

**A Theoretical Study on
the Molecular Mechanism of Ice Melting and
the Local Structure of Aqueous Solutions**

Kenji Mochizuki

Doctor of Philosophy

Department of Functional Molecular Science
School of Physical Science
The Graduate University for Advanced Studies

2013 (School Year)

Contents

1. General Introduction	3
References	8
2. Defect pair separation as the controlling step in homogeneous ice melting	11
2.1 Introduction	12
2.2 Simulation details	14
2.3 Results and discussions	15
2.4 Appendix	24
2.4.1 Supplementary figures	24
2.4.2 Supplementary methods	28
References	37
3. Local structure of methanol-water binary solutions studied by soft X-ray absorption spectroscopy and molecular dynamics simulation	39
3.1 Introduction	40
3.2 Experiments	43
3.3 Results and discussions	44
3.3.1 Oxygen K-edge XAS	44
3.3.2 Carbon K-edge XAS	46
3.3.3 MD simulation	50
3.3.4 Structures of methanol-water mixtures	55
3.4 Conclusions	57
References	59

4. A conformational factorisation approach for estimating the binding free energies of macromolecules	63
4.1 Introduction	64
4.2 Methodology	67
4.2.1 Factorised superposition approach	67
4.2.2 Free energy of local minima	71
4.2.3 Basin-hopping parallel tempering	75
4.3 Application to human aldose reductase	79
4.3.1 Simulation set up	79
4.3.2 Systematic rigidification	82
4.3.3 Sampling local minima	83
4.3.4 Convergence of the free energy with the size of the unconstrained region	86
4.3.5 Computational cost	87
4.3.6 Incorporating solvent effects	88
4.4 Conclusions	89
4.5 Appendix	91
4.5.1 Hessian in the local rigid body coordinates	91
4.5.2 Supporting tables and figures	93
References	96
 5. Summary	 99

Chapter 1

General Introduction

Water is the most ubiquitous substance on earth and is vital for the activity of life, playing critical roles in the biological process and in the climate of the globe. Water has various anomalous physical properties,¹ which arise from the nature of the hydrogen bond (HB). For example, density of water is larger than that of ice, and becomes the maximum at 277K and decrease with decreasing temperature. Water has the very high melting ($T_M=273.15\text{K}$) and boiling temperatures ($T_B=373.15\text{K}$) in comparison with the isoelectronic species, such as hydrogen sulfide H_2S ($T_M=187.65\text{K}$ and $T_B=212.45\text{K}$). Thus, all three phases are present in the ambient condition. The solid phase, ice, exhibits more than ten kinds of crystal structures, depending on the temperature and pressure.² In addition, recently a new ice phase, plastic ice, has been discovered by computer simulation,³ and unknown phases of ice might still exist.

The anomalies of liquid water and ice arise from the intermediate strength and highly directional nature of HB. HB has the specific attraction character between electronegative atoms (O,F,N) and hydrogen atom, and its strength is in between that of van der Waals interaction and of covalent bond. A water molecule, having H-O-H angle 104.5° , can make the HBs with four surrounding water molecules and form the tetrahedral three-dimensional HB network, whose ideal tetrahedral angle is 109.5° .

The structure of liquid water was first observed by Bernal and Fowler in 1933⁴ by using X-ray diffraction method. Based on their experimental data, they proposed that a water molecule is tetra-coordinated in liquid phase.⁴ In 1938 Morgan and Warren claimed that many HBs connecting neighboring molecules must be broken in order to interpret their radial distribution function data, if the firm tetrahedral structure is assumed.⁵ In 1951, Pople suggested that the majority of HBs are distorted rather than broken.⁶ This Pople's model was extended to the random network model by Bernal.⁷ Since then, various models had been proposed.⁸

In 1971, the first molecular dynamics (MD) study of liquid water performed by Rahmann and Stillinger gave the direct molecular level insight.⁹ They showed that the liquid structure consists of a highly strained random HB network which bears little structural resemblance to known aqueous crystals. By the development of computer simulations, the cooperative motions and fluctuations in liquid water associated with the HB rearrangement dynamics were investigated.¹⁰⁻¹²

The HB structure of the ice (ice Ih) was first proposed by Pauling in 1935. The oxygen atoms are arranged on a hexagonal lattice. Water molecules are linked to one another by HBs, each molecule offering its hydrogen atoms to two other molecules and accepting HBs from other two. Thus, the strong direction dependence of HB requires water molecule in ice to occupy the

tetrahedral lattice position with a specific HB direction.

