

# Phase diagrams—why they matter and how to predict them

Pin Yu Chew and Aleks Reinhardt

*Yusuf Hamied Department of Chemistry, University of Cambridge, Lensfield Road, Cambridge, CB2 1EW, United Kingdom*

(Dated: 28 November 2022)

Understanding the thermodynamic stability and metastability of materials can help us to gauge for example whether crystalline polymorphs in pharmaceutical formulations are likely to be durable. It can also help us to design experimental routes to novel phases with potentially interesting properties. In this article, we provide an overview of how thermodynamic phase behaviour can be quantified both in computer simulations and using machine-learning approaches to determining phase diagrams, as well as combinations of the two. We review the basic workflow of free-energy computations for condensed phases, including some practical implementation advice, ranging from the Frenkel–Ladd approach to thermodynamic integration and to direct-coexistence simulations. We illustrate the applications of such methods on a range of systems from materials chemistry to biological phase separation. Finally, we outline some challenges, questions and practical applications of phase-diagram determination which we believe are likely to be possible to address in the near future using such state-of-the-art free-energy calculations, which may provide fundamental insight into separation processes using multicomponent solvents.

## I. INTRODUCTION

Knowing a material’s thermodynamic stability under different conditions, as summarised in a phase diagram, is of considerable practical importance with numerous technological applications, for instance in separation processes.<sup>1</sup> In simple cases, an experimental determination of phase diagrams is often sufficient, but as the building blocks become more complex, it is not always clear which polymorphs may be stable, and if they are not, what the limits of their metastability may be. For example, the phase diagram of water has been studied for over a hundred years with progressively more solid phases discovered and characterised,<sup>2,3</sup> while the famous case of HIV drug ritonavir suddenly converting into a largely insoluble most stable polymorph is thought not to be an isolated case amongst pharmaceutical compounds.<sup>4</sup>

Questions of thermodynamic stability and metastability can in principle be addressed using statistical mechanics and computer simulations, as first shown by Hoover and Ree in the context of hard-sphere melting back in 1968.<sup>5</sup> On the one hand, computing phase diagrams in silico can be very helpful under conditions that are difficult to achieve experimentally, such as when studying high-pressure behaviour relevant to planetary cores<sup>6–9</sup> or the properties of synthetic elements,<sup>10</sup> and can thus help guide experimental efforts. On the other hand, phase diagrams are very sensitive to the potentials used to describe the building blocks themselves<sup>11–14</sup> and a potential that captures phase behaviour accurately will often also describe other properties of a substance well. Classical potentials can therefore be parameterised by fine-tuning their ability to reproduce phase behaviour; such an approach has fruitfully been used to design ever better water<sup>15–25</sup> and protein<sup>26–28</sup> models.

Indeed, when computing phase diagrams with computer simulations, the first step is generally to choose how to simulate the material of interest. This can range from quantum-mechanically accurate potentials<sup>29</sup> based on either wavefunction or density-functional theory (DFT) electronic structure calculations, to classical potentials parameterised either from quantum simulations<sup>30</sup> or from experiment,<sup>31</sup> to coarse-grained<sup>32–35</sup> and ‘toy model’ potentials.<sup>36–44</sup> Such an approach has been used to investigate the phase behaviour of water,<sup>7–9,45,46</sup> gallium,<sup>47</sup> supercritical hydrogen<sup>6</sup> and titanium dioxide.<sup>48</sup>

The computation of phase diagrams has been thoroughly

reviewed in an excellent paper by Vega and co-workers.<sup>18</sup> In this perspective, we briefly summarise the fundamental principles, as well as review some more recent work on determining phase diagrams in computer simulations, and identify some of the trickier aspects of the process and caveats involved. We also provide a short overview of some machine-learning methods that have recently been used to predict phase diagrams without the need for expensive molecular simulations. Finally, we speculate about some possible future applications of the methods that have been developed to investigate phase behaviour.

## II. STRATEGIES FOR COMPUTING CHEMICAL POTENTIALS

When the temperatures, pressures and chemical potentials of two phases are equal, they are at equilibrium and are said to coexist. It is straightforward in computer simulations to fix the temperature and pressure, but the chemical potential is often more difficult to compute because it entails a ‘thermal’ (i.e. non-mechanical) component, namely the entropy, that cannot be estimated by sampling. Perhaps the simplest approach to computing the phase behaviour of a system is therefore to assume that the chemical potential is completely dominated by the enthalpic term, and that the entropy differences between phases are negligible. Such an approach is advantageous because the enthalpy is a simple mechanical observable that can readily be determined, and is often not unreasonable when dealing with phases of a similar structure and thus similar entropy. In the context of solid phases, the phase behaviour is in many cases completely dominated by enthalpic terms and the entropy can be relatively unimportant.<sup>12,49</sup> For example, many such ‘absolute-zero’ phase diagrams have been reported for titanium dioxide.<sup>50–57</sup> However, when competing phases have relatively similar enthalpies, entropy differences can be crucial in controlling the phase behaviour.<sup>58</sup> Sometimes, approximations can be made to estimate the free energy as a function of temperature.<sup>56,58–63</sup>

Even if an approximate zero-temperature phase diagram is to be calculated, a further challenge remains when predicting phase diagrams of solid phases: how can we ensure that all potentially stable polymorphs have been identified? We can consider any phases known experimentally, but a significant advantage of computer simulation is precisely that we may be

88 able to predict phases that are not yet known experimentally and  
 89 thus help steer experimental exploration of new phases. There  
 90 are numerous approaches that can help us to identify potentially  
 91 stable crystal phases.<sup>64–68</sup> Most such approaches usually entail  
 92 some stochastic step to identify potentially unseen phases and  
 93 a subsequent energy minimisation step, often followed by  
 94 identification of the crystal space-group symmetry, for example  
 95 with tools like FINDSYM.<sup>69</sup> Broadly speaking, the larger a  
 96 crystal’s unit cell, the more difficult it is to find by random  
 97 searching. Since energy minimisation is usually involved,  
 98 these approaches are often able to identify polymorphs that are  
 99 competitive at very low temperatures, but may be less adept at  
 100 identifying high-temperature phases that are only entropically  
 101 stabilised. In random searches, low-enthalpy structures have  
 102 been shown to occur more frequently and so have a larger basin  
 103 of attraction;<sup>70</sup> indeed, the volume of the basin of attraction of  
 104 each minimum can be approximated by the number of times  
 105 each structure is found in a random search, and this can provide  
 106 a crude first approximation of the relative entropies of the  
 107 competing polymorphs.<sup>49</sup>

## 108 A. Direct estimation of chemical potentials

109 If we wish to compute the chemical potential in computer  
 110 simulations, we can note that, although the chemical potential  
 111 cannot directly be determined as a canonical average over the  
 112 phase space, there are several approaches that we can take  
 113 to compute it. Using the Widom insertion method, we can  
 114 relate the excess chemical potential to the thermal average of a  
 115 Boltzmann factor for the change in energy  $\Delta U(N, N + 1)$  when  
 116 an additional particle is randomly inserted into a system of  $N$   
 117 particles,<sup>71</sup> i.e.

$$\mu^{\text{ex}} = -k_{\text{B}}T \ln \left[ \int \langle \exp(-\beta \Delta U(N, N + 1)) \rangle_N d\mathbf{r}_{N+1} \right]. \quad (1)$$

118 This and similar methods, such as the Gibbs ensemble ap-  
 119 proach,<sup>72,73</sup> can be used to determine the phase behaviour of  
 120 systems that are sufficiently dilute to enable additional particles  
 121 to be inserted with a non-negligible probability. Phase switch-  
 122 ing<sup>74,75</sup> is an approach based on a similar idea, but is also  
 123 applicable to solid systems.

## 124 B. Thermodynamic integration

125 For denser systems, however, we can also note that although  
 126 the chemical potential is a thermal quantity, its derivatives are  
 127 mechanical observables. For example, from the Gibbs–Duhem  
 128 relation  $d\mu = v dP - s dT$  (where  $\mu$  is the chemical potential,  
 129  $v$  is the volume per particle,  $P$  is the pressure,  $s$  is the entropy  
 130 per particle and  $T$  is the absolute temperature), we can write  
 131  $(\partial\mu/\partial P)_T = v$ . By integrating this equation numerically from  
 132 some initial pressure  $P_0$ , we obtain

$$\mu(P) = \mu(P_0) + \int_{P_0}^P v(P') dP'. \quad (2)$$

133 The analogous derivative with respect to temperature  
 134  $((\partial\mu/\partial T)_P = s)$  is unhelpful, since the entropy is also not

135 a mechanical observable, but we can instead use the Gibbs–  
 136 Helmholtz equation,

$$\left( \frac{\partial(G/T)}{\partial T} \right)_{N, P} = -H/T^2, \quad (3)$$

137 where  $G$  is the Gibbs energy,  $N$  is the number of particles and  $H$   
 138 is the enthalpy. Since for a one-component system  $G = N\mu$ , this  
 139 thus relates the derivative of the chemical potential to another  
 140 mechanical observable, the enthalpy. Similar derivatives can  
 141 be constructed where variables other than the pressure or the  
 142 temperature are held fixed, e.g. along iso- $(\beta P)$  lines.<sup>76,77</sup>

## 143 1. Choice of reference state

144 In principle, if we can determine how the enthalpy and the  
 145 density of a system change as a function of the thermodynamic  
 146 variables of interest, we can also determine how the chemical  
 147 potential of a phase changes relative to the starting point.  
 148 However, this does not yet give us an ‘absolute’ chemical  
 149 potential of the phase in question, since the starting point is  
 150 arbitrary. We cannot in general integrate different phases to  
 151 the same starting point, since paths along which we perform  
 152 thermodynamic integration must be reversible to maintain  
 153 numerical stability. In order to ensure that chemical potentials  
 154 of different phases are expressed relative to a common origin,  
 155 there are several possible strategies. The simplest strategy would  
 156 be to determine a coexistence point of each phase of interest with  
 157 another common phase (e.g. the liquid or the vapour); this is the  
 158 basic idea behind direct-coexistence simulations (Sec. II D). It  
 159 is also possible to perform thermodynamic integration along a  
 160 non-physical reversible path from a liquid to a solid,<sup>78</sup> such that  
 161 the chemical potential is well defined throughout the process.

162 Alternatively, we can relate the chemical potential to a  
 163 state whose free energy can be computed analytically. There  
 164 are not many such states, but examples include the perfect  
 165 gas and the Einstein and Debye crystals.<sup>18,79–82</sup> However, such  
 166 reference states with analytically computable partition functions  
 167 are simple non-interacting systems, and to relate them to an  
 168 interacting system, we need to find a reversible path from the  
 169 reference state to the conditions of interest. For fluid systems,  
 170 this can often readily be achieved simply by reducing the density  
 171 until the particles are on average too far to interact; if the system  
 172 has long-ranged electrostatic interactions, these can often be  
 173 switched off gradually (see Sec. II B 2 below). In principle,  
 174 the Helmholtz energy could be integrated in volume using the  
 175 standard relation  $(\partial A/\partial V)_{N, T} = -P$  expressed as an integral  
 176 in density,

$$\beta a(\rho_1) = \beta a(\rho_0) + \int_{\rho_0}^{\rho_1} \frac{\beta P}{\rho^2} d\rho, \quad (4)$$

177 where  $\beta = 1/k_{\text{B}}T$  and  $a = A/N$  is the Helmholtz energy per  
 178 particle. We could then use the ideal Helmholtz energy for  
 179 the reference state  $\beta a(\rho_0) = \beta a^{\text{id}} = \ln(\rho\Lambda^3) - 1$ , where  $\Lambda$  is  
 180 the de Broglie thermal wavelength. The difficulty with this  
 181 expression is that the integrand scales as  $1/\rho$  at low densities  
 182 (for which  $P \rightarrow P^{\text{id}} = \rho k_{\text{B}}T$ ), whilst  $\lim_{x \rightarrow 0} \ln x = -\infty$ . To  
 183 maintain numerical stability, it is therefore usual to compute  
 184 the excess Helmholtz energy,  $a^{\text{ex}} \equiv a - a^{\text{id}}$ , using the relation<sup>18</sup>

$$\beta a^{\text{ex}}(\rho_1) = \int_{\rho_0}^{\rho_1} \left[ \frac{\beta P}{\rho^2} - \frac{1}{\rho} \right] d\rho, \quad (5)$$

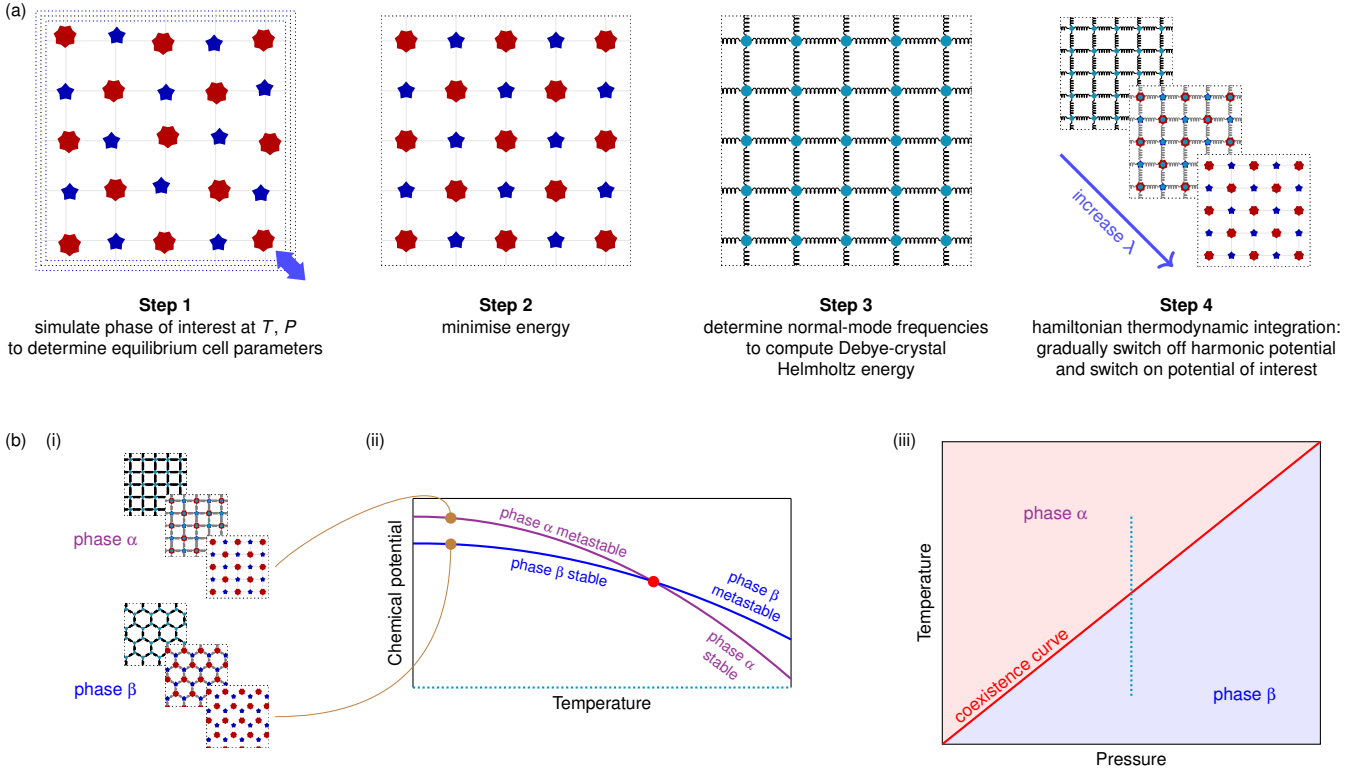


Figure 1. (a) A schematic illustration of the steps involved in determining the initial chemical potential of a solid phase. (b) A schematic of the procedure to determine the phase diagram. The process in panel (a) can be applied to different phases, here labelled  $\alpha$  and  $\beta$  [panel (i)]. This leads to two initial chemical potentials at the initial conditions investigated, shown as brown points in panel (ii). Thermodynamic integration at a constant pressure from these initial points can be used to compute the chemical potential at other temperatures [panel (ii)]. In this illustration, at low temperatures, phase  $\beta$  is thermodynamically stable and phase  $\alpha$  is metastable, while at higher temperatures the converse holds. The point labelled in red is the point at which the temperatures, pressures and chemical potentials of the two phases are equal, and is therefore part of the coexistence curve [panel (iii)]. The data in panel (ii) correspond to a slice of panel (iii), as indicated by the dotted cyan line. The rest of the coexistence curve can be obtained for example by repeating the entire calculation at other pressures; by using thermodynamic integration as a function of pressure and then repeating step (ii) to give another coexistence point; or by using a method such as Gibbs–Duhem integration.

185 with  $a^{\text{id}}(\rho_1)$  added on at the density of interest  $\rho_1$ , rather than  
 186 in the limit of low density, in the final step. In this integration,  
 187 as the density tends to zero, the integrand tends to the second  
 188 virial coefficient,<sup>18,38</sup> which can be computed independently in  
 189 a Monte Carlo (MC) calculation.

## 190 2. Artificial thermodynamic integration

191 For solids and systems with long-ranged interactions, another  
 192 approach is to define a reference potential  $U_0$  corresponding  
 193 to a simpler system whose partition function is calculable, and  
 194 relate it to  $U_1$ , the potential of interest, often with a simple  
 195 linear scaling such as  $U(\lambda) = \lambda U_1 + (1 - \lambda)U_0$ .<sup>83</sup> When the  
 196 parameter  $\lambda$  is zero, this potential is equivalent to the reference  
 197 potential, whilst when it is unity, it is equivalent to the potential  
 198 of interest. The canonical partition function of this potential,  
 199  $Q(\lambda) = \int \exp(-\beta U(\mathbf{r}^N; \lambda)) \mathbf{d}\mathbf{r}^N$ , is of course also a function  
 200 of  $\lambda$ , and is in principle not feasible to compute analytically.  
 201 However, as we have noted, derivatives of the free energy are  
 202 often mechanical observables, and using the product rule and

203 the bridge relation  $A = -k_B T \ln Q$ , we can readily compute that

$$\left( \frac{\partial A(\lambda)}{\partial \lambda} \right)_{NVT} = \frac{\int \left( \frac{\partial U(\mathbf{r}^N; \lambda)}{\partial \lambda} \right) \exp[-\beta U(\mathbf{r}^N; \lambda)] \mathbf{d}\mathbf{r}^N}{\int \exp[-\beta U(\mathbf{r}^N; \lambda)] \mathbf{d}\mathbf{r}^N} \quad (6)$$

$$= \left\langle \frac{\partial U(\lambda)}{\partial \lambda} \right\rangle_{\lambda}.$$

204 Since by construction we can compute  $A_0 \equiv A(\lambda = 0)$  analytic-  
 205 ally, we can compute the Helmholtz energy of the potential of  
 206 interest as

$$A_1 = A_0 + \int_0^1 \langle U_1 - U_0 \rangle_{\lambda} \mathbf{d}\lambda, \quad (7)$$

207 provided the path is reversible. This approach is known as  
 208 artificial or hamiltonian thermodynamic integration,<sup>79</sup> and with  
 209 it, we can compute an initial Helmholtz energy of the system  
 210 with the potential of interest. If we add a suitable  $PV$  term,<sup>80</sup>  
 211 we can also determine the Gibbs energy and hence the absolute  
 212 chemical potential.

