, et al., Anton 3: Twenty microseconds of molecular dynamics simulation before lunch in *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis*, SC '21. (Association for Computing Machinery, New York, NY, USA), (2021).
21. G Torrie, J Valleau, Nonphysical sampling distributions in monte carlo free-energy estimation: Umbrella sampling. *J. Comput. Phys.* **23**, 187–199 (1977).

22. A Laio, M Parrinello, Escaping free-energy minima. *Proc. Natl. Acad. Sci.* **99**, 12562–12566 (2002).
23. RJ Allen, PB Warren, PR ten Wolde, Sampling rare switching events in biochemical networks. *Phys. Rev. Lett.* **94**, 018104 (2005).
24. T Lazaridis, M Karplus, Orientational correlations and entropy in liquid water. *The J. Chem. Phys.* **105**, 4294–4316 (1996).
25. P Rosales-Pelaez, I Sanchez-Burgos, C Valeriani, C Vega, E Sanz, Seeding approach to nucleation in the *nvt* ensemble: The case of bubble cavitation in overstretched lennard jones fluids. *Phys. Rev. E* **101**, 022611 (2020).
26. VE Bazterra, MB Ferraro, JC Facelli, Modified genetic algorithm to model crystal structures. i. benzene, naphthalene and anthracene. *The J. Chem. Phys.* **116**, 5984–5991 (2002).
27. Y Wang, JML Ribeiro, P Tiwary, Machine learning approaches for analyzing and enhancing molecular dynamics simulations. *Curr. Opin. Struct. Biol.* **61**, 139–145 (2020).
28. M Ceriotti, Unsupervised machine learning in atomistic simulations, between predictions and understanding. *The J. chemical physics* **150**, 150901 (2019).
29. M Hoffmann, et al., Deeptime: a python library for machine learning dynamical models from time series data. *Mach. Learn. Sci. Technol.* **3**, 015009 (2021).
30. M Ghorbani, S Prasad, JB Klauda, BR Brooks, Graphvampnet, using graph neural networks and variational approach to markov processes for dynamical modeling of biomolecules. *The J. Chem. Phys.* **156**, 184103 (2022).
31. S Banik, et al., Cegan: Crystal edge graph attention network for multiscale classification of materials environment (2022).
32. RS Hong, et al., Insights into the polymorphic structures and enantiotropic layer-slip transition in paracetamol form iii from enhanced molecular dynamics. *Cryst. Growth & Des.* **21**, 886–896 (2021).
33. Y Mori, Ki Okazaki, T Mori, K Kim, N Matubayasi, Learning reaction coordinates via cross-entropy minimization: Application to alanine dipeptide. *The J. Chem. Phys.* **153**, 054115 (2020).
34. J Rogal, E Schneider, ME Tuckerman, Neural-network-based path collective variables for enhanced sampling of phase transformations. *Phys. Rev. Lett.* **123**, 245701 (2019).
35. RS DeFever, S Sarupria, Contour forward flux sampling: Sampling rare events along multiple collective variables. *The J. Chem. Phys.* **150**, 024103 (2019).
36. C Velez-Vega, EE Borrero, FA Escobedo, Kinetics and mechanism of the unfolding native-to-loop transition of trp-cage in explicit solvent via optimized forward flux sampling simulations. *The J. Chem. Phys.* **133**, 105103 (2010).
37. A Barducci, G Bussi, M Parrinello, Well-tempered metadynamics: A smoothly converging and tunable free-energy method. *Phys. Rev. Lett.* **100**, 020603 (2008).
38. F Giberti, M Salvalaglio, M Parrinello, Metadynamics studies of crystal nucleation. *IUCrJ* **2**, 256–266 (2015).
39. T Karmakar, M Ivernizzi, V Rizzi, M Parrinello, Collective variables for the study of crystallisation. *Mol. Phys.* **0**, e1893848 (2021).
40. GC Sosso, et al., Crystal nucleation in liquids: Open questions and future challenges in molecular dynamics simulations. *Chem. reviews* **116**, 7078–7116 (2016).