How does liquid water find out this specific ice structure from infinite number of possible conformations without any surface/interface? In 2002, M. Matsumoto et.al, have first performed the MD simulation of the freezing process and found that the crystallization proceeds in two steps: First, a metastable amorphous cluster appears, then it transforms into crystalline structure.¹³ The subsequent computational studies for the spontaneous crystallization of methane hydrate also showed the similar mechanism.^{14,15}

Despite the fact that the computer simulations have revealed the homogeneous freezing process of ice at molecular level, the opposite process, homogeneous melting of ice, remains poorly understood. The thermally induced homogeneous melting of solids is fairly well understood, and involves the formation and growth of melting nuclei.¹⁶⁻²⁰ But in the case of water, resilient HBs render ice melting more complex. We know that the first defects appearing during homogeneous ice melting are pairs of five- and seven-membered rings, which appear and disappear repeatedly and randomly in space and time in the crystalline ice structure.²¹⁻²³ However, the accumulation of these defects to form an aggregate is nearly additive in energy, and results in a steep free energy increase that suppresses further growth. In Chapter 2, I report that molecular dynamics simulations of homogeneous ice melting identify as a crucial first step not the formation but rather the spatial separation of a defect pair. I find that once it is separated, the defect pair, either an interstitial (I) and a vacancy (V) defect pair (a Frenkel pair), or an L and a D defect pair (a Bjerrum pair),²⁴ is entropically stabilized, or ‘entangled’. In this state, the defects with threefold HB coordination persist and grow, and thereby prepare the system for subsequent rapid melting.

Besides its crucial roles in pure water or ice, the HB structure and dynamics around surfaces of amphiphile molecules are the keys to understand processes involving hydrophobic/hydrophilic interactions in very different contexts, e.g., aqueous solution of alcohols, molecular self-assembling, protein folding, ligand/protein docking, as well as nanoscale water confinement and surface wetting.²⁵⁻²⁷ When the amphiphile molecules are mixed with water, the hydrophilic region of these molecules makes the stable interaction with water while the hydrophobic region tends to adhere to each other to minimize their exposure to water. The smallest amphiphile molecule is methanol, which contains both hydrophilic group (OH) and hydrophobic group (CH₃). It is well known that water mixes with methanol at any mole fraction but the entropy increase of the water-methanol mixture is far less than that expected for an ideal solution of randomly mixed molecules.²⁸ This feature of the mixing entropy had been explained by, for example, using clathrate-like structure

models.²⁸ Recently a lot of experimental²⁹ and theoretical studies³⁰ have been performed on this system and showed that water and methanol does mix but incompletely at the molecular level. It has not been well explored, however, how the microscopic structure of methanol-water mixture changes against mole fraction. In Chapter 3, I present the results on the local structure of methanol-water binary solution at different concentrations investigated by the O and C K-edge X-ray Absorption Spectroscopy (XAS) and the MD simulation. I have found that methanol-water mixture exhibits three different kinds of the local structures against the methanol molar fraction, X ; the slopes of the XAS C K-edge spectral intensity against molar fraction changes at $X=0.3$ and $X=0.7$. It is found that the HB network structure among water molecules becomes non-percolated above $X=0.3$, and the HB network among water molecules and alcohol molecules form cluster segments and the HB among these cluster segments are broken above $X=0.7$. On the other hand the pre-edge feature in the O K-edge XAS is found to show almost linear dependence on the concentration.

In addition to these studies on water and aqueous solutions, a new approach for estimating the binding free energies of ligand/protein docking is investigated. Calculating the binding affinities using atomistic simulations can provide detailed molecular level insight into molecular recognition mechanism, contributions of water or HBs network in biomolecule system, and help to inform fields such as structure-based drug design^{31,32} and self-assembly.^{33,34} There are a lot of approaches to estimate the binding affinity, such as docking and scoring approach, thermodynamic integration, free energy perturbation and many other methods. Molecular Mechanics/Poisson Boltzmann Surface Area (MM-PBSA), which relies on MD simulation of only the free and bound states, is seen as one of the most applicative approaches to estimate the free energy of small molecules.³⁵⁻³⁷ However, for protein size systems, it hardly harvests adequate conformations in current computer power because the conformational space for sampling increases exponentially with the system size and in addition the trajectory is easily trapped in local minima of the potential. The superposition approach provides an alternative formulation for global thermodynamics within the energy landscape framework, which harvests the local energy minima and use the harmonic approximation around each minimum.³⁸ This approach is faster than the MM-PBSA and has been successfully applied to various small host-guest systems. However, for macromolecules, it has the same problems as the MM-PBSA approach. In Chapter 4, I present a conformational factorization method to improve the sampling efficiency based on the superposition approach. In this method, the number of minima needed to be sampled is greatly reduced by fixing local configurations which are sufficiently distant from the binding site. Furthermore, the basin-hopping parallel tempering³⁹ and the local rigid body framework⁴⁰ are also employed in this method to sample potential energy

minima. I benchmark this approach for human aldose reductase (PDB code 2INE). When varying the size of the rigid region, the free energy difference converges for factorization of groups at a distance of 14 Å from the binding site, which corresponds to 80% of the protein being locally rigidified.