## 213 3. Einstein and Debye crystals

214 For crystalline phases, a possible reference system may be the  
 215 Einstein crystal, a very simple model of a crystal where particles

are tethered to their lattice positions by harmonic springs. It is straightforward to compute the canonical partition function of such a crystal,<sup>18,79</sup> even for non-spherical molecules.<sup>84</sup> The spring constants can then gradually be made weaker whilst the potential of interest is switched on using hamiltonian thermodynamic integration, as discussed above. In practice, to reduce the scope for numerical issues in the integration, an additional constraint of a fixed centre of mass is introduced; this can be corrected for analytically, but with reasonable system sizes, the correction arising from this constraint is in any case very small.<sup>80</sup> An alternative approach is to constrain the position of a single molecule.<sup>85</sup> Another possible crystalline reference state is the Debye crystal, where harmonic springs connect atoms in such a way that their frequencies correspond to the normal-mode frequencies of the original crystal.

In the normal-mode approximation,<sup>86</sup> we first expand the potential energy in a Taylor series to second (harmonic) order,

$$U(\mathbf{R}) \approx U_{\text{harm}}(\mathbf{R}) = U(\mathbf{R}_{\min}) + \frac{1}{2} \sum_{i,j} \delta R_i \delta R_j \left( \frac{\partial^2 U}{\partial R_i \partial R_j} \right)_{\min}, \quad (8)$$

where  $U(\mathbf{R})$  is the potential energy of the system as a function of all the  $3N$  co-ordinates of each atom, denoted by  $\mathbf{R}$ . The second derivatives are evaluated at the energy minimum, and make up matrix elements  $H_{ij}$  of the hessian matrix  $\mathbf{H}$ . Since heavier atoms move less than lighter ones in each normal mode, it is convenient to introduce mass-weighted displacements with components  $q_i = (m_i)^{1/2} \delta R_i$ , and a mass-weighted hessian matrix with matrix elements  $K_{ij} = U_{ij} / (m_i m_j)^{1/2}$ . In this representation, the potential energy, relative to the minimum energy, can be written to harmonic order as

$$U_{\text{harm}}(\mathbf{q}) - U(\mathbf{R}_{\min}) = \frac{1}{2} \sum_{i,j} q_i K_{ij} q_j. \quad (9)$$

Although all the  $3N$  harmonic oscillators are coupled, since the matrix  $\mathbf{K}$  is symmetric, we can find a similarity transform to diagonalise it,  $\tilde{\mathbf{C}}\mathbf{K}\mathbf{C}$ , where  $\mathbf{C}$  is the matrix of normalised eigenvectors of  $\mathbf{K}$ , and  $\tilde{\mathbf{C}}$  is its transpose. In this ‘normal mode’ representation,  $\tilde{\mathbf{Q}} = \tilde{\mathbf{q}}\mathbf{C}$ , the potential energy is made up of  $3N$  uncoupled harmonic oscillator terms,

$$U_{\text{harm}}(\mathbf{Q}) - U(\mathbf{R}_{\min}) = \frac{1}{2} \sum_i \lambda_i Q_i^2, \quad (10)$$

where  $\lambda_i = \omega_i^2$  are the eigenvalues of the matrix  $\mathbf{K}$ , and  $\omega_i$  are the angular frequencies of the corresponding oscillators. In the Debye-crystal approach, in the first instance we therefore need to compute the normal-mode frequencies. To achieve this, we first equilibrate a system at the conditions of interest [Fig. 1(a), Step 1]. We then minimise the potential energy of an example configuration of the target solid at the correct density for the temperature of interest [Fig. 1(a), Step 2]. We can achieve this using steepest descent, conjugate gradient minimisation<sup>87</sup> or similar optimisation methods.<sup>88</sup> To determine the matrix of second partial derivatives [Fig. 1(a), Step 3], we can use a finite-difference approach; namely, for each pair of degrees of freedom  $i$  and  $j$ , we can estimate the hessian matrix element  $H_{ij}$

by moving the particles in question by some small distance  $\pm\delta$ ,

$$H_{ij} \approx \frac{1}{4\delta^2} [U(R_i + \delta, R_j + \delta) + U(R_i - \delta, R_j - \delta) - U(R_i + \delta, R_j) - U(R_i, R_j + \delta)], \quad (11)$$

where the remaining degrees of freedom are unchanged. Once the full mass-weighted hessian matrix is obtained, we find its eigenvalues  $\lambda_i = \omega_i^2$ , i.e. the squares of the normal-mode angular frequencies. Since each of the normal-mode harmonic oscillators is uncoupled by construction, we can easily find the corresponding partition function and in turn the Helmholtz energy. In particular, apart from the three translational modes with zero eigenvalues, we can obtain the classical harmonic Helmholtz energy of the harmonic crystal by summing over the Helmholtz energies of each normal-mode oscillator [Fig. 1(a), Step 3],

$$\beta A_{\text{harm}} = \beta U(\mathbf{R}_{\min}) + \sum_{i=1}^{3N-3} \ln \frac{\hbar \omega_i}{k_B T}. \quad (12)$$

Finally, we can use the hessian matrix to define the energy function of the Debye crystal, i.e. Eq. (8), and use hamiltonian thermodynamic integration [Eq. 7] to couple it to the overall potential energy [Fig. 1(a), Step 4]. Since the Debye crystal energy is just the harmonic part of the overall potential energy, at least at reasonably low temperatures the difference with the true potential is small, and so the final hamiltonian thermodynamic integration is usually smooth and results in a small anharmonic correction to the free energy. The Debye-crystal route has therefore been suggested to offer an approach that is often numerically better behaved than the Einstein crystal route,<sup>89,90</sup> although the final free energy should in principle be the same with either method.

Once the absolute chemical potential is known at one set of conditions, we can use thermodynamic integration to determine the form of the chemical potential as a function of, for example, the pressure or the temperature. Coexistence points can then be identified by finding where chemical potential curves of different phases cross [Fig. 1(b)].

At sufficiently low temperatures, the entropic part of the chemical potential is always negligible. In principle, if all solid phases of interest were stable at low enough temperatures, one could simply use thermodynamic integration to determine the change in  $\Delta G$  as a function of temperature to determine high-temperature phase behaviour. However, simulations at very low temperature are often slow and require long equilibration times. The harmonic approximation is usually reasonable<sup>80</sup> at somewhat higher temperatures, where equilibration is usually less problematic, and so it is again often possible to use thermodynamic integration alone, without the full Einstein-crystal formalism. This approach however does not work if the phase of interest is not (meta)stable at temperatures at which the approximation is reasonable.

In analytical reference states for both fluid and solid phases, the canonical partition function and hence the chemical potential depend on the de Broglie wavelength once momentum degrees of freedom have been integrated over. In principle, in classical statistical mechanics, the choice of the de Broglie wavelength cannot affect thermodynamic properties, as it shifts the free energy by the same amount across all phases under given conditions. Even though  $\Lambda^2 = h^2 / 2\pi m k_B T$  depends on temperature, it is in practice often chosen to be a

fixed value, such as  $\Lambda = 1 \text{ \AA}$ , for numerical convenience.<sup>18</sup> If this is the case, care must be taken to exclude the translational kinetic energy from other thermodynamic integrations (e.g. in the enthalpy when integrating Eq. (3)) in order to ensure thermodynamic consistency.<sup>12</sup> Finally, it ought to be borne in mind that the above discussion applies to systems governed by classical mechanics. If the momentum degrees of freedom are not factorisable, i.e. in systems where quantum effects may be significant, the translational part of the partition function can no longer straightforwardly be integrated out. The quantum kinetic energy depends on the environment<sup>91–93</sup> and is therefore no longer phase-independent. When exploring the role of nuclear quantum effects on phase behaviour, it is thus often convenient to compute the chemical potentials of the corresponding classical system and then add a quantum correction in a separate step.<sup>91,94</sup>

### C. Estimating the density of states

Instead of using thermodynamic integration to known reference states, it is also possible to estimate the density of states in a Wang–Landau simulation<sup>95–98</sup> and in turn to compute the free energy of a system. In a usual Wang–Landau calculation in the canonical ensemble, the density of states of a given system,  $g(E)$ , is approximated by constructing a histogram with a random walk biased towards visiting previously unsampled states. The initial estimate for the density of states can be  $g(E) = 1$  if no information is known about the system, although a reasonable initial guess can speed up convergence in some cases.<sup>99,100</sup> A Monte Carlo acceptance probability between states 1 and 2 is computed as  $\min[1, g(E_1)/g(E_2)]$ . Each time a state with energy  $E$  is visited, its corresponding density-of-states histogram value is multiplied by a modifier  $f > 1$  so that  $g(E) \leftarrow g(E)f$ , and a histogram of how many times each energy state was visited is incremented by one. The process is repeated until the histogram shows a uniform sampling of relevant states. Since  $g(E)$  changes between Monte Carlo steps, the Wang–Landau algorithm does not obey detailed balance,<sup>99</sup> although if  $f$  is progressively decreased towards unity to minimise the final error, this is not usually a significant concern. The choice of how  $f$  is scaled can significantly affect the algorithm’s performance.<sup>101</sup> Finally, thermodynamic properties such as the Helmholtz energy can be computed as<sup>96</sup>

$$A = -k_B T \ln Q \approx -k_B T \ln \left[ \sum_{\text{bins } i} g(E_i) \exp(-\beta E_i) \right], \quad (13)$$

where  $Q$  is the canonical partition function that we approximate by summing over the bins of the histogram. All other thermodynamic quantities can then be computed from the Helmholtz energy. A similar approach can be used in the grand ensemble.<sup>98</sup>

The Wang–Landau sampling approach and its analogues have been very successful in modelling the behaviour of lattice models such as the Ising model,<sup>96,97</sup> the XY model<sup>102</sup> and even lattice-based liquid crystals<sup>103</sup> and alloys,<sup>104</sup> as well as in modelling the phase behaviour of fluids<sup>98,105,106</sup> and solutions.<sup>107</sup> In lattice models with discrete energy levels, the minimum and maximum energy to consider is usually well defined. By contrast, in systems with continuous energy levels, the Wang–Landau approach of binning energies into histograms

and imposing a minimum and maximum energy to consider can lead to inefficiencies and errors; systematic discretisation methods<sup>108</sup> have been proposed to address this problem.

Moreover, since solid-state densities of states are usually much narrower than those of liquid phases, approaches based on densities of states are likely to find supercooled liquid structures instead of solid phases in a random walk.<sup>109</sup> It has been proposed that starting from low-energy structures and partitioning the energy space upwards permits the calculation of the density of states of known solid structures.<sup>110</sup> It is also possible to bias the sampling of solid states with a suitable order parameter, for example using umbrella sampling.<sup>111</sup>

Another strategy for estimating the density of states is known as nested sampling.<sup>112–114</sup> The approach relies on estimating the degeneracy of a given energy level by generating a pool of  $K$  random configurations, identifying the configuration with the highest energy, and replacing it with a new random configuration that has a lower energy. In each step of the procedure, we can estimate, relative to an unimportant origin, the volume of phase space at the maximum energy as  $\Phi(i) = [K/(K+1)]^i$ ,<sup>114</sup> where  $i$  is the iteration number. The density of states can then be computed as  $g(i) \approx \Phi(i) - \Phi(i-1)$ . This approach does not require the binning of energies; however, its ability to find solid phases is similarly limited as the Wang–Landau-based approaches mentioned above. In particular, a random sampling of position co-ordinates corresponds effectively to sampling at a very high temperature, where the large entropy of vapour and liquid states completely dominates the system’s behaviour, and it is very unlikely, unless a very large number of configurations is sampled, that potential energy wells of solid phases will be adequately represented. However, such states are usually possible to obtain by energy minimisation; a way of computing the density of states that captures solid phases that dominate the phase behaviour at low temperatures entails a combination of optimisation and nested sampling.<sup>115,116</sup>

### D. Direct coexistence

As an alternative to computing ‘absolute’ chemical potentials relative to known analytical models, we can instead determine phase coexistence analogously to its experimental determination: by simulating the two coexisting phases explicitly in a ‘direct-coexistence’ simulation.<sup>117,118</sup> In this approach, we typically simulate a system in a slab geometry in an elongated periodic box with two explicit interfaces between the phases in question, one on each side of the two coexisting bulk phases. We usually run simulations at a given pressure over a range of temperatures and then bracket regions where a given phase shrinks or grows to determine the coexistence temperature, although simulations in other ensembles (e.g. microcanonical or canonical) are also possible.<sup>119</sup> This approach was introduced in the 1970s as a way of investigating both the interfacial structure<sup>117</sup> and the thermodynamic behaviour<sup>118</sup> of Lennard-Jones (LJ) particles. As computational power has increased and much larger system sizes could be simulated, the method has been used in progressively more contexts since, from hard spheres<sup>120–122</sup> to metals<sup>123,124</sup> to water.<sup>125,126</sup>

The main advantage of direct-coexistence simulations is their relative simplicity: often, they are very easy to implement, since in principle the system is just evolved using normal MC or molecular dynamics (MD) algorithms. However, there are

also several possible disadvantages. The simulation relies on the stable phase growing rapidly at the expense of the unstable phase when away from their coexistence temperature. If the two phases in question have very similar chemical-potential gradients, the driving force for the phase transition may not be sufficient to enable the coexistence point to be bracketed accurately. Moreover, even if the driving force is significant but the phase transition is dynamically slow, direct-coexistence methods are unlikely to be productive.

There are further practical considerations that may limit the applicability of direct-coexistence simulations. When computing chemical potentials directly, pure phases are simulated, and relatively small system sizes are often sufficient to obtain accurate results. By contrast, in direct-coexistence simulations, an explicit interface forms between the bulk phases. Although in the thermodynamic limit, the role of the interface is immaterial, typical simulation-box sizes are much smaller and the proportion of particles at the interface is non-negligible. The interface can thus dominate the system's apparent behaviour. With large enough system sizes, such finite-size effects in theory disappear, and so probing the behaviour as a function of simulation-box size is especially important in such simulations. In practice, even if the system is in principle large enough to simulate coexistence, large blocks of the two coexisting phases must still be brought into direct contact and the interface properly equilibrated at each set of conditions being considered, which can be a laborious process. The larger the system size, the longer such initial equilibration will take, while for smaller systems, an unfortunate initial, pre-equilibration choice of interface structuring may result in such large restructuring of the coexisting phases that they disappear altogether even if they are not far from coexistence. When explicit interfaces are present, particular care must be exercised in interpreting simulation results when using truncated potentials.<sup>14,127</sup>

When one of the phases in question is a solid, further difficulties arise, since the simulation box must be compatible with the unit cell at the conditions studied to prevent placing stress on the solid.<sup>128–130</sup> Practically, this can be addressed by (i) determining the crystal's unit cell parameters as a function of pressure and temperature, (ii) scaling the box with an interface to these dimensions at the pressure and temperature studied, and (iii) applying a barostat only in the direction orthogonal to the interface.<sup>18</sup> This approach relies on the coexisting phase being able to adapt to the box shape and size consistent with the solid phase, which in general will only be possible with a fluid phase. However, this does not preclude direct-coexistence simulations from being used in determining solid–solid coexistence lines as long as both solids can coexist with a fluid phase, even if they are only metastable.<sup>126</sup>

Despite some of these potential complications with using direct-coexistence simulations in determining phase diagrams, direct-coexistence simulations have especially recently been very popular for studying what is effectively vapour–liquid coexistence in biomolecular mixtures.<sup>26–28,131–135</sup> In such systems, proteins and nucleic acids are typically modelled using fast, coarse-grained potentials with implicit solvents, enabling any interfaces to equilibrate readily and coexistence properties to be determined easily. Typically, such simulations are performed in the canonical ensemble at a fixed overall density held somewhere in the range where phase separation occurs; this typically corresponds to very low coexistence pressures.<sup>14</sup> Since phase diagrams are increasingly being reported in the

experimental literature,<sup>136</sup> the computation of phase diagrams, primarily using direct coexistence, has been instrumental in both designing and benchmarking numerous coarse-grained potentials<sup>26,28,137</sup> which have subsequently been used to probe the fundamentals of protein phase behaviour, including the effects of salt concentration,<sup>138</sup> valency,<sup>134,139–141</sup> composition (of both the individual proteins and of the entire mixture) and patterning/distribution of amino acids,<sup>140,142,143</sup> as well as evolving protein sequences to promote or inhibit their capacity to phase-separate.<sup>133,135</sup> As such models evolve and account for more complex systems, it may well be that the limitations of direct-coexistence simulations become more apparent, especially in the light of large systems sizes that may need to be simulated to avoid artefacts. Other free-energy methods, such as thermodynamic integration, may thus become more broadly used in this field in due course.

## E. Interface pinning

In direct-coexistence simulations, the coexistence point is usually determined by varying the thermodynamic parameters until one phase grows at the expense of the other. However, rather than needing to rely on potentially slow equilibration under different conditions, we can ensure that the system remains in a two-phase state by applying a bias to the system in an approach known as interface pinning.<sup>144–146</sup> The bias applied is related to the free-energy difference between the two phases.

An order parameter must first be identified that can distinguish between the two phases of interest. Often, this is achieved by probing the local environment around each particle, for example with spherical-harmonic functions.<sup>147–153</sup> The spherical harmonics form a complete orthonormal system and any real function on the unit sphere – such as the particle neighbour density – can be expressed using the spherical harmonics as a basis set in a Laplace series.<sup>154</sup> It can often be computationally more efficient to use real spherical harmonics instead of their more usual complex analogues.<sup>155</sup> In typical simulation studies, since order-parameter calculations often need to be repeated many times, typically only one term in the series is picked that can best distinguish between the phases of interest<sup>147</sup> to help reduce the computational expense. Other order parameters may also be suitable,<sup>156–158</sup> depending on the system studied, and for simple solid–fluid interfaces, a simple density-field approach is often sufficient.<sup>145</sup>

As long as an order parameter has been found that can distinguish between two phases, say  $\alpha$  and  $\beta$ , we can define the total number of particles as  $N = N_\alpha + N_\beta + N_{\text{int}}$ , where  $N_{\text{int}}$  is the number of particles at the interface. In a slab geometry, like with direct-coexistence simulations, the interfacial area is minimised by forming two planar interfaces across the smallest box dimension, and as long as the interface remains planar, the number of interfacial particles and the resulting interfacial contribution to the free energy should not change as the interface moves. The total Gibbs energy is thus given by  $G = N_\alpha \mu_\alpha + N_\beta \mu_\beta + G_{\text{int}}$  or equivalently  $G = N_\alpha \Delta\mu + \text{constant}$ , where  $\Delta\mu = \mu_\alpha - \mu_\beta$  and the constant accounts for terms that do not depend on  $N_\alpha$ . The chemical potential difference  $\Delta\mu$  is the quantity we wish to determine; it is zero at coexistence, while for non-zero values the sign tells us which phase,  $\alpha$  or  $\beta$ , is thermodynamically more stable.