41. S Hussain, A Haji-Akbari, Studying rare events using forward-flux sampling: Recent breakthroughs and future outlook. *The J. Chem. Phys.* **152**, 060901 (2020).
42. KE Blow, D Quigley, GC Sosso, The seven deadly sins: when computing crystal nucleation rates, the devil is in the details (2021).
43. B Peters, BL Trout, Obtaining reaction coordinates by likelihood maximization. *The J. Chem. Phys.* **125**, 054108 (2006).
44. A Ma, AR Dinner, Automatic method for identifying reaction coordinates in complex systems. *The J. Phys. Chem. B* **109**, 6769–6779 (2005).
45. EE Borrero, FA Escobedo, Reaction coordinates and transition pathways of rare events via forward flux sampling. *The J. chemical physics* **127**, 164101 (2007).
46. Z Zou, ST Tsai, P Tiwary, Toward automated sampling of polymorph nucleation and free energies with the sgooop and metadynamics. *The J. Phys. Chem. B* **125**, 13049–13056 (2021).
47. CR Schwantes, VS Pande, Improvements in markov state model construction reveal many non-native interactions in the folding of ntl9. *J. Chem. Theory Comput.* **9**, 2000–2009 (2013) PMID: 23750122.
48. EE Borrero, FA Escobedo, Reaction coordinates and transition pathways of rare events via forward flux sampling. *The J. chemical physics* **127**, 164101 (2007).
49. RS DeFever, S Sarupria, Nucleation mechanism of clathrate hydrates of water-soluble guest molecules. *The J. Chem. Phys.* **147**, 204503 (2017).
50. NER Zimmermann, B Vorselaars, D Quigley, B Peters, Nucleation of nacl from aqueous solution: Critical sizes, ion-attachment kinetics, and rates. *J. Am. Chem. Soc.* **137**, 13352–13361 (2015).
51. A Arjun, PG Bolhuis, Molecular understanding of homogeneous nucleation of co2 hydrates using transition path sampling. *The J. Phys. Chem. B* **125**, 338–349 (2021).
52. LC Jacobson, M Matsumoto, V Molinero, Order parameters for the multistep crystallization of clathrate hydrates. *The J. Chem. Phys.* **135**, 074501 (2011).
53. M Badin, R Martoňák, Nucleating a different coordination in a crystal under pressure: A study of the *b1*–*b2* transition in nacl by metadynamics. *Phys. Rev. Lett.* **127**, 105701 (2021).
54. A Samanta, ME Tuckerman, TQ Yu, W E, Microscopic mechanisms of equilibrium melting of a solid. *Science* **346**, 729–732 (2014).
55. H Song, L Vogt-Maranto, R Wiscons, AJ Matzger, ME Tuckerman, Generating cocrystal polymorphs with information entropy driven by molecular dynamics-based enhanced sampling. *The J. Phys. Chem. Lett.* **11**, 9751–9758 (2020) PMID: 33141590.
56. D Wang, P Tiwary, State predictive information bottleneck. *The J. Chem. Phys.* **154**, 134111 (2021).
57. S Mehdi, D Wang, S Pant, P Tiwary, Accelerating all-atom simulations and gaining mechanistic understanding of biophysical systems through state predictive information bottleneck. *J. Chem. Theory Comput.* **18**, 3231–3238 (2022).
58. PR ten Wolde, D Frenkel, Computer simulation study of gas–liquid nucleation in a lennard-jones system. *J. Chem. Phys.* **109**, 9901–9918 (1998).
59. PM Piaggi, O Valsson, M Parrinello, Enhancing entropy and enthalpy fluctuations to drive crystallization in atomistic simulations. *Phys. Rev. Lett.* **119**, 015701 (2017).
60. PM Piaggi, M Parrinello, Predicting polymorphism in molecular crystals using orientational entropy. *Proc. Natl. Acad. Sci.* **115**, 10251–10256 (2018).
61. PJ Steinhardt, DR Nelson, M Ronchetti, Bond-orientational order in liquids and glasses. *Phys. Rev. B* **28**, 784 (1983).