References

- [1] Eisenberg, D., Kauzmann, W., *The Structure and Properties of Water*, Oxford (2005)
- [2] Petrenko, V. F., Whitworth, R. W., *Physics of Ice*, Clarendon Press (1999)
- [3] Takii, Y., Koga, K., Tanaka, H., *J. Chem. Phys.*, **128**, 204501 (2008)
- [4] Bernal, J. D., Fowler, R. H. A., *J. Chem. Phys.*, **1**, 515-548 (1933)
- [5] Morgan, J., Warren, B. E., *J. Chem. Phys.*, **6**, 666-673 (1938)
- [6] Pople, J. A., *Proc. Roy. Soc. A*, **205**, 163-178 (1951)
- [7] Bernal, J. D., *Proc. Roy. Soc. A*, **280**, 299-322 (1964)
- [8] Frank, H. S., Wen, W. Y., *Discuss. Faraday Soc.*, **24**, 133-140 (1957)
- [9] Rahman, A., *J. Chem. Phys.*, **55**, 3336 (1971)
- [10] Stillinger, F. H., Weber, T. A., *Phys. Rev. A*, **25**, 978-989 (1982)
- [11] Tanaka, H., Ohmine, I., *J. Chem. Phys.*, **91**, 6318-6327 (1989)
- [12] Ohmine, I., Tanaka, H., *J. Chem. Phys.*, **93**, 8138-8147 (1990)
- [13] Matsumoto, M., Saito, S., Ohmine, I., *Nature*, **416**, 409-413 (2002)
- [14] Jacobson, L. C., Hujo, W., Molinero, V., *J. Am. Chem. Soc.*, **132**, 11806-11811 (2010)
- [15] Walsh, M. R. et al., *Science*, **326**, 1095-1098 (2009)
- [16] Iglev, H. et al., *Nature*, **439**, 183-186 (2006)
- [17] Fecht, H. J., *Nature*, **356**, 133-135 (1992)
- [18] Cahn, R. W., *Nature*, **413**, 582-583 (2001)
- [19] Forsblom, M., Grimvall, G., *Nature Mater.*, **4**, 388-390 (2005)
- [20] Jin, Z. et al., *Phys. Rev. Lett.*, **87**, 055703 (2001)
- [21] Tanaka, H., Mohanty, J., *J. Am. Chem. Soc.*, **124**, 8085-8089 (2002)
- [22] Grishina, N., Buch, V., *J. Chem. Phys.*, **120**, 5217-5225 (2004)
- [23] Donadio, D., Raiteri, P., Parrinello, M., *J. Phys. Chem. B*, **109**, 5421-5424 (2005)
- [24] Bjerrum, N., *Science*, **115**, 385-390 (1952)
- [25] Safran, S. A., *Statistical Thermodynamics of Surfaces, Interfaces, and Membranes*, Westview Press (2003)
- [26] Tanford, C., *Science*, **200**, 1012-1018 (1978)
- [27] Chandler, D., *Nature*, **437**, 640-647 (2005)
- [28] Frank, H. S., Evans, M. W., *J. Chem. Phys.*, **13**, 507-532 (1945)
- [29] Dixit, S. et al., *Nature*, **416**, 829-832 (2002)
- [30] da Silva, J. A. B. et al., *Phys. Chem. Chem. Phys.*, **13**, 6452-6461 (2011)
- [31] Jorgensen, W. L., *Science*, **303**, 1813-1818 (2004)

- [32] Michel, J., Foloppe, N., Essex, J. W., *Mol. Inf.*, **29**, 570-578 (2010)
- [33] Johnson, R. R. et al., *Nano Lett.*, **9**, 537-541 (2009)
- [34] Ercolani, G., *J. Am. Chem. Soc.*, **125**, 16097-16103 (2003)
- [35] Brown, S. P., Muchmore, S. W., *J. Chem Inf. Model.*, **46**, 999-1005 (2006)
- [36] Gouda, H. et al., *Biopolymers*, **68**, 16-34 (2003)
- [37] Srinivasan, J. et al., *J. Am. Chem. Soc.*, **120**, 9401-9409 (1998)
- [38] Wales, D. *Energy Landscapes*, Cambridge University Press (2003)
- [39] Strodel, B. et al., *J. Am. Chem. Soc.*, **132**, 13300-13312 (2010)
- [40] Kusumaatmaja, H. et al., *J. Chem. Theory Comput.*, **8**, 5159-5165 (2012).