To compute  $\Delta\mu$ , we can apply an additional bias potential to the system that depends on the global order parameter of the simulation box. In the following, we will use as an example  $N_\alpha$ , the number of particles in phase  $\alpha$ , as the order parameter for simplicity; however, it is not actually necessary for us to be able to determine the phase of each molecule individually: a global order parameter that changes between the bulk phase  $\alpha$  and phase  $\beta$  is sufficient.<sup>144,145</sup>

One possibility is to compute the free energy directly using a method such as umbrella sampling.<sup>111</sup> If we apply a bias to the global order parameter such that we force the number of particles of phase  $\alpha$  to change (i.e. the interface moves in a direction we bias it towards) and repeat the procedure at different target values of  $N_\alpha$ , and then account for the bias, we find  $G(N_\alpha)$  relative to an arbitrary origin for a range of  $N_\alpha$ , which permits us to determine  $\Delta\mu$  as the gradient. Alternatively, suppose the bias is of the harmonic form  $U_{\text{bias}}(N_\alpha) = (\kappa/2)[N_\alpha - a]^2$ , with  $a$  chosen to be close to the number of particles in phase  $\alpha$  at the start of the simulation, and  $\kappa$  a tuneable parameter. The associated Gibbs energy is  $G_{\text{biased}} = G(N_\alpha) + U_{\text{bias}}(N_\alpha) = (\kappa/2)[N_\alpha - a + \Delta\mu/\kappa]^2 + \text{constant}$ , where we have completed the square and ignored terms independent of  $N_\alpha$ .<sup>144,145</sup> The Boltzmann probability for observing a particular number of particles in phase  $\alpha$  is  $P(N_\alpha) \propto \exp\{-\beta G_{\text{biased}}\}$ . Assuming that  $\kappa$  is chosen to be sufficiently large so that the interface does not move appreciably, and that the number of particles is sufficiently large to treat it as a continuous variable, we can normalise this probability by extending the integration limits in  $N_\alpha$  to  $\pm\infty$  without introducing significant error. The probability  $P(N_\alpha)$  then becomes of a standard gaussian form. Finally, we can find

$$\langle N_\alpha \rangle = \int_{-\infty}^{\infty} N_\alpha P(N_\alpha) dN_\alpha = a - \Delta\mu/\kappa. \quad (14)$$

We can therefore find  $\Delta\mu$  simply by determining  $\langle N_\alpha \rangle$  in a simulation with a bias.<sup>144,145</sup>

By changing the thermodynamic conditions and repeating this calculation, we can determine the point at which  $\Delta\mu = 0$ , i.e. where the phases  $\alpha$  and  $\beta$  coexist. Since the two-phase system is stable if a sufficient bias is applied, such simulations are at equilibrium even when the phases are not at coexistence, and are therefore less sensitive to how they are initialised than direct-coexistence simulations. However, as with direct-coexistence simulations, since an explicit interface is present in the simulation box, the method usually requires considerably larger system sizes than approaches based on thermodynamic integration, and finite-size effects must be investigated carefully.<sup>144</sup> Similarly, because of the explicit interface required, it is not usually possible to use this approach to study solid–solid phase coexistence directly.

## F. Determining coexistence curves

Coexistence curves on a phase diagram correspond to loci of points where two (or more) phases have the same temperatures, pressures and chemical potentials [see Fig. 1(b)(iii)]. We can determine a single coexistence point using the methods outlined so far, either by explicitly evaluating chemical potentials or by using a proxy method such as direct coexistence. Other points may then be determined by repeating the procedure at different initial conditions. Alternatively, the coexistence curve

can itself be integrated. Perhaps the simplest approach is to run a series of simulations for a pair of pure phases that are known to coexist at a given point and, starting from this known point, numerically integrate the Clapeyron equation,

$$\frac{dP}{dT} = \frac{\Delta H}{T\Delta V}. \quad (15)$$

This is often known as Gibbs–Duhem integration,<sup>18,159,160</sup> although, potentially confusingly, it entails integrating the Clapeyron equation rather than the Gibbs–Duhem equation. Such an approach is often reasonable for integrating over relatively small ranges of temperature or pressure, but there is no inbuilt error-checking mechanism: small errors in the enthalpy or volume change are cumulative. Other approaches, such as multistate reweighting, are far more accurate.<sup>161–163</sup> Although their implementation may be more challenging, several open-source tools exist that make their use straightforward.

## G. Analysing errors

Since there are many potential sources of error in free-energy calculations, both numerical and systematic, it can be difficult to determine and propagate errors.<sup>14,18</sup> It is usual in such calculations to employ a posteriori consistency checks<sup>12,18,41,89</sup> to ensure thermodynamic consistency. This is especially important if coexistence curves are directly integrated with methods like Gibbs–Duhem integration, but consistency checks are a valuable tool for finding implementation errors in all cases. For example, if one performs an Einstein-crystal calculation at two different pressures and temperatures for a given phase, it should be possible to integrate the chemical potential from one of these points to the other along many different isotherms and isobars, and all such reversible paths should give the same chemical potential (within numerical accuracy) as the Einstein-crystal calculation itself. At least a few coexistence points should also ideally be checked to ensure that coexistence temperatures and pressures determined from chemical-potential crossovers can be reproduced in direct-coexistence simulations. Of course it is not impossible to estimate errors in the chemical potential itself; for example, we can compute the chemical potential in several completely independent simulations and determine the corresponding standard deviation,<sup>8</sup> or use polynomial fitting to isotherm or isobar data to give prediction bands for the chemical potential.<sup>9</sup> However, such error estimation cannot account for any systematic errors, and consistency checks are usually rather more revealing than numerical error analysis.

## H. Capturing long-timescale entropies

One significant potential difficulty with calculating chemical potentials by thermodynamic integration from known reference states is that some features of phases may be impossible to equilibrate at computationally accessible time- and length-scales. A well-studied example of this is proton disorder in ice. Many ice phases exist as proton-ordered and proton-disordered analogues,<sup>8,164–170</sup> with the former dominating at low temperatures where the stable phase is determined primarily by low enthalpy, and the latter dominating at high temperatures at which the higher entropic favourability of (partially) disordered

659 phases takes over. However, different manifestations of proton  
 660 disorder, which contribute to the overall entropy of such phases,  
 661 are not readily accessible to computer simulations. Moreover,  
 662 if a particular proton-disordered configuration that is not fully  
 663 representative of the phase in question is chosen, this can signi-  
 664 ficantly affect the calculated phase behaviour; the calculated  
 665 phase diagram of such phases can easily be wrong without any  
 666 clear signs that anything is amiss.<sup>126</sup> Although in principle  
 667 a direct-coexistence simulation should result in correct phase  
 668 behaviour, in practice the proton disorder in the crystalline  
 669 phases is largely locked in at all computationally accessible  
 670 time scales. In calculations of water phase diagrams, proton  
 671 disorder is therefore usually accounted for by assuming that  
 672 the experimentally determined proton disorder is correct<sup>126,171</sup>  
 673 and adding a suitable analytically computed residual entropy to  
 674 the chemical potential.

675 Another example of a system whose entropy is difficult to  
 676 capture are quasicrystals, which have no long-ranged transla-  
 677 tional order, and so cannot be described by a periodic lattice,  
 678 but which often have long-ranged orientational ordering.<sup>172</sup>  
 679 They are often thought of as projections of a higher-dimensional  
 680 crystal into a lower dimension. Quasicrystals often, though  
 681 not always,<sup>173</sup> arise because of a competition between multiple  
 682 lengthscales, either in multicomponent mixtures of differently  
 683 sized particles or in single-component mixtures with different  
 684 intrinsic lengthscales<sup>174,175</sup> or in systems with explicit orienta-  
 685 tional order.<sup>77,176,177</sup> A knowledge of the phase behaviour of  
 686 quasicrystals from simple models can help guide experimental  
 687 design; for example, computer simulations of quasicrystals  
 688 made from DNA nanostar motifs<sup>178</sup> have in part led to such  
 689 soft binary quasicrystals being realised experimentally.<sup>179</sup>

690 When it comes to determining a quasicrystal’s stability,  
 691 the configurational disorder clearly entropically stabilises the  
 692 quasicrystal, but is not straightforward to determine directly.  
 693 Although simple models sometimes have such fast dynamics  
 694 that an equation of state that accounts for bulk quasicrystalline  
 695 phase behaviour can readily be determined in brute-force simu-  
 696 lations,<sup>180</sup> in other systems, the timescale at which quasicrys-  
 697 talline structures can rearrange themselves can vastly exceed  
 698 the simulation time available. To determine thermodynamic  
 699 stability, sometimes the free energy of an approximant crys-  
 700 tal – i.e. a periodic crystal with a large unit cell featuring  
 701 motifs identified in the quasicrystalline phase of interest – is  
 702 computed.<sup>181–183</sup> Although the quasicrystal’s configurational  
 703 entropy is not captured,<sup>184</sup> it can be subsequently added using a  
 704 random-tiling approximation.<sup>185,186</sup> Another possible approach  
 705 is to use the thermodynamic integration formalism introduced  
 706 above to compute the free energy of a particular quasicrystal  
 707 configuration and account for the configurational entropy using  
 708 an approximation of uncorrelated phason flips.<sup>187,188</sup> However,  
 709 in simulations of systems with intermediately fast dynamics, a  
 710 convenient approach would again be to use a direct-coexistence  
 711 simulation with an explicit interface with the fluid phase, in  
 712 a system that is sufficiently large that the free-energy cost of  
 713 defects arising from the periodic box are negligible. Once one  
 714 point of coexistence is known, the chemical potential of the  
 715 quasicrystal must by construction be equal to that of the fluid  
 716 phase, which can easily be determined using thermodynamic  
 717 integration from an ideal gas. The chemical potential of the  
 718 quasicrystal can then be integrated along isotherms or isobars  
 719 to other conditions of interest, and the full phase diagram can  
 720 be determined this way.<sup>77</sup> Such an approach is only feasible if

721 simulations are fast enough to be tractable even with systems  
 722 so large that the quasicrystal–fluid equilibrium is dominated by  
 723 bulk terms, which should be explicitly checked. More broadly,  
 724 a method can work very well for one system, whilst being  
 725 inappropriate to use for another, highlighting the importance  
 726 of having a range of tools at our disposal when investigating  
 727 different systems.

### 728 III. MACHINE-LEARNING APPROACHES TO PREDICTING 729 PHASE-DIAGRAMS

730 As an alternative to free-energy calculations, machine-  
 731 learning methods have been gaining traction in recent years  
 732 due to their potential for high-accuracy predictions with low  
 733 computational cost. Machine-learning models can often make  
 734 accurate predictions quickly once trained on sufficient data, and  
 735 can moreover usually be updated or retrained relatively easily  
 736 when more data become available. However, on the flip side,  
 737 the process of developing and training such models can often  
 738 be complex and time consuming. One first has to decide on  
 739 the choice of model architecture, as well as the type and source  
 740 of data to use to train the model. These data need to be either  
 741 generated or extracted from repositories and then processed  
 742 into a suitable and consistent format to be used as training  
 743 inputs for the model. To this end, suitable descriptors that in  
 744 some way quantify features relating to the properties of interest  
 745 must be identified. Finally, the training process itself requires  
 746 considerations such as hyperparameter tuning, cross-validation  
 747 and regularisation techniques to prevent overfitting, as well as  
 748 benchmarking of the model with other methods. The appro-  
 749 priate route to follow in the case of machine-learned methods  
 750 is therefore often somewhat less clear than is the case with  
 751 traditional methods.

752 Here, we provide a brief summary of some machine-learning  
 753 methods that have recently been applied to predict phase dia-  
 754 grams.

#### 755 A. Machine-learned potentials

756 One possible route that fruitfully exploits such techniques is  
 757 to use machine-learned potentials (MLPs)<sup>6–9,45–48,190–193</sup> that  
 758 can simulate systems with the accuracy of ab initio methods  
 759 but at a fraction of the computational cost [Fig. 2(a)]. Such  
 760 approaches permit the determination of phase diagrams at the  
 761 level of the underlying electronic structure theory; indeed, des-  
 762 pite some possible limitations of machine-learned approaches in  
 763 quantifying longer-ranged interactions,<sup>45,190</sup> it has been shown  
 764 that the phase diagrams computed with MLPs are about as  
 765 different from the underlying DFT phase diagrams as different  
 766 DFT functionals are to one another.<sup>8</sup> When coupled with MLPs  
 767 to make the computations tractable,<sup>45</sup> thermodynamic phase  
 768 behaviour could thus be used to obtain better approximate DFT  
 769 functionals and to understand the nature of the approximations  
 770 made in such functionals.

771 One advantage of such MLPs is that they provide an intuitive  
 772 molecular-level description of the material in question, as well  
 773 as likely mechanisms by which phase transformations can occur.  
 774 However, the full machinery outlined so far in this article is  
 775 required to determine the phase behaviour of such potentials.

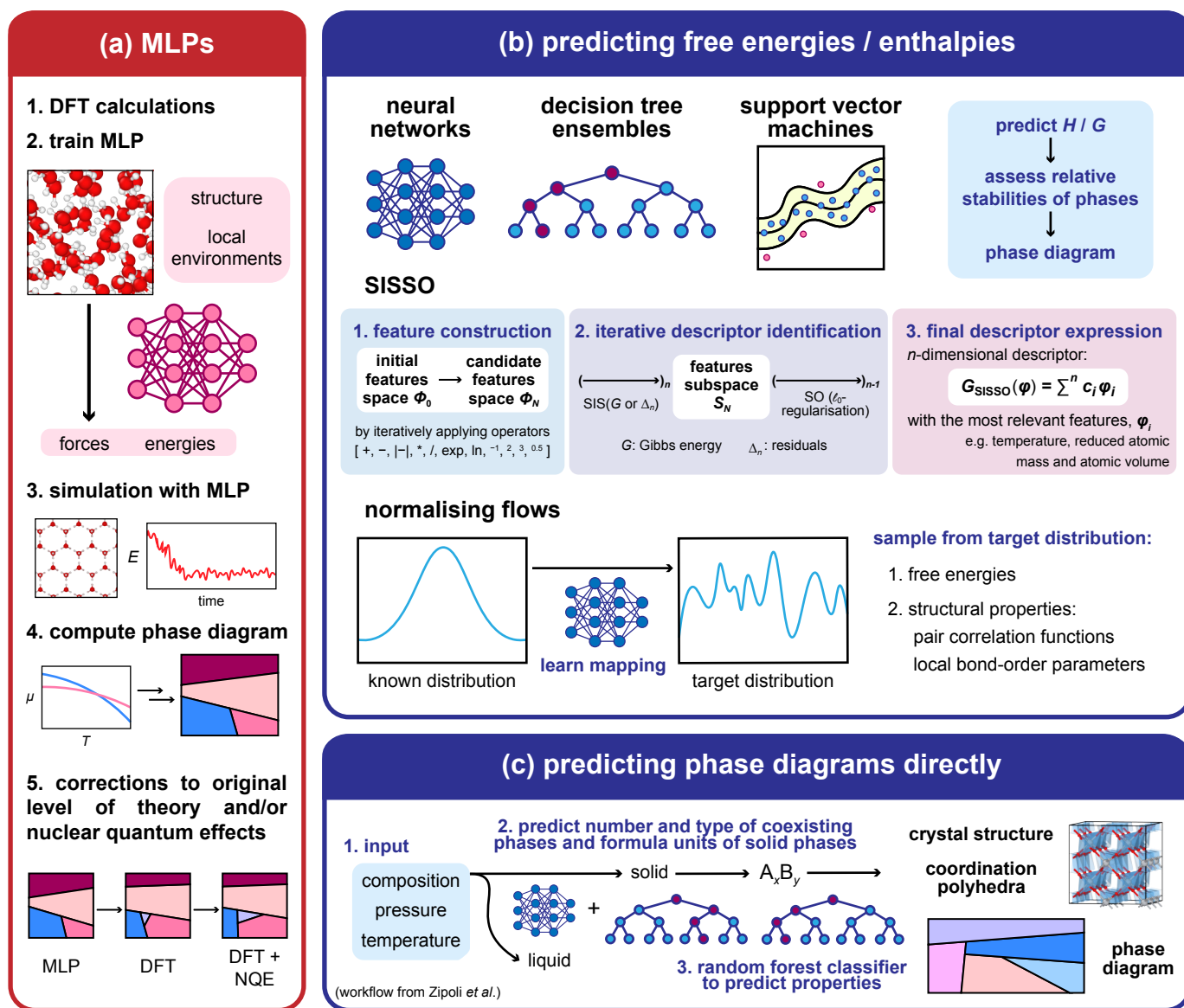


Figure 2. Summary of different machine-learning approaches to determining phase diagrams. (a) Machine-learned potentials (MLPs). (b) Predicting the enthalpies or free energies of different phases, which then allows one to assess their relative stabilities at different conditions (e.g. via a convex hull construction) to construct the phase diagram. This can be done via a variety of methods, ranging from more common machine-learning algorithms like neural networks and decision tree ensembles, to the sure-independence-screening sparsifying operator (SISSO) approach and normalising flow models. (c) Obtaining phase diagrams by directly predicting the number and types of coexisting phases at a given condition. The workflow presented here is from Zipoli and co-workers.<sup>189</sup>

## 776 B. Predicting enthalpies and free energies

777 Various machine-learning methods have also been used to  
 778 predict phase diagrams and associated thermodynamic proper-  
 779 ties of materials directly, without requiring expensive molecular  
 780 simulations. The success of such approaches depends heavily  
 781 on the availability and quality of data of the relevant proper-  
 782 ties and existing phase diagrams for a number of materials  
 783 to train the models. The amount of experimental and DFT  
 784 data in various materials repositories that have become widely  
 785 available in recent years has been instrumental in the develop-  
 786 ment of these data-driven approaches using machine-learning  
 787 methods. These include the Materials Project<sup>194</sup> or the Open  
 788 Quantum Materials Database<sup>195</sup> for DFT-calculated quantities  
 789 like the formation enthalpy, the Inorganic Crystal Structure

790 Database,<sup>196</sup> the Cambridge Structural Database or Pearson's  
 791 Crystal Data for crystal structure data, or the NIST 31 database  
 792 for experimental and computational phase diagrams.