62. P Tiwary, M Parrinello, A time-independent free energy estimator for metadynamics. *The J. Phys. Chem. B* **119**, 736–742 (2015).
63. M Salvalaglio, P Tiwary, GM Maggioni, M Mazzotti, M Parrinello, Overcoming time scale and finite size limitations to compute nucleation rates from small scale well tempered metadynamics simulations. *The J. Chem. Phys.* **145**, 211925 (2016).
64. KM Bal, Nucleation rates from small scale atomistic simulations and transition state theory. *The J. Chem. Phys.* **155**, 144111 (2021).
65. D Moroni, PR ten Wolde, PG Bolhuis, Interplay between structure and size in a critical crystal nucleus. *Phys. Rev. Lett.* **94**, 235703 (2005).
66. P Rein ten Wolde, MJ Ruiz-Montero, D Frenkel, Numerical calculation of the rate of crystal nucleation in a lennard-jones system at moderate undercooling. *The J. Chem. Phys.* **104**, 9932–9947 (1996).
67. P Rein ten Wolde, D Frenkel, Homogeneous nucleation and the ostwald step rule. *Phys. Chem. Chem. Phys.* **1**, 2191–2196 (1999).
68. I Volkov, M Cieplak, J Koplik, JR Banavar, Molecular dynamics simulations of crystallization of hard spheres. *Phys. Rev. E* **66**, 061401 (2002).
69. S Auer, D Frenkel, Prediction of absolute crystal-nucleation rate in hard-sphere colloids. *Nature* **409**, 1020–1023 (2001).
70. M Chopra, M Müller, JJ de Pablo, Order-parameter-based monte carlo simulation of crystallization. *The J. Chem. Phys.* **124**, 134102 (2006).
71. C Desgranges, J Delhomelle, Molecular mechanism for the cross-nucleation between polymorphs. *J. Am. Chem. Soc.* **128**, 10368–10369 (2006).
72. R Radhakrishnan, BL Trout, Nucleation of hexagonal ice (ih) in liquid water. *J. Am. Chem. Soc.* **125**, 7743–7747 (2003).
73. R Radhakrishnan, BL Trout, Nucleation of crystalline phases of water in homogeneous and inhomogeneous environments. *Phys. Rev. Lett.* **90**, 158301 (2003).
74. D Quigley, PM Rodger, Metadynamics simulations of ice nucleation and growth. *The J. Chem. Phys.* **128**, 154518 (2008).
75. AV Brukhno, J Anwar, R Davidchack, R Handel, Challenges in molecular simulation of homogeneous ice nucleation. *J. Physics: Condens. Matter* **20**, 494243 (2008).
76. A Reinhardt, JPK Doye, EG Noya, C Vega, Local order parameters for use in driving homogeneous ice nucleation with all-atom models of water. *The J. Chem. Phys.* **137**, 194504 (2012).
77. W Lechner, C Dellago, Accurate determination of crystal structures based on averaged local bond order parameters. *The J. Chem. Phys.* **129**, 114707 (2008).
78. PM Piaggi, M Parrinello, Calculation of phase diagrams in the multithermal-multibarc ensemble. *The J. Chem. Phys.* **150**, 244119 (2019).
79. W Humphrey, A Dalke, K Schulten, Vmd: Visual molecular dynamics. *J. Mol. Graph.* **14**, 33–38 (1996).
80. A Stukowski, Visualization and analysis of atomistic simulation data with OVITO-the Open Visualization Tool. *MODELLING AND SIMULATION IN MATERIALS SCIENCE AND ENGINEERING* **18** (2010).
81. ER Beyerle, S Mehdi, P Tiwary, Quantifying energetic and entropic pathways in molecular systems. *The J. Phys. Chem. B* **126**, 3950–3960 (2022) PMID: 35605180.
82. DP Kingma, M Welling, Auto-encoding variational bayes. *2nd Int. Conf. on Learn. Represent. ICLR 2014 - Conf. Track Proc.* pp. 1–14 (2014).