Chapter 2

Defect pair separation as the controlling step in homogeneous ice melting

K. Mochizuki, M. Matsumoto and I. Ohmine

Nature, **498**, 350-354 (2013)

2.1 Introduction

Upon heating, ice melts to water. This familiar phase transition is usually initiated at the surface of impurities at the melting point. Such a process is called heterogeneous melting. On the other hand, if the surface melting is suppressed by some means or ice is heated from inside, for example, by laser irradiation, ice spontaneously melts by thermal fluctuation with considerable superheating, which is called homogeneous melting.¹ How does superheated ice start melting under such an ideal condition? It sounds a trivial question, but is not quite simple in reality. Study on crystal melting give us important and fundamental information on the transition to the disordered state.

Considerable experimental effort has been invested to study the bulk melting.²⁻⁴ It is found that many substances can be superheated beyond their melting points.² For ice, Laubereau et al observed its bulk melting by laser-induced temperature jump and found that the maximum superheating is $330 \pm 10 \text{ K}$ to persist over the monitored time interval of 1.3 ns,¹ although, at the limit of superheating, the ultrafast melting process takes place.³ Even at present, it is not possible to directly observe the very molecular process in initial stage of the melting at moderately superheated temperature, experimentally. Because the initial melting embryo appears randomly after a long induction time and then grows too fast to be observed by any experimental mean.

There have been intensive theoretical investigations on the melting mechanism of solids. Various theoretical criteria for melting have been employed: for examples, Lindemann criterion⁵ and Born criterion.⁶ Lindemann proposed that melting is caused when the thermal displacement of the atoms exceeds a certain fraction of the nearest neighbor distance. Born proposed that a “rigidity catastrophe” with a vanishing elastic shear modulus occurs in the melting. By using molecular dynamics simulation, Jin et al. examined these criteria and found that Lindemann and Born criteria response simultaneously when both are applied to the bulk melting of LJ molecules at the limit of superheating.⁷ Moreover, they argued that thermally destabilized particles form liquid embryos and the coalescence of such embryos sets off the homogeneous melting. Forsblom found that the aggregation of point defects, interstitials and vacancies, initiates melting in aluminum.^{8,9} Although there are the previous theories which have argued for melting mechanisms governed by dislocations (defects),^{10,11} Forsblom first claimed that very small number of defect, which is an aggregation of 6-7 interstitial and 3-4 vacancies, leads the system to melt.

Many studies have been also carried out for ice melting. Buch revealed that specific kinds of defects, one is called 5+7 defect consisting of 5- and 7-member rings, and another is L+D defect,

and their complex appear spontaneously by thermal fluctuation, and predicted that a 5+7 defect is likely to serve as a nucleation center for melting.¹² Parrinello argued that the aggregation of 5+7 defects is the general disordering processes in tetrahedral network-forming materials, that are water and silicon.¹³

Does the simple aggregation of the initial defects, such as 5+7 and L+D defect, actually lead to melting? Once the defect is transformed from the quite stable tetrahedral structure of ice, the system contains the large distortion and the potential energy increases by 20-35 kJ mol⁻¹. However, the system gains less entropic stabilization when the defect is formed, because the strong direction dependence of HB does not permit the large variety of conformations. Thus, the simple aggregation of the initial defects likely forms the sharp free energy well and the system hardly escapes from the global minimum, that is ice. Thus, while the initial defects in ice are well defined, the embryo formation of melting nucleus transformed from these initial defects and the growth of the embryo are poorly understood.

In this chapter, I investigate the homogeneous melting of ice Ih at superheated temperature by using molecular dynamics simulation. Whole process between perfect ice and critical nucleus is studied by various theoretical analyses in molecular level. Especially, I focus on the following questions with dynamical aspect; What and why is the trigger to start melting in the perfect crystals? How does the embryo grow by only thermal fluctuation? When dose liquid appear and what role dose the liquid play? What is responsible for the temperature difference ?