793 One approach to constructing phase diagrams is to predict the  
 794 relevant thermodynamic quantities of different phases directly  
 795 from their composition, and sometimes structure, and assess the  
 796 relative stability of the phases using these quantities [Fig. 2(b)].  
 797 Such an approach has, for example, been used in materials  
 798 discovery to identify if materials with certain elemental com-  
 799 positions are synthetically accessible by developing models  
 800 to predict their formation enthalpies.<sup>197</sup> Models are trained  
 801 on existing data of thermodynamic quantities obtained either  
 802 experimentally or from DFT calculations, and then extrapolated  
 803 to make predictions for new materials. The stability of a certain  
 804 material relative to others with similar elemental compositions

can then be assessed with a convex hull construction.<sup>198</sup> Although differences between their relative entropies have often been approximated to be negligible when dealing with solid–solid transitions, models have also been trained to predict Gibbs energies.<sup>63,199</sup> For example, neural networks have frequently been used to predict formation enthalpies<sup>197,200,201</sup> and free energies<sup>63</sup> of materials. In general, neural networks are made up of interconnected nodes in a layered structure: the input layer first takes in and processes the input information before passing it on to the next layer, and the subsequent hidden layers analyse and process the output from the previous layer before it reaches the output layer, where the final ‘result’ is given. Connections between nodes are usually associated with a given weight and bias, and during the training process the neural network learns to adjust these parameters continuously via feedback loops to minimise the discrepancy between the predicted and actual quantity to improve predictions. Neural networks are thus able to learn complex relationships between the input data and the target quantity of interest, and they have the ability to make very accurate predictions, especially when trained on large amounts of data. Other models including decision tree ensembles<sup>202,203</sup> and support vector regressors<sup>204–206</sup> have also been used to predict thermodynamic quantities. Such machine-learning approaches have been shown to be able to predict the formation enthalpies of materials with an accuracy comparable to DFT; however, predicting relative stabilities of related compounds from their predicted formation enthalpies using convex hull constructions is less accurate, as the advantageous systematic error cancellation in DFT predictions does not apply to these models.<sup>207</sup> Additionally, the majority of such models only use compositional information and hence cannot make predictions for polymorphic transitions. Models that include both compositional and structural descriptors have been shown to produce more accurate stability predictions;<sup>197,207</sup> however, the structure is not always reported in the underlying data set.

One common issue with machine-learning models is their interpretability. For example, in neural network models, as the input data are passed through multiple hidden layers and processed, the predictions often become difficult to understand at a physical level and we are often not able to get a sense of what kinds of features make an important contribution to the accuracy of the predictions, and hence have physical significance.<sup>197</sup> In order to gain a better understanding of what properties are the most relevant for predicting thermodynamic quantities, we could make use of alternative methods which are explicitly able to select features to construct descriptors that result in the best predictions. An example is the sure-independence-screening sparsifying operator (SISSO)<sup>208</sup> method, which searches a space of mathematical expressions of selected features to find an optimal solution for an accurate descriptor of the quantity we want to predict. In this approach, to obtain a range of candidate descriptors, a combination of different mathematical operations is first recursively applied to an initial feature space consisting of properties that are potentially relevant for capturing the predicted quantity, while ensuring that only sensible combinations with physical units are permitted. From the set of candidate descriptors, a subspace of the best-performing features with the highest correlation with the target property can then be selected. Within the selected subspace,  $l_0$ -norm regularisation (or similar approximations) can be used to identify the best one-dimensional descriptor, which is then used to predict the

target property and the associated residuals for the training set. This descriptor identification process can then be repeated to consider higher-dimensional solutions with each additional iteration, with the residuals from the previous iteration as the new target property. The goal is to identify the lowest-dimensional solution with acceptable errors below a desired threshold, and the features which are ultimately selected in the descriptor would be those that are the most important for capturing the property of interest.

The SISSO approach has been used by Bartel and co-workers to derive an accurate descriptor of the Gibbs energy of inorganic crystalline solids.<sup>199</sup> As a starting point for the initial feature space, they consider several quantities that are potentially relevant for predicting the Gibbs energy of a material. These include the atomic volume and band gap derived from the Materials Project database, experimental formation enthalpy and temperature, as well as tabulated elemental properties (ionisation energy, electron affinity, covalent radius, electronegativity and atomic mass), which were then transformed into compound-specific versions by finding suitable stoichiometrically weighted arithmetic and geometric means. Despite the larger set of properties being considered as input features for the model, the SISSO-learned descriptor depends only on temperature, reduced atomic mass and atomic volume, hence showing that these properties are the most significant in predicting the Gibbs energy. Their descriptor is benchmarked against Gibbs energies calculated via the quasiharmonic approximation (i.e. assuming harmonic normal modes [cf. Sec. II B 3], but with a changeable lattice constant) and then used to predict reaction energetics and equilibrium product distributions, as well as to assess temperature-dependent stability via a convex hull construction using the Gibbs formation energies of the materials. However, since this descriptor only predicts the free energies for solid phases, it is unable to determine the melting point of solid–liquid transitions directly. Transitions between different polymorphs are also beyond the scope of this predictor, since the only input feature that describes structure is the atomic volume, and errors from the descriptor are typically larger than the free-energy differences between polymorphs.<sup>199</sup>

In addition to training models to predict energies or free energies directly, machine-learning methods can also allow us to develop models to sample complex probability distributions, such as the Boltzmann distribution of a given system with some known potential energy function, from which we can sample to compute equilibrium energies and structures. Sampling from such distributions is traditionally done using methods like MC and MD, which can become expensive for complicated systems or when ergodicity is locally broken and rare events need to be simulated. Recently, alternative methods like flow-based sampling have been applied to various systems.<sup>209–212</sup> A normalising flow model<sup>213</sup> transforms an analytically tractable base distribution, such as a gaussian distribution, which we can easily sample from and evaluate probability densities for, into a more complex target distribution of the system of interest, through a series of invertible and differentiable functions. We can sample from the target distribution  $\rho_x$  by first sampling from the initial known base distribution  $\rho_z$ , and then applying this series of functions, i.e. taking  $x = f(z)$ . The probability density of such a transformed sample can be calculated using the product of the density of the original sample under the base distribution  $\rho_z$  and the associated change in volume from the sequence of transformations, which is the product of the

absolute values of the jacobian determinants for each inverse transformation, as given in the change of variable formula

$$\rho_x(x) = \rho_z(z) \left| \det J_f(z) \right|^{-1}, \quad (16)$$

where  $J_f$  is the jacobian matrix of  $f$ . In theory, as long as the two distributions are of the same dimension, any target distribution  $\rho_x$  can be obtained from any base distribution  $\rho_z$  given a sufficiently complex transformation function  $f$ . Such a function can comprise of a sequence of  $N$  invertible functions,

$$f = f_N \circ f_{N-1} \circ \dots \circ f_1; \quad (17)$$

this function composition successively constructs more complex functions from the previous one, and the base distribution can be said to ‘flow’ through the series of functions.

Deep-learning algorithms can help us find a suitable composition of functions which can transform the base distribution into the target distribution. This approach has been used by Wirnsberger and co-workers to estimate the free energy and structural properties of atomic solids from a flow-based model.<sup>210,214</sup>

Here, the flow-based model is trained to approximate the Boltzmann distribution of various atomic solids, starting from a base distribution of a lattice model with spherically truncated gaussian noise added to the atom at each lattice site, followed by a random permutation of all atoms. In this manner, we can target specific crystal structures (such as cubic and hexagonal ice) by using an appropriate choice of lattice for the base distribution to guide the model towards the phase of interest, without changing the potential energy function or using ground-truth samples. To transform the base distribution into the target distribution, each  $f_n$  in the series of functions transforms element-wise one or two co-ordinates of all the atoms as a function of the remaining atom co-ordinates, and each  $f_n$  is parameterised by a separate neural network whose parameters are trained by minimising the difference between the target and base distribution (as calculated by the Kullback–Leibler divergence<sup>215</sup> as the loss function). The trained flow-based model is able to reproduce structural properties and free-energy estimates with good accuracy: energy histograms, pair correlation functions and local bond-order parameters computed from sampling from the trained model for the truncated and shifted LJ face-centred cubic phase and for ice I using the mW water model were consistent with traditional methods, and free-energy estimates obtained using ‘learned’ free-energy perturbation methods with the trained model were also shown to be in good agreement with multistate-reweighting methods for both ice and LJ systems. Unlike traditional methods like thermodynamic integration or multistate reweighting, where samples from intermediate thermodynamic states are needed, this flow-based approach to obtain free energies does not require sampling from the target distribution for training the model or computing free-energy estimates. However, although once trained, generating samples and computing probability densities and various estimates from the trained model is efficient since samples can be obtained in parallel unlike in trajectory-based methods, one downside to this approach is that constructing and training the model is difficult and computationally expensive.<sup>210</sup>

### C. Predicting phase diagrams directly

The methods discussed above involve first predicting thermodynamic quantities like free energies from the trained models,

and the energies can then be used to assess the relative stabilities of the different phases considered to obtain a prediction for the phase diagram. An alternative approach may be to use machine-learning models to construct phase diagrams directly by training the models to predict the number and types of coexisting phases at a given condition, with or without predicting, estimating or computing free energies of the phases as an intermediate step.<sup>63,189,216</sup> This approach has recently been adopted by Zipoli and co-workers, who used machine-learning models to predict the polymorphs that are thermodynamically stable at certain conditions<sup>189</sup> [Fig. 2(c)]. For a certain condition defined by the overall chemical composition, temperature and pressure, we can first predict the types of phases (i.e. solid, liquid or vapour) that are present and the formula units of any solid phases using neural networks or random forest classifiers and regressors trained on data extracted from experimental and computational phase diagrams. Using the chemical formulae of the solid phases obtained from the first step, we can then use classifiers trained on crystal structure data of inorganic crystalline materials to predict structural properties of the solid phases in question (e.g. Bravais lattice types, crystal system, space group, structure type and local atomic environments characterised by the co-ordination polyhedra present). Such a data-driven approach with the workflow proposed by Zipoli and co-workers in principle has the potential to be able to build up phase diagrams directly and relatively quickly, using just the overall elemental composition, temperature and pressure to make predictions from the trained model, without any sort of feature construction required. However, the different components of the overall workflow may have varying levels of success: for example, even though structural properties of a particular solid phase can be predicted using classifiers from their chemical composition with good levels of accuracy,<sup>189,217</sup> obtaining the formula units of the different phases from the overall composition of the mixture when there are multiple phases present has proved to be much more difficult.<sup>189</sup>

Although machine-learning models for predicting phase behaviour and thermodynamic quantities have mostly focussed on inorganic crystalline solids, such approaches have also been used on soft materials, for example to predict the melting temperatures and the ternary phase diagrams of lipid mixtures, by predicting the types of coexisting phases from the mole fraction and properties of the lipids present.<sup>218</sup> Other recent approaches in predicting phase behaviour for biomolecules include sequence-based predictions of the propensity for proteins to phase separate.<sup>219–222</sup> In addition to the protein sequence, some models also take experimental conditions (protein concentration, salt concentration, temperature, pH and presence of crowding agents) into account when making predictions of phase-separation propensity;<sup>223</sup> however, such models are in general not able to predict quantitative measures such as the critical temperature or the physicochemical properties of the coexisting phases directly.

## IV. FUTURE OUTLOOK

As evidenced by the large amount of work done in this field, the computational prediction of phase diagrams has already been very successful. Considerable effort has gone into developing a range of methods for studying phase coexistence, and even more to applications to a wide range of systems, from

1043 inorganic materials to biological systems. A knowledge of 1105  
 1044 the thermodynamic phase behaviour of a given computational 1106  
 1045 model is useful both on the applied side, for example for 1107  
 1046 understanding the underlying physics of a given system that 1108  
 1047 may be difficult to study experimentally, and from the point 1109  
 1048 of view of model development and refinement. The various 1110  
 1049 methods outlined in this paper all have their advantages, but 1111  
 1050 also some drawbacks.

1051 For example, some machine-learning methods are able to pro- 1113  
 1052 duce predicted phase diagrams quickly, and so may be useful in 1114  
 1053 guiding experiments towards promising phases with interesting 1115  
 1054 properties. However, the ‘best’ approach in terms of the model 1116  
 1055 architecture and training process is often problem-specific and 1117  
 1056 might not be straightforward to decide on a priori. Indeed, 1118  
 1057 often several models are first trained and the performance of 1119  
 1058 the different models is then used to decide what kind of model 1120  
 1059 architecture might work better for this one particular context 1121  
 1060 considered. Despite the large amount of data available, the 1122  
 1061 accuracy of machine-learning-based models is also limited by 1123  
 1062 the quality and diversity of the data. It is often difficult to 1124  
 1063 train such models in an unbiased way, since it is not uncom- 1125  
 1064 mon for data-sets to contain missing or inaccurate information, 1126  
 1065 and unbalanced data sets with some underrepresented features 1127  
 1066 might bias the model. Additionally, a large amount of data 1128  
 1067 is needed to train a transferable model that can be applied to 1129  
 1068 different system types whilst still achieving a high accuracy of 1130  
 1069 predictions compared to more system-specific models, in which 1131  
 1070 a smaller subset of components reduces the types of correla- 1132  
 1071 tions that the model has to learn. Splitting a large data set into 1133  
 1072 different subsets of similar materials could improve accuracy at 1134  
 1073 the cost of generalisability. For soft materials in particular, one 1135  
 1074 challenge in developing appropriate models is that fewer data 1136  
 1075 are available compared to inorganic materials in general, and 1137  
 1076 there is a wider variety of structures and behaviours that soft ma- 1138  
 1077 terials can exhibit, so approaches will need to be more tailored 1139  
 1078 to specific contexts and system types. The effects of entropy 1140  
 1079 may also be more significant for soft-matter systems, which 1141  
 1080 may prove more challenging to capture with simple machine- 1142  
 1081 learning approaches. Moreover, a data-driven approach only 1143  
 1082 sees its input data – for example structures and energies from 1144  
 1083 DFT calculations or from experiment –, but does not know 1145  
 1084 anything about the underlying interactions between particles. 1146  
 1085 Although predictions of behaviours under given conditions can 1147  
 1086 often be surprisingly good, whether they are good for the right 1148  
 1087 reasons is less clear: they are entirely phenomenological and 1149  
 1088 even a significantly improved understanding or description of 1150  
 1089 the physics of the building blocks would not directly help with 1151  
 1090 the accuracy of the predictions of such models.

1091 By contrast, the traditional methods of statistical and classical 1153  
 1092 thermodynamics, such as thermodynamic integration, Gibbs– 1154  
 1093 Duhem integration, histogram reweighting or Widom insertion, 1155  
 1094 tend to be much slower because they require a number of 1156  
 1095 well-equilibrated simulations to be run. The results obtained 1157  
 1096 with such methods can be drastically improved with a better 1158  
 1097 description of the building blocks, indeed perhaps even with 1159  
 1098 a machine-learned interatomic potential, but the methods are 1160  
 1099 often both computationally and labour intensive, and are not 1161  
 1100 easy to automate for larger-scale applications. A combination 1162  
 1101 of data-driven methods for an initial screening followed by 1163  
 1102 traditional statistical-mechanical approaches to obtain accurate 1164  
 1103 results may be the key to future applications of computational 1165  
 1104 thermodynamics.

1105 Although we have focussed primarily on predictions of phase 1106  
 1107 diagrams in this article, several of the approaches we have 1108  
 1109 outlined can provide details that go beyond mere thermodynamic 1110  
 1111 stability. In particular, by determining the chemical potentials 1112  
 1113 of potentially competitive phases, we can determine regions of 1114  
 1115 metastability in addition to pure thermodynamic stability. This 1116  
 1117 can be especially useful in understanding the likely experimental 1118  
 1119 routes to synthesising metastable phases. For example, it may 1120  
 1121 be more productive to try to synthesise a given metastable 1122  
 1123 phase under conditions where the thermodynamically stable 1124  
 1124 phase is only very marginally more stable.<sup>8</sup> Metastable phases 1125  
 1126 can sometimes be relatively more stable in the initial stages 1127  
 1127 of a phase transformation, even if the final bulk structure is a 1128  
 1128 different one. An investigation of such metastable phases can 1129  
 1129 thus help to reveal the possible mechanisms by which phases can 1130  
 1130 transform into one another. For example, the metastable ‘ice 0’ 1131  
 1131 structure was first proposed in the context of the mechanism of 1132  
 1132 ice I nucleation.<sup>224</sup> We expect that, as calculations of chemical 1133  
 1133 potentials become more routine, the importance of metastable 1134  
 1134 phases will become clearer, and many more such pathways are 1135  
 1135 likely to be identified.

1136 Of course understanding the thermodynamics, stability and 1137  
 1137 metastability is only the first step in understanding a phase 1138  
 1138 transition. Studying the way phase transitions occur under 1139  
 1139 different conditions is perhaps an even more challenging prob- 1140  
 1140 lem.<sup>225</sup> Broadly speaking, relatively close to coexistence, there 1141  
 1141 are two main mechanisms by which phase transitions occur: 1142  
 1142 heterogeneous nucleation, where an external nucleation seed is 1143  
 1143 available, and homogeneous nucleation, where a spontaneous 1144  
 1144 fluctuation has to overcome a free-energy barrier. Even though 1145  
 1145 the phase diagram of substances like water close to atmospheric 1146  
 1146 conditions is well known, a huge amount of work has gone into 1147  
 1147 understanding the dynamics of ice I nucleation,<sup>153,224,226–252</sup> 1148  
 1148 and there is doubtless more work to be done. Indeed, Oxtoby 1149  
 1149 remarked back in 1998 that the study of nucleation is ‘one 1150  
 1150 of the few areas of science in which agreement of predicted 1151  
 1151 and measured rates to within several orders of magnitude is 1152  
 1152 considered a major success’,<sup>253</sup> and more than two decades 1153  
 1153 later, the same could be argued to hold. Nucleation is a rare 1154  
 1154 event and this makes it difficult to study in computer simula- 1155  
 1155 tions. Not only does it happen infrequently, but in order to 1156  
 1156 track its progress, and potentially to drive it with rare-event 1157  
 1157 methods, as with the interface pinning method, a suitable order 1158  
 1158 parameter must first be identified to distinguish between the 1159  
 1159 two phases in question, as discussed in Sec. II E. As systems 1160  
 1160 become more complicated, distinguishing between phases be- 1161  
 1161 comes more difficult, and machine-learning approaches might 1162  
 1162 again be helpful.<sup>254–258</sup> For example, Statt and co-workers 1163  
 1163 have used unsupervised machine-learning techniques for local 1164  
 1164 environments<sup>259</sup> to identify and classify different aggregate 1165  
 1165 morphologies formed from model copolymer sequences<sup>260</sup> 1166  
 1166 which the more conventional local order parameters failed to 1167  
 1167 distinguish.<sup>142</sup> Techniques such as principal component ana- 1168  
 1168 lysis have also been used to find order parameters for detecting 1169  
 1169 the freezing transition of hard spheres and ellipses, liquid– 1170  
 1170 vapour phase separation of patchy particles and compositional 1171  
 1171 demixing in the Widom–Rowlinson model.<sup>261,262</sup> Such order 1172  
 1172 parameters will likely prove to be a useful tool in investigating 1173  
 1173 nucleation mechanisms in more complex systems, including 1174  
 1174 possible transitions to metastable phases discussed above.