83. I Higgins, et al., beta-VAE: Learning basic visual concepts with a constrained variational framework in International Conference on Learning Representations. (2017).
84. AA Alemi, I Fischer, JV Dillon, K Murphy, Deep variational information bottleneck. *CoRR abs/1612.00410* (2016).
85. N Tishby, FC Pereira, W Bialek, The information bottleneck method (2000).
86. F Noé, I Horenko, C Schütte, JC Smith, Hierarchical analysis of conformational dynamics in biomolecules: Transition networks of metastable states. *J. Chem. Phys.* **126**, 155102 (2007).
87. P Deuffhard, W Huisinga, A Fischer, C Schütte, Identification of almost invariant aggregates in reversible nearly uncoupled Markov chains. *Linear Algebr. Its Appl.* **315**, 39–59 (2000).
88. M Salvalaglio, C Perego, F Giberti, M Mazzotti, M Parrinello, Molecular-dynamics simulations of urea nucleation from aqueous solution. *Proc. Natl. Acad. Sci.* **112**, E6–E14 (2015).
89. M Salvalaglio, M Mazzotti, M Parrinello, Urea homogeneous nucleation mechanism is solvent dependent. *Faraday Discuss.* **179**, 291–307 (2015).
90. MJ Abraham, et al., Gromacs: High performance molecular simulations through multi-level parallelism from laptops to supercomputers. *SoftwareX* **1**, 19–25 (2015).
91. GA Tribello, M Bonomi, D Branduardi, C Camilloni, G Bussi, Plumed 2: New feathers for an old bird. *Comp. Phys. Comm.* **185**, 604–613 (2014).
92. M Bonomi, G Bussi, CC Camilloni, Promoting transparency and reproducibility in enhanced molecular simulations. *Nat. Methods.* **16**, 670–673 (2019).
93. HP Weber, WG Marshall, V Dmitriev, High-pressure polymorphism in deuterated urea. *Acta Crystallogr. Sect. A* **58**, c174 (2002).
94. FJ Lamelas, ZA Dreger, YM Gupta, Raman and x-ray scattering studies of high-pressure phases of urea. *The J. Phys. Chem. B* **109**, 8206–8215 (2005).
95. A Olejniczak, K Ostrowska, A Katusiak, H-bond breaking in high-pressure urea. *The J. Phys. Chem. C* **113**, 15761–15767 (2009).
96. M Donnelly, et al., Urea and deuterium mixtures at high pressures. *The J. Chem. Phys.* **142**, 124503 (2015).
97. K Dziubek, M Citroni, S Fanetti, AB Cairns, R Bini, High-pressure high-temperature structural properties of urea. *The J. Phys. Chem. C* **121**, 2380–2387 (2017).
98. K Roszak, A Katusiak, Giant anomalous strain between high-pressure phases and the

- mesomers of urea. *The J. Phys. Chem. C* **121**, 778–784 (2017).
99. F Safari, M Tkacz, A Katrusiak, High-pressure sorption of hydrogen in urea. *The J. Phys. Chem. C* **125**, 7756–7762 (2021).
 100. F Giberti, M Salvalaglio, M Mazzotti, M Parrinello, Insight into the nucleation of urea crystals from the melt. *Chem. Eng. Sci.* **121**, 51–59 (2015).
 101. T Mandal, RG Larson, Nucleation of urea from aqueous solution: Structure, critical size, and rate. *The J. Chem. Phys.* **146**, 134501 (2017).
 102. NF Francia, LS Price, J Nyman, SL Price, M Salvalaglio, Systematic finite-temperature reduction of crystal energy landscapes. *Cryst. Growth & Des.* **20**, 6847–6862 (2020).
 103. C Shang, XJ Zhang, ZP Liu, Crystal phase transition of urea: what governs the reaction kinetics in molecular crystal phase transitions. *Phys. Chem. Chem. Phys.* **19**, 32125–32131 (2017).