1175 In this article, we have largely considered simple systems 1176  
 1176 made up of only one or two components. Even for such systems,

1167 it is challenging to determine their phase behaviour well. The  
 1168 primary reason for this is still the quality of interatomic and  
 1169 intermolecular potentials, which are generally not applicable  
 1170 far beyond the range over which they were parameterised. We  
 1171 anticipate that ever-better interaction potentials will become  
 1172 available in the near future as machine-learned potentials are  
 1173 better parameterised.<sup>263</sup> However, there are still limitations in  
 1174 what can be simulated. For example, computer simulations  
 1175 often comprise relatively small systems, and it can be difficult  
 1176 to study defect formation. Schottky-style defects can be ac-  
 1177 counted for,<sup>80</sup> as can interstitial defects of neutral atoms;<sup>264,265</sup>  
 1178 however, interstitial defects involving charged species pose  
 1179 more of a challenge. Experimentally, it is possible to determine  
 1180 the energy of interstitial formation, but not usually the free  
 1181 energy: the effect of such defects on the phase behaviour is  
 1182 consequently very difficult to ascertain using anything other  
 1183 than computer simulations. Refining the methodology to facil-  
 1184 itate the study of such defects would enable us to determine  
 1185 accurate high-pressure phase behaviour of technologically and  
 1186 geologically important materials even at high temperatures at  
 1187 which entropically driven defects are likely to play a significant  
 1188 role.

1189 Much recent work has been done on extending the ‘standard’  
 1190 free-energy methods to study both defects and interfaces.<sup>266,267</sup>  
 1191 Yeandel and co-workers have, for example, used the Einstein  
 1192 crystal framework to compute interfacial free-energy densities  
 1193 of solid–liquid mixtures, including systems where the solid con-  
 1194 tains miscible species that diffuse in the liquid.<sup>267</sup> In particular,  
 1195 they used the fact that Einstein crystal particles do not interact  
 1196 with one another and therefore the precise location of each indi-  
 1197 vidual molecule in either the bulk crystal phase or a crystalline  
 1198 slab in contact with the liquid is immaterial. This simplifies  
 1199 the calculation of interfacial properties considerably, and is  
 1200 another illustration of the power of thermodynamic integration  
 1201 and the construction of clever thermodynamic cycles, which  
 1202 should enable a progressively simpler calculation of interfacial  
 1203 properties at coexistence, and in particular the thermodynamics  
 1204 of systems in contact with non-pure solutions.<sup>267</sup>

1205 Finally, there has been significant recent progress in simulat-  
 1206 ing multicomponent systems.<sup>268</sup> Compositional phase diagrams  
 1207 involving fluid phases could be obtained using chemical po-  
 1208 tential measurements,<sup>269,270</sup> and such ideas may be important  
 1209 when investigating liquid–solid solutions. For example, deep eu-  
 1210 tectic solvents are thought to be promising for green-chemistry  
 1211 separation methods.<sup>271–273</sup> Computational modelling of such  
 1212 systems has often focussed on small systems that could be  
 1213 treated with density-functional theory, empirical equations of  
 1214 state or statistical associated fluid theory;<sup>274,275</sup> as discussed  
 1215 above, a determination of accurate phase diagrams for such  
 1216 systems could help in the development of transferable empirical  
 1217 potentials, allowing new insights to be gained into these techno-  
 1218 logically important systems from a molecular perspective. This  
 1219 challenging problem is only just beginning to be addressed.<sup>276</sup>

## 1220 ACKNOWLEDGMENTS

1221 This work was supported by the University of Cambridge  
 1222 Ernest Oppenheimer Fund and the Winton Programme for the  
 1223 Physics of Sustainability.

## 1224 REFERENCES

- 1225 <sup>1</sup>C. J. King, *Separation processes*, 2nd ed. (McGraw-Hill, 1980).
- 1226 <sup>2</sup>C. G. Salzmann, P. G. Radaelli, B. Slater, and J. L. Finney, ‘The poly-  
 1227 morphism of ice: Five unresolved questions,’ *Phys. Chem. Chem. Phys.* **13**,  
 1228 18468–18480 (2011).
- 1229 <sup>3</sup>C. G. Salzmann, ‘Advances in the experimental exploration of water’s phase  
 1230 diagram,’ *J. Chem. Phys.* **150**, 060901 (2019).
- 1231 <sup>4</sup>M. A. Neumann and J. van de Streek, ‘How many ritonavir cases are there  
 1232 still out there?’ *Faraday Discuss.* **211**, 441–458 (2018).
- 1233 <sup>5</sup>W. G. Hoover and F. H. Ree, ‘Melting transition and communal entropy for  
 1234 hard spheres,’ *J. Chem. Phys.* **49**, 3609–3617 (1968).
- 1235 <sup>6</sup>B. Cheng, G. Mazzola, C. J. Pickard, and M. Ceriotti, ‘Evidence for  
 1236 supercritical behaviour of high-pressure liquid hydrogen,’ *Nature* **585**,  
 1237 217–220 (2020).
- 1238 <sup>7</sup>B. Cheng, M. Bethkenhagen, C. J. Pickard, and S. Hamel, ‘Phase behaviours  
 1239 of superionic water at planetary conditions,’ *Nat. Phys.* **17**, 1228–1232  
 1240 (2021).
- 1241 <sup>8</sup>A. Reinhardt and B. Cheng, ‘Quantum-mechanical exploration of the phase  
 1242 diagram of water,’ *Nat. Commun.* **12**, 588 (2021).
- 1243 <sup>9</sup>A. Reinhardt, M. Bethkenhagen, F. Coppari, M. Millot, S. Hamel, and  
 1244 B. Cheng, ‘Thermodynamics of high-pressure ice phases explored with  
 1245 atomistic simulations,’ *Nat. Commun.* **13**, 4707 (2022).
- 1246 <sup>10</sup>E. Florez, O. R. Smits, J.-M. Mewes, P. Jerabek, and P. Schwerdtfeger, ‘From  
 1247 the gas phase to the solid state: The chemical bonding in the superheavy  
 1248 element flerovium,’ *J. Chem. Phys.* **157**, 064304 (2022).
- 1249 <sup>11</sup>D. Weidler and J. Gross, ‘Transferable anisotropic united-atom force field  
 1250 based on the Mie potential for phase equilibria: Aldehydes, ketones, and  
 1251 small cyclic alkanes,’ *Ind. Eng. Chem. Res.* **55**, 12123–12132 (2016).
- 1252 <sup>12</sup>A. Reinhardt, ‘Phase behavior of empirical potentials of titanium dioxide,’  
 1253 *J. Chem. Phys.* **151**, 064505 (2019).
- 1254 <sup>13</sup>X. Wang, S. Ramírez-Hinestrosa, J. Dobnikar, and D. Frenkel, ‘The Lennard-  
 1255 Jones potential: when (not) to use it,’ *Phys. Chem. Chem. Phys.* **22**, 10624–  
 1256 10633 (2020).
- 1257 <sup>14</sup>D. Atherton, A. Michaelides, and S. J. Cox, ‘Can molecular simulations  
 1258 reliably compare homogeneous and heterogeneous ice nucleation?’ *J. Chem.*  
 1259 *Phys.* **156**, 164501 (2022).
- 1260 <sup>15</sup>J. L. F. Abascal and C. Vega, ‘A general purpose model for the condensed  
 1261 phases of water: TIP4P/2005,’ *J. Chem. Phys.* **123**, 234505 (2005).
- 1262 <sup>16</sup>C. Vega, E. Sanz, and J. L. F. Abascal, ‘The melting temperature of the most  
 1263 common models of water,’ *J. Chem. Phys.* **122**, 114507 (2005).
- 1264 <sup>17</sup>C. Vega, J. L. F. Abascal, E. Sanz, L. G. MacDowell, and C. McBride, ‘Can  
 1265 simple models describe the phase diagram of water?’ *J. Phys.: Condens.*  
 1266 *Matter* **17**, S3283 (2005).
- 1267 <sup>18</sup>C. Vega, E. Sanz, J. L. F. Abascal, and E. G. Noya, ‘Determination of phase  
 1268 diagrams via computer simulation: Methodology and applications to water,  
 1269 electrolytes and proteins,’ *J. Phys.: Condens. Matter* **20**, 153101 (2008).
- 1270 <sup>19</sup>C. Vega, J. L. F. Abascal, M. M. Conde, and J. L. Aragonés, ‘What ice can  
 1271 teach us about water interactions: A critical comparison of the performance  
 1272 of different water models,’ *Faraday Discuss.* **141**, 251–276 (2009).
- 1273 <sup>20</sup>C. Vega and J. L. F. Abascal, ‘Simulating water with rigid non-polarizable  
 1274 models: A general perspective,’ *Phys. Chem. Chem. Phys.* **13**, 19663–19688  
 1275 (2011).
- 1276 <sup>21</sup>P. T. Kiss and A. Baranyai, ‘Sources of the deficiencies in the popular SPC/E  
 1277 and TIP3P models of water,’ *J. Chem. Phys.* **134**, 054106 (2011).
- 1278 <sup>22</sup>E. G. Noya, C. Menduiña, J. L. Aragonés, and C. Vega, ‘Equation of state,  
 1279 thermal expansion coefficient, and isothermal compressibility for ices I<sub>h</sub>, II,  
 1280 III, V, and VI, as obtained from computer simulation,’ *J. Phys. Chem. C*  
 1281 **111**, 15877–15888 (2007).
- 1282 <sup>23</sup>H. L. Pi, J. L. Aragonés, C. Vega, E. G. Noya, J. L. F. Abascal, M. A.  
 1283 Gonzalez, and C. McBride, ‘Anomalies in water as obtained from computer  
 1284 simulations of the TIP4P/2005 model: Density maxima, and density, iso-  
 1285 thermal compressibility and heat capacity minima,’ *Mol. Phys.* **107**, 365–374  
 1286 (2009).
- 1287 <sup>24</sup>C. McBride, E. G. Noya, J. L. Aragonés, M. M. Conde, and C. Vega, ‘The  
 1288 phase diagram of water from quantum simulations,’ *Phys. Chem. Chem.*  
 1289 *Phys.* **14**, 10140–10146 (2012).
- 1290 <sup>25</sup>L.-P. Wang, T. Head-Gordon, J. W. Ponder, P. Ren, J. D. Chodera, P. K.  
 1291 Eastman, T. J. Martinez, and V. S. Pande, ‘Systematic improvement of a  
 1292 classical molecular model of water,’ *J. Phys. Chem. B* **117**, 9956–9972  
 1293 (2013).
- 1294 <sup>26</sup>J. A. Joseph, A. Reinhardt, A. Aguirre, P. Y. Chew, K. O. Russell, J. R.  
 1295 Espinosa, A. Garaizar, and R. Collepardo-Guevara, ‘Physics-driven coarse-  
 1296 grained model for biomolecular phase separation with near-quantitative

- accuracy,' *Nat. Comput. Sci.* **1**, 732–743 (2021).
- <sup>27</sup>T. Dannenhoffer-Lafage and R. B. Best, 'A data-driven hydrophobicity scale for predicting liquid–liquid phase separation of proteins,' *J. Phys. Chem. B* **125**, 4046–4056 (2021).
- <sup>28</sup>G. Tesei, T. K. Schulze, R. Crehuet, and K. Lindorff-Larsen, 'Accurate model of liquid–liquid phase behavior of intrinsically disordered proteins from optimization of single-chain properties,' *Proc. Natl Acad. Sci. U.S.A.* **118**, e2111696118 (2021).
- <sup>29</sup>R. Iftimie, P. Minary, and M. E. Tuckerman, 'Ab initio molecular dynamics: Concepts, recent developments, and future trends,' *Proc. Natl. Acad. Sci. U. S. A.* **102**, 6654–6659 (2005).
- <sup>30</sup>F. Ercolessi and J. B. Adams, 'Interatomic potentials from first-principles calculations: The force-matching method,' *EPL–Europhys. Lett.* **26**, 583–588 (1994).
- <sup>31</sup>W. L. Jorgensen, J. Chandrasekhar, J. D. Madura, R. W. Impey, and M. L. Klein, 'Comparison of simple potential functions for simulating liquid water,' *J. Chem. Phys.* **79**, 926–935 (1983).
- <sup>32</sup>D. Reith, M. Pütz, and F. Müller-Plathe, 'Deriving effective mesoscale potentials from atomistic simulations,' *J. Comput. Chem.* **24**, 1624–1636 (2003).
- <sup>33</sup>F. Romano, J. Russo, and H. Tanaka, 'Novel stable crystalline phase for the Stillinger-Weber potential,' *Phys. Rev. B* **90**, 014204 (2014).
- <sup>34</sup>L. Lu, J. F. Dama, and G. A. Voth, 'Fitting coarse-grained distribution functions through an iterative force-matching method,' *J. Chem. Phys.* **139**, 121906 (2013).
- <sup>35</sup>T. D. Potter, J. Tasche, and M. R. Wilson, 'Assessing the transferability of common top-down and bottom-up coarse-grained molecular models for molecular mixtures,' *Phys. Chem. Chem. Phys.* **21**, 1912–1927 (2019).
- <sup>36</sup>E. G. Noya, C. Vega, J. P. K. Doye, and A. A. Louis, 'Phase diagram of model anisotropic particles with octahedral symmetry,' *J. Chem. Phys.* **127**, 054501 (2007).
- <sup>37</sup>F. Romano, E. Sanz, and F. Sciortino, 'Role of the range in the fluid–crystal coexistence for a patchy particle model,' *J. Phys. Chem. B* **113**, 15133 (2009).
- <sup>38</sup>F. Romano, E. Sanz, and F. Sciortino, 'Phase diagram of a tetrahedral patchy particle model for different interaction ranges,' *J. Chem. Phys.* **132**, 184501 (2010).
- <sup>39</sup>F. Romano, E. Sanz, and F. Sciortino, 'Crystallization of tetrahedral patchy particles *in silico*,' *J. Chem. Phys.* **134**, 174502 (2011).
- <sup>40</sup>F. Romano and F. Sciortino, 'Two dimensional assembly of triblock Janus particles into crystal phases in the two bond per patch limit,' *Soft Matter* **7**, 5799–5804 (2011).
- <sup>41</sup>A. Reinhardt, A. J. Williamson, J. P. K. Doye, J. Carrete, L. M. Varela, and A. A. Louis, 'Re-entrant phase behavior for systems with competition between phase separation and self-assembly,' *J. Chem. Phys.* **134**, 104905 (2011).
- <sup>42</sup>G. Doppelbauer, E. G. Noya, E. Bianchi, and G. Kahl, 'Self-assembly scenarios of patchy colloidal particles,' *Soft Matter* **8**, 7768–7772 (2012).
- <sup>43</sup>E. G. Noya, I. Kolovos, G. Doppelbauer, G. Kahl, and E. Bianchi, 'Phase diagram of inverse patchy colloids assembling into an equilibrium lamellar phase,' *Soft Matter* **10**, 8464–8474 (2014).
- <sup>44</sup>P. Teixeira and J. Tavares, 'Phase behaviour of pure and mixed patchy colloids — Theory and simulation,' *Curr. Op. Colloid Interface Sci.* **30**, 16–24 (2017).
- <sup>45</sup>B. Cheng, E. A. Engel, J. Behler, C. Dellago, and M. Ceriotti, 'Ab initio thermodynamics of liquid and solid water,' *Proc. Natl Acad. Sci. U. S. A.* **116**, 1110–1115 (2019).
- <sup>46</sup>L. Zhang, H. Wang, R. Car, and W. E., 'Phase diagram of a deep potential water model,' *Phys. Rev. Lett.* **126**, 236001 (2021).
- <sup>47</sup>H. Niu, L. Bonati, P. M. Piaggi, and M. Parrinello, 'Ab initio phase diagram and nucleation of gallium,' *Nat. Commun.* **11**, 2654 (2020).
- <sup>48</sup>J. G. Lee, C. J. Pickard, and B. Cheng, 'High-pressure phase behaviors of titanium dioxide revealed by a  $\delta$ -learning potential,' *J. Chem. Phys.* **156**, 074106 (2022).
- <sup>49</sup>A. Reinhardt, C. J. Pickard, and B. Cheng, 'Predicting the phase diagram of titanium dioxide with random search and pattern recognition,' *Phys. Chem. Chem. Phys.* **22**, 12697–12705 (2020).
- <sup>50</sup>V. Swamy, J. D. Gale, and L. S. Dubrovinsky, 'Atomistic simulation of the crystal structures and bulk moduli of TiO<sub>2</sub> polymorphs,' *J. Phys. Chem. Solids* **62**, 887–895 (2001).
- <sup>51</sup>J. Muscat, V. Swamy, and N. M. Harrison, 'First-principles calculations of the phase stability of TiO<sub>2</sub>,' *Phys. Rev. B* **65**, 224112 (2002).
- <sup>52</sup>X. G. Ma, P. Liang, L. Miao, S. W. Bie, C. K. Zhang, L. Xu, and J. J. Jiang, 'Pressure-induced phase transition and elastic properties of TiO<sub>2</sub> polymorphs,' *Phys. Status Solidi B* **246**, 2132–2139 (2009).
- <sup>53</sup>H. Dekura, T. Tsuchiya, Y. Kuwayama, and J. Tsuchiya, 'Theoretical and experimental evidence for a new post-cotunnite phase of titanium dioxide with significant optical absorption,' *Phys. Rev. Lett.* **107**, 045701 (2011).
- <sup>54</sup>V. Swamy and N. C. Wilson, 'First-principles calculations of the pressure stability and elasticity of dense TiO<sub>2</sub> phases using the B3LYP hybrid functional,' *J. Phys. Chem. C* **118**, 8617–8625 (2014).
- <sup>55</sup>Q.-J. Liu, Z. Ran, F.-S. Liu, and Z.-T. Liu, 'Phase transitions and mechanical stability of TiO<sub>2</sub> polymorphs under high pressure,' *J. Alloys Compd.* **631**, 192–201 (2015).
- <sup>56</sup>Y. Luo, A. Benali, L. Shulenburger, J. T. Krogel, O. Heinonen, and P. R. C. Kent, 'Phase stability of TiO<sub>2</sub> polymorphs from diffusion Quantum Monte Carlo,' *New J. Phys.* **18**, 113049 (2016).
- <sup>57</sup>J. Trail, B. Monserrat, P. López Ríos, R. Maezono, and R. J. Needs, 'Quantum Monte Carlo study of the energetics of the rutile, anatase, brookite, and columbite TiO<sub>2</sub> polymorphs,' *Phys. Rev. B* **95**, 121108 (2017).
- <sup>58</sup>Z. Fu, Y. Liang, S. Wang, and Z. Zhong, 'Structural phase transition and mechanical properties of TiO<sub>2</sub> under high pressure,' *Phys. Status Solidi B* **250**, 2206–2214 (2013).
- <sup>59</sup>M. A. Blanco, E. Francisco, and V. Luaña, 'GIBBS: Isothermal-isobaric thermodynamics of solids from energy curves using a quasi-harmonic Debye model,' *Comput. Phys. Commun.* **158**, 57–72 (2004).
- <sup>60</sup>Y. Jing-Xin, F. Min, J. Guang-Fu, and C. Xiang-Rong, 'Phase transition and thermodynamic properties of TiO<sub>2</sub> from first-principles calculations,' *Chinese Phys. B* **18**, 269 (2009).
- <sup>61</sup>Y.-F. Hu, G. Jiang, D.-Q. Meng, and F.-J. Kong, 'Phase transition and thermodynamic properties of TiO<sub>2</sub>,' *Acta Phys.-Chim. Sin.* **26**, 1664–1668 (2010).
- <sup>62</sup>Z.-G. Mei, Y. Wang, S. Shang, and Z.-K. Liu, 'First-principles study of the mechanical properties and phase stability of TiO<sub>2</sub>,' *Comput. Mater. Sci.* **83**, 114–119 (2014).
- <sup>63</sup>S. Srinivasan, R. Batra, D. Luo, T. Loeffler, S. Manna, H. Chan, L. Yang, W. Yang, J. Wen, P. Darancet, and S. K. Sankaranarayanan, 'Machine learning the metastable phase diagram of covalently bonded carbon,' *Nat. Commun.* **13**, 3251 (2022).
- <sup>64</sup>D. J. Wales and H. A. Scheraga, 'Global optimization of clusters, crystals, and biomolecules,' *Science* **285**, 1368–1372 (1999).
- <sup>65</sup>S. M. Woodley and R. Catlow, 'Crystal structure prediction from first principles,' *Nat. Mater.* **7**, 937–946 (2008).
- <sup>66</sup>L. Filion, M. Marechal, B. van Oorschot, D. Pelt, F. Smalenburg, and M. Dijkstra, 'Efficient method for predicting crystal structures at finite temperature: Variable box shape simulations,' *Phys. Rev. Lett.* **103**, 188302 (2009).
- <sup>67</sup>Y. Wang, J. Lv, L. Zhu, and Y. Ma, 'Crystal structure prediction via particle-swarm optimization,' *Phys. Rev. B* **82**, 094116 (2010).
- <sup>68</sup>C. J. Pickard and R. J. Needs, 'Ab initio random structure searching,' *J. Phys.: Condens. Matter* **23**, 053201 (2011).
- <sup>69</sup>H. T. Stokes and D. M. Hatch, 'FINDSYM: program for identifying the space-group symmetry of a crystal,' *J. Appl. Cryst.* **38**, 237–238 (2005).
- <sup>70</sup>V. Stevanović, 'Sampling polymorphs of ionic solids using random superlattices,' *Phys. Rev. Lett.* **116**, 075503 (2016).
- <sup>71</sup>B. Widom, 'Some topics in the theory of fluids,' *J. Chem. Phys.* **39**, 2808–2812 (1963).
- <sup>72</sup>A. Z. Panagiotopoulos, 'Direct determination of phase coexistence properties of fluids by Monte Carlo simulation in a new ensemble,' *Mol. Phys.* **61**, 813–826 (1987).
- <sup>73</sup>A. Z. Panagiotopoulos, N. Quirke, M. Stapleton, and D. J. Tildesley, 'Phase equilibria by simulation in the Gibbs ensemble – Alternative derivation, generalization and application to mixture and membrane equilibria,' *Mol. Phys.* **63**, 527–545 (1988).
- <sup>74</sup>N. B. Wilding and A. D. Bruce, 'Freezing by monte carlo phase switch,' *Phys. Rev. Lett.* **85**, 5138–5141 (2000).
- <sup>75</sup>N. B. Wilding, 'Phase switch Monte Carlo,' *AIP Conf. Proc.* **690**, 349–355 (2003).
- <sup>76</sup>F. Romano, E. Sanz, P. Tartaglia, and F. Sciortino, 'Phase diagram of trivalent and pentavalent patchy particles,' *J. Phys.: Condens. Matter* **24**, 064113 (2012).
- <sup>77</sup>A. Reinhardt, F. Romano, and J. P. K. Doye, 'Computing phase diagrams for a quasicrystal-forming patchy-particle system,' *Phys. Rev. Lett.* **110**, 255503 (2013).
- <sup>78</sup>G. Grochola, 'Constrained fluid  $\lambda$ -integration: Constructing a reversible thermodynamic path between the solid and liquid state,' *J. Chem. Phys.* **120**, 2122–2126 (2004).
- <sup>79</sup>D. Frenkel and A. J. C. Ladd, 'New Monte Carlo method to compute the free energy of arbitrary solids. Application to the fcc and hcp phases of hard spheres,' *J. Chem. Phys.* **81**, 3188–3193 (1984).

- 1447 <sup>80</sup>B. Cheng and M. Ceriotti, ‘Computing the absolute Gibbs free energy in 1522  
1448 atomistic simulations: Applications to defects in solids,’ *Phys. Rev. B* **97**,  
1449 [054102](#) (2018). 1523
- 1450 <sup>81</sup>J. L. Aragones, C. Valeriani, and C. Vega, ‘Note: Free energy calculations 1524  
1451 for atomic solids through the Einstein crystal/molecule methodology using 1525  
1452 GROMACS and LAMMPS,’ *J. Chem. Phys.* **137**, [146101](#) (2012). 1526
- 1453 <sup>82</sup>J. L. Aragones, E. G. Noya, C. Valeriani, and C. Vega, ‘Free energy 1527  
1454 calculations for molecular solids using GROMACS,’ *J. Chem. Phys.* **139**,  
1455 [034104](#) (2013). 1528
- 1456 <sup>83</sup>J. G. Kirkwood, ‘Statistical mechanics of fluid mixtures,’ *J. Chem. Phys.* **3**,  
1457 [300–313](#) (1935). 1529
- 1458 <sup>84</sup>D. Frenkel and B. M. Mulder, ‘The hard ellipsoid-of-revolution fluid I. 1530  
1459 Monte Carlo simulations,’ *Mol. Phys.* **55**, [1171–1192](#) (1985). 1531
- 1460 <sup>85</sup>C. Vega and E. G. Noya, ‘Revisiting the Frenkel-Ladd method to compute 1532  
1461 the free energy of solids: The Einstein molecule approach,’ *J. Chem. Phys.*  
1462 **127**, [154113](#) (2007). 1533
- 1463 <sup>86</sup>J. F. Cornwell, *Group theory in physics, Volume 1*, Techniques of physics 7  
1464 (Academic Press, London, 1986). 1534
- 1465 <sup>87</sup>M. Hestenes and E. Stiefel, ‘Methods of conjugate gradients for solving 1535  
1466 linear systems,’ *J. Res. Natl. Bur. Stand. (U.S.)* **49**, [409](#) (1952). 1536
- 1467 <sup>88</sup>J. Guérolé, W. G. Nöhling, A. Vaid, F. Houllé, Z. Xie, A. Prakash, and  
1468 E. Bitzek, ‘Assessment and optimization of the fast inertial relaxation  
1469 engine (FIRE) for energy minimization in atomistic simulations and its  
1470 implementation in LAMMPS,’ *Comput. Mater. Sci.* **175**, [109584](#) (2020). 1537
- 1471 <sup>89</sup>G. T. Gao, X. C. Zeng, and H. Tanaka, ‘The melting temperature of proton-  
1472 disordered hexagonal ice: A computer simulation of 4-site transferable  
1473 intermolecular potential model of water,’ *J. Chem. Phys.* **112**, [8534–8538](#)  
1474 (2000). 1538
- 1475 <sup>90</sup>S. Habershon and D. E. Manolopoulos, ‘Free energy calculations for a  
1476 flexible water model,’ *Phys. Chem. Chem. Phys.* **13**, [19714–19727](#) (2011). 1539
- 1477 <sup>91</sup>R. Ramírez and C. P. Herrero, ‘Quantum path integral simulation of isotope  
1478 effects in the melting temperature of ice Ih,’ *J. Chem. Phys.* **133**, [144511](#)  
1479 (2010). 1540
- 1480 <sup>92</sup>M. Ceriotti and T. E. Markland, ‘Efficient methods and practical guidelines  
1481 for simulating isotope effects,’ *J. Chem. Phys.* **138**, [014112](#) (2013). 1541
- 1482 <sup>93</sup>M. Ceriotti, W. Fang, P. G. Kusalik, R. H. McKenzie, A. Michaelides,  
1483 M. A. Morales, and T. E. Markland, ‘Nuclear quantum effects in water and  
1484 aqueous systems: Experiment, theory, and current challenges,’ *Chem. Rev.*  
1485 **116**, [7529–7550](#) (2016). 1542
- 1486 <sup>94</sup>B. Cheng, A. T. Paxton, and M. Ceriotti, ‘Hydrogen diffusion and trapping  
1487 in  $\alpha$ -iron: The role of quantum and anharmonic fluctuations,’ *Phys. Rev.*  
1488 *Lett.* **120**, [225901](#) (2018). 1543
- 1489 <sup>95</sup>J. Lee, ‘New Monte Carlo algorithm: Entropic sampling,’ *Phys. Rev. Lett.*  
1490 **71**, [211–214](#) (1993). 1544
- 1491 <sup>96</sup>F. Wang and D. P. Landau, ‘Efficient, multiple-range random walk algorithm  
1492 to calculate the density of states,’ *Phys. Rev. Lett.* **86**, [2050–2053](#) (2001). 1545
- 1493 <sup>97</sup>F. Wang and D. P. Landau, ‘Determining the density of states for classical  
1494 statistical models: A random walk algorithm to produce a flat histogram,’  
1495 *Phys. Rev. E* **64**, [056101](#) (2001). 1546
- 1496 <sup>98</sup>G. Ganzenmüller and P. J. Camp, ‘Applications of Wang–Landau sampling  
1497 to determining phase equilibria in complex fluids,’ *J. Chem. Phys.* **127**,  
1498 [154504](#) (2007). 1547
- 1499 <sup>99</sup>A. J. Williamson, *Methods, rules and limits of successful self-assembly*,  
1500 Ph.D. thesis, University of Oxford, Oxford (2011). 1548
- 1501 <sup>100</sup>V. Egorov and B. Kryzhanovskiy, ‘Influence of initial guess on the conver-  
1502 gence rate and the accuracy of Wang–Landau algorithm,’ *Opt. Mem. Neural*  
1503 *Netw.* **30**, [284–290](#) (2021). 1549
- 1504 <sup>101</sup>P. Poulain, F. Calvo, R. Antoine, M. Broyer, and P. Dugourd, ‘Performances  
1505 of Wang–Landau algorithms for continuous systems,’ *Phys. Rev. E* **73**,  
1506 [056704](#) (2006). 1550
- 1507 <sup>102</sup>S. Sinha and S. K. Roy, ‘Performance of wang–landau algorithm in contin-  
1508 uous spin models and a case study: Modified XY-model,’ *Phys. Lett. A* **373**,  
1509 [308–314](#) (2009). 1551
- 1510 <sup>103</sup>B. K. Latha and V. S. S. Sastry, ‘Phase diagram of a general biaxial nematic  
1511 model based on density of states computation,’ *Liq. Cryst.* **45**, [2197–2213](#)  
1512 (2018). 1552
- 1513 <sup>104</sup>M. Borg, C. Stampfl, A. Mikkelsen, J. Gustafson, E. Lundgren, M. Scheffler,  
1514 and J. N. Andersen, ‘Density of configurational states from first-principles  
1515 calculations: The phase diagram of Al–Na surface alloys,’ *ChemPhysChem*  
1516 **6**, [1923–1928](#) (2005). 1553
- 1517 <sup>105</sup>M. S. Shell, P. G. Debenedetti, and A. Z. Panagiotopoulos, ‘Generalization  
1518 of the Wang–Landau method for off-lattice simulations,’ *Phys. Rev. E* **66**,  
1519 [056703](#) (2002). 1554
- 1520 <sup>106</sup>C. Desgranges and J. Delhommelle, ‘Phase equilibria of molecular fluids  
1521 via hybrid Monte Carlo Wang–Landau simulations: Applications to benzene  
and *n*-alkanes,’ *J. Chem. Phys.* **130**, [244109](#) (2009). 1555
- <sup>107</sup>S. Boothroyd, A. Kerridge, A. Broo, D. Buttar, and J. Anwar, ‘Solubility  
1526 prediction from first principles: a density of states approach,’ *Phys. Chem.*  
1527 *Chem. Phys.* **20**, [20981–20987](#) (2018). 1556
- <sup>108</sup>H. Do, J. D. Hirst, and R. J. Wheatley, ‘Rapid calculation of partition  
1528 functions and free energies of fluids,’ *J. Chem. Phys.* **135**, [174105](#) (2011). 1557
- <sup>109</sup>H. Do and R. J. Wheatley, ‘Density of states partitioning method for  
1529 calculating the free energy of solids,’ *J. Chem. Theory Comput.* **9**, [165–171](#)  
1530 (2013). 1558
- <sup>110</sup>H. Do and R. J. Wheatley, ‘Reverse energy partitioning—an efficient  
1531 algorithm for computing the density of states, partition functions, and free  
1532 energy of solids,’ *J. Chem. Phys.* **145**, [084116](#) (2016). 1559
- <sup>111</sup>G. Torrie and J. Valleau, ‘Nonphysical sampling distributions in Monte  
1533 Carlo free-energy estimation: Umbrella sampling,’ *J. Comput. Phys.* **23**,  
1534 [187–199](#) (1977). 1560
- <sup>112</sup>J. Skilling, ‘Nested sampling,’ *AIP Conf. Proc.* **735**, [395–405](#) (2004). 1561
- <sup>113</sup>L. B. Pártay, A. P. Bartók, and G. Csányi, ‘Efficient sampling of atomic  
1535 configurational spaces,’ *J. Phys. Chem. B* **114**, [10502–10512](#) (2010). 1562
- <sup>114</sup>L. B. Pártay, G. Csányi, and N. Bernstein, ‘Nested sampling for materials,’  
1540 *Eur. Phys. J. B* **94**, [159](#) (2021). 1563
- <sup>115</sup>S. Martiniani, J. D. Stevenson, D. J. Wales, and D. Frenkel, ‘Superposition  
1541 enhanced nested sampling,’ *Phys. Rev. X* **4**, [031034](#) (2014). 1564
- <sup>116</sup>R. J. N. Baldock, L. B. Pártay, A. P. Bartók, M. C. Payne, and G. Csányi,  
1542 ‘Determining pressure-temperature phase diagrams of materials,’ *Phys.*  
1543 *Rev. B* **93**, [174108](#) (2016). 1565
- <sup>117</sup>A. C. L. Opitz, ‘Molecular dynamics investigation of a free surface of liquid  
1544 argon,’ *Phys. Lett. A* **47**, [439–440](#) (1974). 1566
- <sup>118</sup>A. J. C. Ladd and L. V. Woodcock, ‘Triple-point coexistence properties of  
1545 the Lennard-Jones system,’ *Chem. Phys. Lett.* **51**, [155–159](#) (1977). 1567
- <sup>119</sup>J. R. Morris and X. Song, ‘The melting lines of model systems calculated  
1546 from coexistence simulations,’ *J. Chem. Phys.* **116**, [9352–9358](#) (2002). 1568
- <sup>120</sup>R. L. Davidchack and B. B. Laird, ‘Simulation of the hard-sphere crystal–  
1547 melt interface,’ *J. Chem. Phys.* **108**, [9452–9462](#) (1998). 1569
- <sup>121</sup>E. G. Noya, C. Vega, and E. de Miguel, ‘Determination of the melting point  
1548 of hard spheres from direct coexistence simulation methods,’ *J. Chem. Phys.*  
1549 **128**, [154507](#) (2008). 1570
- <sup>122</sup>J. R. Espinosa, E. Sanz, C. Valeriani, and C. Vega, ‘On fluid-solid direct  
1550 coexistence simulations: The pseudo-hard sphere model,’ *J. Chem. Phys.*  
1551 **139**, [144502](#) (2013). 1571
- <sup>123</sup>B. J. Jesson and P. A. Madden, ‘Structure and dynamics at the aluminum  
1552 solid–liquid interface: An *ab initio* simulation,’ *J. Chem. Phys.* **113**, [5935–](#)  
1553 [5946](#) (2000). 1572
- <sup>124</sup>A. B. Belonoshko, R. Ahuja, and B. Johansson, ‘Quasi-*ab initio* molecular  
1554 dynamic study of Fe melting,’ *Phys. Rev. Lett.* **84**, [3638–3641](#) (2000). 1573
- <sup>125</sup>R. García Fernández, J. L. F. Abascal, and C. Vega, ‘The melting point of  
1555 ice I<sub>h</sub> for common water models calculated from direct coexistence of the  
1556 solid–liquid interface,’ *J. Chem. Phys.* **124**, [144506](#) (2006). 1574
- <sup>126</sup>M. M. Conde, M. A. Gonzalez, J. L. F. Abascal, and C. Vega, ‘Determining  
1557 the phase diagram of water from direct coexistence simulations: The phase  
1558 diagram of the TIP4P/2005 model revisited,’ *J. Chem. Phys.* **139**, [154505](#)  
1559 (2013). 1575
- <sup>127</sup>M. Fitzner, L. Joly, M. Ma, G. C. Sosso, A. Zen, and A. Michaelides,  
1560 ‘Communication: Truncated non-bonded potentials can yield unphysical  
1561 behavior in molecular dynamics simulations of interfaces,’ *J. Chem. Phys.*  
1562 **147**, [121102](#) (2017). 1576
- <sup>128</sup>U. Tartaglino and E. Tosatti, ‘Strain effects at solid surfaces near the melting  
1563 point,’ *Surf. Sci.* **532–535**, [623–627](#) (2003). 1577
- <sup>129</sup>J. Q. Broughton and G. H. Gilmer, ‘Molecular dynamics of the crystal–fluid  
1564 interface. V. Structure and dynamics of crystal–melt systems,’ *J. Chem.*  
1565 *Phys.* **84**, [5749–5758](#) (1986). 1578
- <sup>130</sup>J. Wang, S. Yoo, J. Bai, J. R. Morris, and X. C. Zeng, ‘Melting temperature  
1566 of ice I<sub>h</sub> calculated from coexisting solid-liquid phases,’ *J. Chem. Phys.* **123**,  
1567 [036101](#) (2005). 1579
- <sup>131</sup>A. S. Holehouse and R. V. Pappu, ‘Functional implications of intracellular  
1568 phase transitions,’ *Biochemistry* **57**, [2415–2423](#) (2018). 1580
- <sup>132</sup>S. Alberti and D. Dormann, ‘Liquid–liquid phase separation in disease,’  
1569 *Annu. Rev. Genet.* **53**, [171–194](#) (2019). 1581
- <sup>133</sup>S. M. Lichtinger, A. Garaizar, R. Collepardo-Guevara, and A. Reinhardt,  
1570 ‘Targeted modulation of protein liquid–liquid phase separation by evolution  
1571 of amino-acid sequence,’ *PLOS Comput. Biol.* **17**, [e1009328](#) (2021). 1582
- <sup>134</sup>I. Sanchez-Burgos, J. R. Espinosa, J. A. Joseph, and R. Collepardo-Guevara,  
1572 ‘Valency and binding affinity variations can regulate the multilayered organ-  
1573 ization of protein condensates with many components,’ *Biomolecules* **11**,  
1574 [278](#) (2021). 1583
- <sup>135</sup>P. Y. Chew, J. A. Joseph, R. Collepardo-Guevara, and A. Reinhardt, ‘Design-