 104. J Wang, RM Wolf, JW Caldwell, PA Kollman, DA Case, Development and testing of a general amber force field. *J. Comput. Chem.* **25**, 1157–1174 (2004).
 105. O Valsson, P Tiwary, M Parrinello, Enhancing important fluctuations: Rare events and metadynamics from a conceptual viewpoint. *Annu. review physical chemistry* **67**, 159–184 (2016).
 106. GA Tribello, F Giberti, GC Sossio, M Salvalaglio, M Parrinello, Analyzing and driving cluster formation in atomistic simulations. *J. Chem. Theory Comput.* **13**, 1317–1327 (2017).
 107. DL Nelson, MM Cox, *Lehninger Principles of Biochemistry, Fifth Edition*. (Freeman), (2008).
 108. JW Lynch, Molecular structure and function of the glycine receptor chloride channel. *Physiol. Rev.* **84**, 1051–1095 (2004) PMID: 15383648.
 109. MA Razak, PS Begum, B Viswanath, S Rajagopal, Multifarious beneficial effect of nonessential amino acid, glycine: A review. *Oxidative Medicine Cell. Longev.* **2017**, 1716701 (2017).
 110. YG Bushuev, SV Davletbaeva, OI Koifman, Molecular dynamics simulations of aqueous glycine solutions. *Cryst. Eng. Comm.* **19**, 7197–7206 (2017).
 111. E Boldyreva, et al., Polymorphism of glycine, part i. *J. thermal analysis calorimetry* **73**, 409–418 (2003).
 112. Y Yani, PS Chow, RB Tan, Glycine open dimers in solution: New insights into α -glycine nucleation and growth. *Cryst. growth & design* **12**, 4771–4778 (2012).
 113. E Boldyreva, Glycine: The gift that keeps on giving. *Isr. J. Chem.* (2021).
 114. CL Bull, et al., ζ -Glycine: insight into the mechanism of a polymorphic phase transition. *IUCrJ* **4**, 569–574 (2017).
 115. S Prestipino, PV Giaquinta, The entropy multiparticle-correlation expansion for a mixture of spherical and elongated particles. *J. Stat. Mech. Theory Exp.* **2004**, P09008 (2004).
 116. A Baranyai, DJ Evans, Direct entropy calculation from computer simulation of liquids. *Phys. Rev. A* **40**, 3817–3822 (1989).
 117. Y Iitaka, The crystal structure of γ -glycine. *Acta Crystallogr.* **14**, 1–10 (1961).
 118. N Marom, et al., Many-body dispersion interactions in molecular crystal polymorphism. *Angewandte Chemie Int. Ed.* **52**, 6629–6632 (2013).
 119. N Duff, B Peters, Polymorph specific rmsd local order parameters for molecular crystals and nuclei: α -, β -, and γ -glycine. *The J. Chem. Phys.* **135**, 134101 (2011).
 120. Y Iitaka, The crystal structure of β -glycine. *Acta Crystallogr.* **13**, 35–45 (1960).
 121. A Dawson, et al., Effect of high pressure on the crystal structures of polymorphs of glycine. *Cryst. Growth & Des.* **5**, 1415–1427 (2005).
 122. RA Van Santen, The ostwald step rule. *The J. Phys. Chem.* **88**, 5768–5769 (1984).
 123. O Urquidi, J Brazard, N LeMessurier, L Simine, TBM Adachi, In situ optical spectroscopy of crystallization: One crystal nucleation at a time. *Proc. Natl. Acad. Sci.* **119**, e2122990119 (2022).
 124. A Bouchard, GW Hofland, GJ Witkamp, Solubility of glycine polymorphs and recrystallization of β -glycine. *J. Chem. & Eng. Data* **52**, 1626–1629 (2007).
 125. C Parks, et al., Solubility curves and nucleation rates from molecular dynamics for polymorph prediction – moving beyond lattice energy minimization. *Phys. Chem. Chem. Phys.* **19**, 5285–5295 (2017).
 126. J Honeycutt, HC Andersen, The effect of periodic boundary conditions on homogeneous nucleation observed in computer simulations. *Chem. Phys. Lett.* **108**, 535–538 (1984).