- ing multiphase biomolecular condensates by coevolution of protein mixtures,' *bioRxiv* (2022), 10.1101/2022.04.22.489187.
- <sup>136</sup>A. Bremer, M. Farag, W. M. Borchers, I. Peran, E. W. Martin, R. V. Pappu, and T. Mittag, 'Deciphering how naturally occurring sequence features impact the phase behaviours of disordered prion-like domains,' *Nat. Chem.* **14**, 196–207 (2021).
- <sup>137</sup>G. Tesei and K. Lindorff-Larsen, 'Improved predictions of phase behaviour of intrinsically disordered proteins by tuning the interaction range,' *Open Research Europe* **2** (2022), 10.12688/openreseurope.14967.1.
- <sup>138</sup>G. Krainer, T. J. Welsh, J. A. Joseph, J. R. Espinosa, S. Wittmann, E. de Csilléry, A. Sridhar, Z. Toprakcioglu, G. Gudískytė, M. A. Czekalska, W. E. Arter, J. Guillén-Boixet, T. M. Franzmann, S. Qamar, P. S. George-Hyslop, A. A. Hyman, R. Collepardo-Guevara, S. Alberti, and T. P. Knowles, 'Reentrant liquid condensate phase of proteins is stabilized by hydrophobic and non-ionic interactions,' *Nat. Commun.* **12**, 1085 (2021).
- <sup>139</sup>J. R. Espinosa, J. A. Joseph, I. Sanchez-Burgos, A. Garaizar, D. Frenkel, and R. Collepardo-Guevara, 'Liquid network connectivity regulates the stability and composition of biomolecular condensates with many components,' *Proc. Natl. Acad. Sci. U. S. A.* **117**, 13238–13247 (2020).
- <sup>140</sup>E. W. Martin, A. S. Holehouse, I. Peran, M. Farag, J. J. Incicco, A. Bremer, C. R. Grace, A. Soranno, R. V. Pappu, and T. Mittag, 'Valence and patterning of aromatic residues determine the phase behavior of prion-like domains,' *Science* **367**, 694–699 (2020).
- <sup>141</sup>Y. Zhang, B. Xu, B. G. Weiner, Y. Meir, and N. S. Wingreen, 'Decoding the physical principles of two-component biomolecular phase separation,' *eLife* **10**, e62403 (2021).
- <sup>142</sup>A. Statt, H. Casademunt, C. P. Brangwynne, and A. Z. Panagiotopoulos, 'Model for disordered proteins with strongly sequence-dependent liquid phase behavior,' *J. Chem. Phys.* **152**, 75101 (2020).
- <sup>143</sup>Y.-H. Lin, J. P. Brady, J. D. Forman-Kay, and H. S. Chan, 'Charge pattern matching as a "fuzzy" mode of molecular recognition for the functional phase separations of intrinsically disordered proteins,' *New J. Phys.* **19**, 115003 (2017).
- <sup>144</sup>U. R. Pedersen, F. Hummel, G. Kresse, G. Kahl, and C. Dellago, 'Computing Gibbs free energy differences by interface pinning,' *Phys. Rev. B* **88**, 94101 (2013).
- <sup>145</sup>U. R. Pedersen, 'Direct calculation of the solid-liquid Gibbs free energy difference in a single equilibrium simulation,' *J. Chem. Phys.* **139**, 104102 (2013).
- <sup>146</sup>U. R. Pedersen, F. Hummel, and C. Dellago, 'Computing the crystal growth rate by the interface pinning method,' *J. Chem. Phys.* **142**, 44104 (2015).
- <sup>147</sup>P. J. Steinhardt, D. R. Nelson, and M. Ronchetti, 'Bond-orientational order in liquids and glasses,' *Phys. Rev. B* **28**, 784–805 (1983).
- <sup>148</sup>P.-R. ten Wolde, M. J. Ruiz-Montero, and D. Frenkel, 'Simulation of homogeneous crystal nucleation close to coexistence,' *Faraday Discuss.* **104**, 93–110 (1996).
- <sup>149</sup>B. Senger, P. Schaaf, D. S. Corti, R. Bowles, J.-C. Voegel, and H. Reiss, 'A molecular theory of the homogeneous nucleation rate. I. Formulation and fundamental numbers,' *J. Chem. Phys.* **110**, 6421–6437 (1999).
- <sup>150</sup>W. Lechner and C. Dellago, 'Accurate determination of crystal structures based on averaged local bond order parameters,' *J. Chem. Phys.* **129**, 114707 (2008).
- <sup>151</sup>S. Jungblut and C. Dellago, 'Crystallization of a binary Lennard-Jones mixture,' *J. Chem. Phys.* **134**, 104501 (2011).
- <sup>152</sup>E. B. Moore, E. de la Llave, K. Welke, D. A. Scherlis, and V. Molinero, 'Freezing, melting and structure of ice in a hydrophilic nanopore,' *Phys. Chem. Chem. Phys.* **12**, 4124 (2010).
- <sup>153</sup>A. Reinhardt, J. P. K. Doye, E. G. Noya, and C. Vega, 'Local order parameters for use in driving homogeneous ice nucleation with all-atom models of water,' *J. Chem. Phys.* **137**, 194504 (2012).
- <sup>154</sup>G. B. Arfken and H. J. Weber, *Mathematical methods for physicists*, 6th ed. (Elsevier Academic, London, 2005).
- <sup>155</sup>M. A. Blanco, M. Flórez, and M. Bermejo, 'Evaluation of the rotation matrices in the basis of real spherical harmonics,' *J. Mol. Struct.—Theochem.* **419**, 19–27 (1997).
- <sup>156</sup>P.-L. Chau and A. J. Hardwick, 'A new order parameter for tetrahedral configurations,' *Mol. Phys.* **93**, 511–518 (1998).
- <sup>157</sup>J. R. Errington and P. G. Debenedetti, 'Relationship between structural order and the anomalies of liquid water,' *Nature* **409**, 318–321 (2001).
- <sup>158</sup>E. E. Santiso and B. L. Trout, 'A general set of order parameters for molecular crystals,' *J. Chem. Phys.* **134**, 064109 (2011).
- <sup>159</sup>D. A. Kofke, 'Gibbs–Duhem integration: A new method for direct evaluation of phase coexistence by molecular simulation,' *Mol. Phys.* **78**, 1331–1336 (1993).
- <sup>160</sup>D. A. Kofke, 'Direct evaluation of phase coexistence by molecular simulation via integration along the saturation line,' *J. Chem. Phys.* **98**, 4149–4162 (1993).
- <sup>161</sup>M. R. Shirts and J. D. Chodera, 'Statistically optimal analysis of samples from multiple equilibrium states,' *J. Chem. Phys.* **129**, 124105 (2008).
- <sup>162</sup>N. P. Schieber, E. C. Dybeck, and M. R. Shirts, 'Using reweighting and free energy surface interpolation to predict solid-solid phase diagrams,' *J. Chem. Phys.* **148**, 144104 (2018).
- <sup>163</sup>G. Bauer and J. Gross, 'Phase equilibria of solid and fluid phases from molecular dynamics simulations with equilibrium and nonequilibrium free energy methods,' *J. Chem. Theory Comput.* **15**, 3778–3792 (2019).
- <sup>164</sup>L. Pauling, 'The structure and entropy of ice and of other crystals with some randomness of atomic arrangement,' *J. Am. Chem. Soc.* **57**, 2680–2684 (1935).
- <sup>165</sup>J. D. Bernal and R. H. Fowler, 'A theory of water and ionic solution, with particular reference to hydrogen and hydroxyl ions,' *J. Chem. Phys.* **1**, 515–548 (1933).
- <sup>166</sup>J. A. Hayward and J. R. Reimers, 'Unit cells for the simulation of hexagonal ice,' *J. Chem. Phys.* **106**, 1518–1529 (1997).
- <sup>167</sup>V. Buch, P. Sandler, and J. Sadlej, 'Simulations of H<sub>2</sub>O solid, liquid, and clusters, with an emphasis on ferroelectric ordering transition in hexagonal ice,' *J. Phys. Chem. B* **102**, 8641–8653 (1998).
- <sup>168</sup>C. P. Herrero and R. Ramírez, 'Configurational entropy of hydrogen-disordered ice polymorphs,' *J. Chem. Phys.* **140**, 234502 (2014).
- <sup>169</sup>L. G. MacDowell, E. Sanz, C. Vega, and J. L. F. Abascal, 'Combinatorial entropy and phase diagram of partially ordered ice phases,' *J. Chem. Phys.* **121**, 10145–10158 (2004).
- <sup>170</sup>M. Matsumoto, T. Yagasaki, and H. Tanaka, 'GenIce: Hydrogen-disordered ice generator,' *J. Comput. Chem.* **39**, 61–64 (2017).
- <sup>171</sup>J. L. Aragones, L. G. MacDowell, and C. Vega, 'Dielectric constant of ices and water: A lesson about water interactions,' *J. Phys. Chem. A* **115**, 5745–5758 (2011).
- <sup>172</sup>W. Steurer, 'Twenty years of structure research on quasicrystals. Part I. Pentagonal, octagonal, decagonal and dodecagonal quasicrystals,' *Z. Kristallogr.* **219**, 391–446 (2004).
- <sup>173</sup>M. Zu, P. Tan, and N. Xu, 'Forming quasicrystals by monodisperse soft core particles,' *Nat. Commun.* **8**, 2089 (2017).
- <sup>174</sup>C. R. Iacovella, A. S. Keys, and S. C. Glotzer, 'Self-assembly of soft-matter quasicrystals and their approximants,' *Proc. Natl. Acad. Sci. U. S. A.* **108**, 20935–20940 (2011).
- <sup>175</sup>P. F. Damasceno, S. C. Glotzer, and M. Engel, 'Non-close-packed three-dimensional quasicrystals,' *J. Phys.: Condens. Matter* **29**, 234005 (2017).
- <sup>176</sup>M. N. van der Linden, J. P. K. Doye, and A. A. Louis, 'Formation of dodecagonal quasicrystals in two-dimensional systems of patchy particles,' *J. Chem. Phys.* **136**, 054904 (2012).
- <sup>177</sup>D. F. Tracey, E. G. Noya, and J. P. K. Doye, 'Programming patchy particles to form three-dimensional dodecagonal quasicrystals,' *J. Chem. Phys.* **154**, 194505 (2021).
- <sup>178</sup>A. Reinhardt, J. S. Schreck, F. Romano, and J. P. K. Doye, 'Self-assembly of two-dimensional binary quasicrystals: A possible route to a DNA quasicrystal,' *J. Phys.: Condens. Matter* **29**, 014006 (2017).
- <sup>179</sup>L. Liu, Z. Li, Y. Li, and C. Mao, 'Rational design and self-assembly of two-dimensional, dodecagonal DNA quasicrystals,' *J. Am. Chem. Soc.* **141**, 4248–4251 (2019).
- <sup>180</sup>G. Malescio and F. Sciortino, 'Self-assembly of quasicrystals and their approximants in fluids with bounded repulsive core and competing interactions,' *J. Mol. Liq.* **349**, 118209 (2022).
- <sup>181</sup>A. Haji-Akbari, M. Engel, A. S. Keys, X. Zheng, R. G. Petschek, P. Palffy-Muhoray, and S. C. Glotzer, 'Disordered, quasicrystalline and crystalline phases of densely packed tetrahedra,' *Nature* **462**, 773–777 (2009).
- <sup>182</sup>A. Haji-Akbari, M. Engel, and S. C. Glotzer, 'Degenerate quasicrystal of hard triangular bipyramids,' *Phys. Rev. Lett.* **107**, 215702 (2011).
- <sup>183</sup>A. Haji-Akbari, M. Engel, and S. C. Glotzer, 'Phase diagram of hard tetrahedra,' *J. Chem. Phys.* **135**, 194101 (2011).
- <sup>184</sup>K. Jiang and P. Zhang, 'Numerical methods for quasicrystals,' *J. Comput. Phys.* **256**, 428–440 (2014).
- <sup>185</sup>M. Oxborrow and C. L. Henley, 'Random square-triangle tilings: A model for twelvefold-symmetric quasicrystals,' *Phys. Rev. B* **48**, 6966–6998 (1993).
- <sup>186</sup>H. Pattabhiraman, A. P. Gantapara, and M. Dijkstra, 'On the stability of a quasicrystal and its crystalline approximant in a system of hard disks with a soft corona,' *J. Chem. Phys.* **143**, 164905 (2015).
- <sup>187</sup>A. Kiselev, M. Engel, and H.-R. Trebin, 'Confirmation of the random tiling hypothesis for a decagonal quasicrystal,' *Phys. Rev. Lett.* **109**, 225502 (2012).
- <sup>188</sup>M. Engel, 'Entropic stabilization of tunable planar modulated superstructures,' *Phys. Rev. Lett.* **106**, 095504 (2011).