 127. WC Swope, HC Andersen, 10^6 -particle molecular-dynamics study of homogeneous nucleation of crystals in a supercooled atomic liquid. *Phys. Rev. B* **41**, 7042–7054 (1990).
 128. DW Cheong, YD Boon, Comparative study of force fields for molecular dynamics simulations of α -glycine crystal growth from solution. *Cryst. growth & design* **10**, 5146–5158 (2010).
 129. ET Broadhurst, et al., Polymorph evolution during crystal growth studied by 3d electron diffraction. *IUCrJ* **7**, 5–9 (2020).
 130. LJ Little, RP Sear, JL Keddie, Does the γ polymorph of glycine nucleate faster? a quantitative study of nucleation from aqueous solution. *Cryst. Growth & Des.* **15**, 5345–5354 (2015).
 131. G He, et al., Direct growth of γ -glycine from neutral aqueous solutions by slow, evaporation-driven crystallization. *Cryst. Growth & Des.* **6**, 1746–1749 (2006).
 132. CE Hughes, S Hamad, KD Harris, CRA Catlow, PC Griffiths, A multi-technique approach for probing the evolution of structural properties during crystallization of organic materials from solution. *Faraday discussions* **136**, 71–89 (2007).
 133. CE Hughes, KD Harris, The effect of deuteration on polymorphic outcome in the crystallization of glycine from aqueous solution. *New J. Chem.* **33**, 713–716 (2009).
 134. JE Aber, S Arnold, BA Garetz, AS Myerson, Strong dc electric field applied to supersaturated aqueous glycine solution induces nucleation of the γ polymorph. *Phys. Rev. Lett.* **94**, 145503 (2005).
 135. T Karmakar, PM Piaggi, M Parrinello, Molecular dynamics simulations of crystal nucleation from solution at constant chemical potential. *J. Chem. Theory Comput.* **15**, 6923–6930 (2019) PMID: 31657927.
 136. Y Wang, JML Ribeiro, P Tiwary, Past–future information bottleneck for sampling molecular reaction coordinate simultaneously with thermodynamics and kinetics. *Nat. communications* **10**, 1–8 (2019).
 137. WL Jorgensen, J Chandrasekhar, JD Madura, RW Impey, ML Klein, Comparison of simple potential functions for simulating liquid water. *The J. Chem. Phys.* **79**, 926–935 (1983).
 138. DA Case, et al., The amber biomolecular simulation programs. *J. Comput. Chem.* **26**, 1668–1688 (2005).
 139. YG Bushuev, SV Davletbaeva, OI Koifman, Molecular dynamics simulations of aqueous glycine solutions. *CrystEngComm* **19**, 7197–7206 (2017).
 140. HJC Berendsen, JR Grigera, TP Straatsma, The missing term in effective pair potentials. *The J. Phys. Chem.* **91**, 6269–6271 (1987).
 141. JL Derissen, PH Smit, J Voogd, Calculation of the electrostatic lattice energies of α -, β -, and γ -glycine. *The J. Phys. Chem.* **81**, 1474–1476 (1977).
 142. G Bussi, D Donadio, M Parrinello, Canonical sampling through velocity rescaling. *J. Chem. Phys.* **126**, 014101 (2007).
 143. M Parrinello, A Rahman, Polymorphic transitions in single crystals: A new molecular dynamics method. *J. Appl. physics* **52**, 7182–7190 (1981).
 144. U Essmann, L Perera, ML Berkowitz, A smooth particle mesh ewald method. *J. Chem. Phys.* **103**, 8577–8593 (1995).
 145. B Hess, P-lincs: a parallel linear constraint solver for molecular simulation. *J. Chem. Theory* **4**, 116–122 (2007).
 146. U Groemping, Relative importance for linear regression in r: The package relaimpo. *J. Stat. Softw.* **17**, 1–27 (2006).
 147. J Towns, et al., Xsede: Accelerating scientific discovery. *Comput. Sci. Eng.* **16**, 62–74 (2014).