- 1747 <sup>189</sup>F. Zipoli, V. Viterbo, O. Schilter, L. Kahle, and T. Laino, ‘Prediction of  
1748 phase diagrams and associated phase structural properties,’ *Ind. Eng. Chem.*  
1749 *Res* **61**, 8378–8389 (2022).
- 1750 <sup>190</sup>J. Behler, ‘Constructing high-dimensional neural network potentials: A  
1751 tutorial review,’ *Int. J. Quantum Chem.* **115**, 1032–1050 (2015).
- 1752 <sup>191</sup>J. Behler, ‘Perspective: machine learning potentials for atomistic simulations,’  
1753 *J. Chem. Phys.* **145**, 170901 (2016).
- 1754 <sup>192</sup>J. Behler, ‘First principles neural network potentials for reactive simulations  
1755 of large molecular and condensed systems,’ *Angew. Chem., Int. Ed.* **56**,  
1756 12828–12840 (2017).
- 1757 <sup>193</sup>J. Behler, ‘Four generations of high-dimensional neural network potentials,’  
1758 *Chem. Rev.* **121**, 10037–10072 (2021).
- 1759 <sup>194</sup>A. Jain, S. P. Ong, G. Hautier, W. Chen, W. D. Richards, S. Dacek, S. Cholia,  
1760 D. Gunter, D. Skinner, G. Ceder, and K. A. Persson, ‘Commentary: The  
1761 Materials Project: A materials genome approach to accelerating materials  
1762 innovation,’ *APL Mater.* **1**, 011002 (2013).
- 1763 <sup>195</sup>S. Kirklin, J. E. Saal, B. Meredig, A. Thompson, J. W. Doak, M. Aykol,  
1764 S. Rühl, and C. Wolverton, ‘The Open Quantum Materials Database  
1765 (OQMD): assessing the accuracy of DFT formation energies,’ *npj Comput.*  
1766 *Mater.* **1**, 15010 (2015).
- 1767 <sup>196</sup>G. Bergerhoff, R. Hundt, R. Sievers, and I. D. Brown, ‘The Inorganic Crystal  
1768 Structure Data Base,’ *J. Chem. Inf. Comput. Sci.* **23**, 66–69 (1983).
- 1769 <sup>197</sup>G. G. C. Peterson and J. Brgoch, ‘Materials discovery through machine  
1770 learning formation energy,’ *J. Phys. Energy* **3**, 022002 (2021).
- 1771 <sup>198</sup>A. Anelli, E. A. Engel, C. J. Pickard, and M. Ceriotti, ‘Generalized convex  
1772 hull construction for materials discovery,’ *Phys. Rev. Materials* **2**, 103804  
1773 (2018).
- 1774 <sup>199</sup>C. J. Bartel, S. L. Millican, A. M. Deml, J. R. Rumpitz, W. Tumas, A. W.  
1775 Weimer, S. Lany, V. Stevanović, C. B. Musgrave, and A. M. Holder,  
1776 ‘Physical descriptor for the Gibbs energy of inorganic crystalline solids and  
1777 temperature-dependent materials chemistry,’ *Nat. Commun.* **9**, 4168 (2018).
- 1778 <sup>200</sup>D. Jha, L. Ward, A. Paul, W.-k. Liao, A. Choudhary, C. Wolverton, and  
1779 A. Agrawal, ‘ElemNet: Deep learning the chemistry of materials from only  
1780 elemental composition,’ *Sci. Rep.* **8**, 17593 (2018).
- 1781 <sup>201</sup>R. E. A. Goodall and A. A. Lee, ‘Predicting materials properties without  
1782 crystal structure: deep representation learning from stoichiometry,’ *Nat.*  
1783 *Commun.* **11**, 6280 (2020).
- 1784 <sup>202</sup>L. Ward, A. Dunn, A. Faghaninia, N. E. Zimmermann, S. Bajaj, Q. Wang,  
1785 J. Montoya, J. Chen, K. Bystrom, M. Dylla, K. Chard, M. Asta, K. A.  
1786 Persson, G. J. Snyder, I. Foster, and A. Jain, ‘Matminer: An open source  
1787 toolkit for materials data mining,’ *Comput. Mater. Sci.* **152**, 60–69 (2018).
- 1788 <sup>203</sup>S. Gong, S. Wang, T. Xie, W. H. Chae, R. Liu, Y. Shao-Horn, and J. C.  
1789 Grossman, ‘Calibrating DFT formation enthalpy calculations by multifidelity  
1790 machine learning,’ *JACS Au* **2**, 1964–1977 (2022).
- 1791 <sup>204</sup>S. Lotfi, Z. Zhang, G. Viswanathan, K. Fortenberry, A. Mansouri Tehrani,  
1792 and J. Brgoch, ‘Targeting productive composition space through machine-  
1793 learning-directed inorganic synthesis,’ *Matter* **3**, 261–272 (2020).
- 1794 <sup>205</sup>K. K. Yalamanchi, M. Monge-Palacios, V. C. O. van Oudenhoven, X. Gao,  
1795 and S. M. Sarathy, ‘Data science approach to estimate enthalpy of formation  
1796 of cyclic hydrocarbons,’ *J. Phys. Chem. A* **124**, 6270–6276 (2020).
- 1797 <sup>206</sup>S. Ubaru, A. Międlar, Y. Saad, and J. R. Chelikowsky, ‘Formation enthalpies  
1798 for transition metal alloys using machine learning,’ *Phys. Rev. B* **95**, 214102  
1799 (2017).
- 1800 <sup>207</sup>C. J. Bartel, A. Trewartha, Q. Wang, A. Dunn, A. Jain, and G. Ceder, ‘A  
1801 critical examination of compound stability predictions from machine-learned  
1802 formation energies,’ *npj Comput. Mater.* **6**, 97 (2020).
- 1803 <sup>208</sup>R. Ouyang, S. Curtarolo, E. Ahmetcik, M. Scheffler, and L. M. Ghiringhelli,  
1804 ‘SISSO: A compressed-sensing method for identifying the best  
1805 low-dimensional descriptor in an immensity of offered candidates,’ *Phys.*  
1806 *Rev. Materials* **2**, 083802 (2018).
- 1807 <sup>209</sup>F. Noé, S. Olsson, J. Köhler, and H. Wu, ‘Boltzmann generators: Sampling  
1808 equilibrium states of many-body systems with deep learning,’ *Science* **365**,  
1809 eaaw1147 (2019).
- 1810 <sup>210</sup>P. Wirtnsberger, G. Papamakarios, B. Ibarz, S. Racanière, A. J. Ballard,  
1811 A. Pritzel, and C. Blundell, ‘Normalizing flows for atomic solids,’ *Mach.*  
1812 *Learn.: Sci. Technol.* **3**, 025009 (2022).
- 1813 <sup>211</sup>M. S. Albero, G. Kanwar, and P. E. Shanahan, ‘Flow-based generative  
1814 models for Markov chain Monte Carlo in lattice field theory,’ *Phys. Rev. D*  
1815 **100**, 034515 (2019).
- 1816 <sup>212</sup>G. Kanwar, M. S. Albero, D. Boyda, K. Cranmer, D. C. Hackett, S. Racan-  
1817 ière, D. J. Rezende, and P. E. Shanahan, ‘Equivariant flow-based sampling  
1818 for lattice gauge theory,’ *Phys. Rev. Lett.* **125**, 121601 (2020).
- 1819 <sup>213</sup>E. G. Tabak and C. V. Turner, ‘A family of nonparametric density estimation  
1820 algorithms,’ *Comm. Pure Appl. Math.* **66**, 145–164 (2013).
- 1821 <sup>214</sup>P. Wirtnsberger, A. J. Ballard, G. Papamakarios, S. Abercrombie, S. Racan-  
1822 ière, A. Pritzel, D. Jimenez Rezende, and C. Blundell, ‘Targeted free energy  
1823 estimation via learned mappings,’ *J. Chem. Phys.* **153**, 144112 (2020).
- 1824 <sup>215</sup>S. Kullback and R. A. Leibler, ‘On information and sufficiency,’ *Ann. Math.*  
1825 *Stat.* **22**, 79–86 (1951).
- 1826 <sup>216</sup>G. Deffrennes, K. Terayama, T. Abe, and R. Tamura, ‘A machine learn-  
1827 ing-based classification approach for phase diagram prediction,’ *Mater. Des.*  
1828 **215**, 110497 (2022).
- 1829 <sup>217</sup>Y. Zhao, Y. Cui, Z. Xiong, J. Jin, Z. Liu, R. Dong, and J. Hu, ‘Machine  
1830 learning-based prediction of crystal systems and space groups from inorganic  
1831 materials compositions,’ *ACS Omega* **5**, 3596–3606 (2020).
- 1832 <sup>218</sup>M. Aghaaminiha, S. A. Ghanadian, E. Ahmadi, and A. M. Farnoud, ‘A  
1833 machine learning approach to estimation of phase diagrams for three-  
1834 component lipid mixtures,’ *Biochim. Biophys. Acta, Biomembr.* **1862**,  
1835 183350 (2020).
- 1836 <sup>219</sup>G. van Mierlo, J. R. Jansen, J. Wang, I. Poser, S. J. van Heeringen, and  
1837 M. Vermeulen, ‘Predicting protein condensate formation using machine  
1838 learning,’ *Cell Rep.* **34**, 108705 (2021).
- 1839 <sup>220</sup>K. L. Saar, A. S. Morgunov, R. Qi, W. E. Arter, G. Krainer, A. A. Lee, and  
1840 T. P. J. Knowles, ‘Learning the molecular grammar of protein condensates  
1841 from sequence determinants and embeddings,’ *Proc. Natl. Acad. Sci. U. S. A.*  
1842 **118**, e2019053118 (2021).
- 1843 <sup>221</sup>X. Chu, T. Sun, Q. Li, Y. Xu, Z. Zhang, L. Lai, and J. Pei, ‘Prediction  
1844 of liquid–liquid phase separating proteins using machine learning,’ *BMC*  
1845 *Bioinform.* **23**, 72 (2022).
- 1846 <sup>222</sup>M. Hardenberg, A. Horvath, V. Ambrus, M. Fuxreiter, and M. Vendruscolo,  
1847 ‘Widespread occurrence of the droplet state of proteins in the human  
1848 proteome,’ *Proc. Natl. Acad. Sci. U. S. A.* **117**, 33254–33262 (2020).
- 1849 <sup>223</sup>D. Raimondi, G. Orlando, E. Michiels, D. Pakravan, A. Bratek-Skicki,  
1850 L. Van Den Bosch, Y. Moreau, F. Rousseau, and J. Schymkowitz, ‘In silico  
1851 prediction of in vitro protein liquid–liquid phase separation experiments  
1852 outcomes with multi-head neural attention,’ *Bioinformatics* **37**, 3473–3479  
1853 (2021).
- 1854 <sup>224</sup>J. Russo, F. Romano, and H. Tanaka, ‘New metastable form of ice and its  
1855 role in the homogeneous crystallization of water,’ *Nat. Mater.* **13**, 733–739  
1856 (2014).
- 1857 <sup>225</sup>K. E. Blow, D. Quigley, and G. C. Sosso, ‘The seven deadly sins: When  
1858 computing crystal nucleation rates, the devil is in the details,’ *J. Chem. Phys.*  
1859 **155**, 040901 (2021).
- 1860 <sup>226</sup>L. Vrbka and P. Jungwirth, ‘Homogeneous freezing of water starts in the  
1861 subsurface,’ *J. Phys. Chem. B* **110**, 18126–18129 (2006).
- 1862 <sup>227</sup>D. Quigley and P. M. Rodger, ‘Metadynamics simulations of ice nucleation  
1863 and growth,’ *J. Chem. Phys.* **128**, 154518 (2008).
- 1864 <sup>228</sup>A. V. Brukhno, J. Anwar, R. Davidchack, and R. Handel, ‘Challenges in  
1865 molecular simulation of homogeneous ice nucleation,’ *J. Phys.: Condens.*  
1866 *Matter* **20**, 494243 (2008).
- 1867 <sup>229</sup>E. Pluhařová, L. Vrbka, and P. Jungwirth, ‘Effect of surface pollution on  
1868 homogeneous ice nucleation: A molecular dynamics study,’ *J. Phys. Chem. C*  
1869 **114**, 7831–7838 (2010).
- 1870 <sup>230</sup>P. Pirzadeh and P. G. Kusalik, ‘On understanding stacking fault formation  
1871 in ice,’ *J. Am. Chem. Soc.* **133**, 704–707 (2011).
- 1872 <sup>231</sup>T. Li, D. Donadio, G. Russo, and G. Galli, ‘Homogeneous ice nucleation  
1873 from supercooled water,’ *Phys. Chem. Chem. Phys.* **13**, 19807 (2011).
- 1874 <sup>232</sup>E. B. Moore and V. Molinero, ‘Is it cubic? Ice crystallization from deeply  
1875 supercooled water,’ *Phys. Chem. Chem. Phys.* **13**, 20008–20016 (2011).
- 1876 <sup>233</sup>T. L. Malkin, B. J. Murray, A. V. Brukhno, J. Anwar, and C. G. Salzmann,  
1877 ‘Structure of ice crystallized from supercooled water,’ *Proc. Natl. Acad. Sci.*  
1878 *U. S. A.* **109**, 1041–1045 (2012).
- 1879 <sup>234</sup>S. J. Cox, S. M. Kathmann, J. A. Purton, M. J. Gillan, and A. Michaelides,  
1880 ‘Non-hexagonal ice at hexagonal surfaces: The role of lattice mismatch,’  
1881 *Phys. Chem. Chem. Phys.* **14**, 7944–7949 (2012).
- 1882 <sup>235</sup>B. J. Murray, D. O’Sullivan, J. D. Atkinson, and M. E. Webb, ‘Ice nucleation  
1883 by particles immersed in supercooled cloud droplets,’ *Chem. Soc. Rev.* **41**,  
1884 6519–6554 (2012).
- 1885 <sup>236</sup>A. Reinhardt and J. P. K. Doye, ‘Free energy landscapes for homogeneous  
1886 nucleation of ice for a monatomic water model,’ *J. Chem. Phys.* **136**, 054501  
1887 (2012).
- 1888 <sup>237</sup>G. Bullock and V. Molinero, ‘Low-density liquid water is the mother  
1889 of ice: On the relation between mesostructure, thermodynamics and ice  
1890 crystallization in solutions,’ *Faraday Discuss.* **167**, 371–388 (2013).
- 1891 <sup>238</sup>A. Reinhardt and J. P. K. Doye, ‘Note: Homogeneous TIP4P/2005 ice  
1892 nucleation at low supercooling,’ *J. Chem. Phys.* **139**, 096102 (2013).
- 1893 <sup>239</sup>E. Sanz, C. Vega, J. R. Espinosa, R. Caballero-Bernal, J. L. F. Abascal, and  
1894 C. Valeriani, ‘Homogeneous ice nucleation at moderate supercooling from  
1895 molecular simulation,’ *J. Am. Chem. Soc.* **135**, 15008–15017 (2013).
- 1896 <sup>240</sup>A. Reinhardt and J. P. K. Doye, ‘Effects of surface interactions on hetero-

- geneous ice nucleation for a monatomic water model,' *J. Chem. Phys.* **141**, 084501 (2014).
- <sup>241</sup>J. R. Espinosa, E. Sanz, C. Valeriani, and C. Vega, 'Homogeneous ice nucleation evaluated for several water models,' *J. Chem. Phys.* **141**, 18c529 (2014).
- <sup>242</sup>L. Lupi, A. Hudait, and V. Molinero, 'Heterogeneous nucleation of ice on carbon surfaces,' *J. Am. Chem. Soc.* **136**, 3156–3164 (2014).
- <sup>243</sup>L. Ickes, A. Welti, C. Hoose, and U. Lohmann, 'Classical nucleation theory of homogeneous freezing of water: thermodynamic and kinetic parameters,' *Phys. Chem. Chem. Phys.* **17**, 5514–5537 (2015).
- <sup>244</sup>T. L. Malkin, B. J. Murray, C. G. Salzmann, V. Molinero, S. J. Pickering, and T. F. Whale, 'Stacking disorder in ice *i*,' *Phys. Chem. Chem. Phys.* **17**, 60–76 (2015).
- <sup>245</sup>J. R. Espinosa, J. M. Young, H. Jiang, D. Gupta, C. Vega, E. Sanz, P. G. Debenedetti, and A. Z. Panagiotopoulos, 'On the calculation of solubilities via direct coexistence simulations: Investigation of nacl aqueous solutions and lennard-jones binary mixtures,' *J. Chem. Phys.* **145**, 154111 (2016).
- <sup>246</sup>B. Cheng, C. Dellago, and M. Ceriotti, 'Theoretical prediction of the homogeneous ice nucleation rate: disentangling thermodynamics and kinetics,' *Phys. Chem. Chem. Phys.* **20**, 28732–28740 (2018).
- <sup>247</sup>M. Fitzner, G. C. Sosso, S. J. Cox, and A. Michaelides, 'Ice is born in low-mobility regions of supercooled liquid water,' *Proc. Natl. Acad. Sci. U. S. A.* **116**, 2009–2014 (2019).
- <sup>248</sup>S. Hussain and A. Haji-Akbari, 'Role of nanoscale interfacial proximity in contact freezing in water,' *J. Am. Chem. Soc.* **143**, 2272–2284 (2021).
- <sup>249</sup>G. C. Sosso, P. Sudera, A. T. Backes, T. F. Whale, J. Fröhlich-Nowoisky, M. Bonn, A. Michaelides, and E. H. G. Backus, 'The role of structural order in heterogeneous ice nucleation,' *Chem. Sci.* **13**, 5014–5026 (2022).
- <sup>250</sup>M. B. Davies, M. Fitzner, and A. Michaelides, 'Accurate prediction of ice nucleation from room temperature water,' *Proc. Natl. Acad. Sci. U. S. A.* **119**, e2205347119 (2022).
- <sup>251</sup>F. Martelli and J. C. Palmer, 'Signatures of sluggish dynamics and local structural ordering during ice nucleation,' *J. Chem. Phys.* **156**, 114502 (2022).
- <sup>252</sup>F. N. Isenrich, N. Shardt, M. Rösch, J. Nette, S. Stavarakis, C. Marcolli, Z. A. Kanji, A. J. deMello, and U. Lohmann, 'The Microfluidic Ice Nuclei Counter Zürich (MINCZ): a platform for homogeneous and heterogeneous ice nucleation,' *Atmos. Meas. Tech.* **15**, 5367–5381 (2022).
- <sup>253</sup>D. W. Oxtoby, 'Homogeneous nucleation: theory and experiment,' *Acc. Chem. Res.* **10**, 897 (1998).
- <sup>254</sup>P. Geiger and C. Dellago, 'Neural networks for local structure detection in polymorphic systems,' *J. Chem. Phys.* **139**, 164105 (2013).
- <sup>255</sup>E. Boattini, M. Ram, F. Smallenburg, and L. Filion, 'Neural-network-based order parameters for classification of binary hard-sphere crystal structures,' *Mol. Phys.* **116**, 3066–3075 (2018).
- <sup>256</sup>C. Dietz, T. Kretz, and M. H. Thoma, 'Machine-learning approach for local classification of crystalline structures in multiphase systems,' *Phys. Rev. E* **96**, 011301 (2017).
- <sup>257</sup>R. S. DeFever, C. Targonski, S. W. Hall, M. C. Smith, and S. Sarupria, 'A generalized deep learning approach for local structure identification in molecular simulations,' *Chem. Sci.* **10**, 7503–7515 (2019).
- <sup>258</sup>A. J. Mukhtyar and F. A. Escobedo, 'Computing free energy barriers for the nucleation of complex network mesophases,' *J. Chem. Phys.* **156**, 034502 (2022).
- <sup>259</sup>W. F. Reinhart, 'Unsupervised learning of atomic environments from simple features,' *Comput. Mater. Sci.* **196**, 110511 (2021).
- <sup>260</sup>A. Statt, D. C. Kleeblatt, and W. F. Reinhart, 'Unsupervised learning of sequence-specific aggregation behavior for a model copolymer,' *Soft Matter* **17**, 7697–7707 (2021).
- <sup>261</sup>R. B. Jadrlich, B. A. Lindquist, and T. M. Truskett, 'Unsupervised machine learning for detection of phase transitions in off-lattice systems. I. Foundations,' *J. Chem. Phys.* **149**, 194109 (2018).
- <sup>262</sup>R. B. Jadrlich, B. A. Lindquist, W. D. Piñeros, D. Banerjee, and T. M. Truskett, 'Unsupervised machine learning for detection of phase transitions in off-lattice systems. II. Applications,' *J. Chem. Phys.* **149**, 194110 (2018).
- <sup>263</sup>B. Monserrat, J. G. Brandenburg, E. A. Engel, and B. Cheng, 'Liquid water contains the building blocks of diverse ice phases,' *Nat. Commun.* **11**, 5757 (2020).
- <sup>264</sup>S. Chiesa, P. M. Derlet, and S. L. Dudarev, 'Free energy of a  $\langle 110 \rangle$  dumbbell interstitial defect in bcc Fe: Harmonic and anharmonic contributions,' *Phys. Rev. B* **79**, 214109 (2009).
- <sup>265</sup>J. Luo, C. Zhou, Q. Li, and L. Liu, 'Thermodynamic formation properties of point defects in germanium crystal,' *Materials* **15**, 4026 (2022).
- <sup>266</sup>R. K. R. Addula and S. N. Punnathanam, 'Computation of solid–fluid interfacial free energy in molecular systems using thermodynamic integration,' *J. Chem. Phys.* **153**, 154504 (2020).
- <sup>267</sup>S. R. Yeandel, C. L. Freeman, and J. H. Harding, 'A general method for calculating solid/liquid interfacial free energies from atomistic simulations: Application to  $\text{CaSO}_4 \cdot x\text{H}_2\text{O}$ ,' *J. Chem. Phys.* **157**, 084117 (2022).
- <sup>268</sup>R. S. DeFever and E. J. Maginn, 'Computing the liquidus of binary monatomic salt mixtures with direct simulation and alchemical free energy methods,' *J. Phys. Chem. A* **125**, 8498–8513 (2021).
- <sup>269</sup>A. Rahbari, R. Hens, D. Dubbeldam, and T. J. H. Vlugt, 'Improving the accuracy of computing chemical potentials in CFMCM simulations,' *Mol. Phys.* **117**, 3493–3508 (2019).
- <sup>270</sup>B. Cheng, 'Computing chemical potentials of solutions from structure factors,' *J. Chem. Phys.* **157**, 121101 (2022).
- <sup>271</sup>A. P. Abbott, D. Boothby, G. Capper, D. L. Davies, and R. K. Rasheed, 'Deep eutectic solvents formed between choline chloride and carboxylic acids: versatile alternatives to ionic liquids,' *J. Am. Chem. Soc.* **126**, 9142–9147 (2004).
- <sup>272</sup>N. Singh, S. Das, N. Singh, and T. Agrawal, 'Computer simulation, thermodynamic and microstructural studies of benzamide–benzoic acid eutectic system,' *J. Cryst. Growth* **310**, 2878–2884 (2008).
- <sup>273</sup>A. González de Castilla, J. P. Bittner, S. Müller, S. Jakobtorweihen, and I. Smirnova, 'Thermodynamic and transport properties modeling of deep eutectic solvents: A review on gE-models, equations of state, and molecular dynamics,' *J. Chem. Eng. Data* **65**, 943–967 (2020).
- <sup>274</sup>P. V. A. Pontes, E. A. Crespo, M. A. R. Martins, L. P. Silva, C. M. S. S. Neves, G. J. Maximo, M. D. Hubinger, E. A. C. Batista, S. P. Pinho, J. A. Coutinho, G. Sadowski, and C. Held, 'Measurement and PC-SAFT modeling of solid-liquid equilibrium of deep eutectic solvents of quaternary ammonium chlorides and carboxylic acids,' *Fluid Phase Equilib.* **448**, 69–80 (2017).
- <sup>275</sup>M. A. R. Martins, E. A. Crespo, P. V. A. Pontes, L. P. Silva, M. Bülow, G. J. Maximo, E. A. C. Batista, C. Held, S. P. Pinho, and J. A. P. Coutinho, 'Tunable hydrophobic eutectic solvents based on terpenes and monocarboxylic acids,' *ACS Sustainable Chem. Eng.* **6**, 8836–8846 (2018).
- <sup>276</sup>P. Vainikka, S. Thallmair, P. C. T. Souza, and S. J. Marrink, 'Martini 3 coarse-grained model for type III deep eutectic solvents: Thermodynamic, structural, and extraction properties,' *ACS Sustainable Chem. Eng.* **9**, 17338–17350 (2021).