2.2 Simulation details

Molecular dynamics trajectory calculations were performed on a system containing 896 water molecules in an almost cubic cell (edge lengths are 3.142, 3.110 and 2.932 nm) with periodic boundary conditions. I use the TIP4P water model,¹⁴ which is one of the most successful models in terms of reproducing the thermodynamic properties of water.¹⁵ Intermolecular interaction is smoothly truncated from 1.0530 nm to 1.1638 nm. A preparatory calculation of the proton-disordered ice Ih structure evolving for 1 ns at 250 K is followed by step-by-step increase of the temperature; a pre-melting equilibration run is performed for 1 ns at 270 K after a 1 ns run at 260 K, then the melting process is observed for several nanoseconds at 275 K. More than 10 μ s of trajectories in total are used for statistical analyses. Even though 275 K is higher than the melting temperature $T_m = 232$ K¹⁶ predicted by the TIP4P model, 275 K is found to be about the lowest temperature at which the system melts.¹⁷ At a higher temperature, for example 300 K, the crystal collapses as soon as it is heated up, while at 275 K the system shows an induction time of a few nanoseconds before melting starts.

Hundreds of trajectories are started with different initial proton order configurations and different initial velocities applied to water molecular motions in ice. Temperature is controlled by the Nosé-Hoover thermostat.¹⁸ The density, instead of the pressure, of the system is kept constant in this work. The density is set to 0.935 g cm⁻³, which is the density of ice in this model at the melting point under atmospheric pressure. I mainly focus on the molecular mechanism of the initial stage of melting, when the volume of the system has not yet changed substantially.

In order to find energy barriers required for the structural changes of the HB network, reaction coordinate analyses¹⁹ are performed on inherent structures of the system. The inherent structures are obtained by applying the conjugate gradient method²⁰ to the instantaneous structures visited by the trajectories.

2.3 Result and Discussion

Ice melting trajectory calculations were performed by modelling the superheating of crystalline ice Ih²¹ and following the melting process. Fig.1a plots the potential energy per molecule along a typical trajectory. Fig.1b shows the corresponding change in the total number of ‘off-lattice’ water molecules (n), and in the size of the largest cluster of such molecules (n_{LC}) that grows into the melting nucleus. (For snapshots of the corresponding HB structures, see Appendix Fig.A1.) I define water molecules as ‘off-lattice’ if they are more than 0.1 nm away from the nearest lattice point of the ice structure (Appendix Fig.A2), and consider off-lattice molecules within a distance of 0.6 nm as adjacent to each other and forming a cluster, and thereby determine n_{LC} .

In the quiescent period (<2,150 ps, Fig.1), I see in the ice HB network structural defects of five- and seven-membered rings (‘5+7 defect’) and/or pairs of L and D defects (‘L+D complex’) randomly scattered in space.^{12,13,22} (In the ordered ice structure, there is one proton between two oxygen atoms; in the L and D defects, there are respectively no protons and two protons between two oxygen atoms.) Although the appearance of these 5+7 defects and/or L+D complexes is the first step in ice melting, their simple accumulation will not result in melting: because they retain fourfold HB coordination, individual water molecules in the 5+7 defects are strongly restricted in their motion, and entropy will therefore not increase rapidly with the energy increase. Hence, the free energy will rise sharply with an increase in the number of these defects. In the trajectories involving only 5+7 defects and/or L+D complexes, I indeed found that n hardly exceeds 15, and that the system repeatedly exhibits intermittent creation and annihilation of small-sized melting clusters.

In trajectories resulting in melting, the defect growing to form the melting nucleus is either an I defect spatially separated from its accompanying V defect, or a D defect spatially separated from its accompanying L defect (see Appendix Fig.A3 for defect structures). In these separated structures, the I defect encompasses an additional lattice water molecule, while the D defect has two hydrogen atoms between two oxygen atoms and thus breaks the Bernal-Fowler ice rule.²³ Separated defects form occasionally during the recrystallization that occurs after thermal fluctuations have created several 5+7 defects and/or L+D complexes. Although defects usually appear and disappear rapidly, recrystallization occasionally fails to revert back to ice obeying the Bernal-Fowler rule and instead yields separated I and V (or D and L) defects while the rest of the region recovers the original crystalline structure. Clusters of 5+7 defects (or L+D complexes) with $n_{LC} > 5$ accumulate in our simulations every 275ps on average, with 14% forming a separated defect pair and the rest recrystallizing. Once a separated defect pair is created (for example, at 2,150 ps in Fig.1), it

undergoes facile dislocation on the lattice (see below), and the distance between I and V (or D and L) defects increases rapidly. Such a separated pair is hard to annihilate because its recombination would require the right sequence of many HB alternations, and I thus refer to it as entangled. A typical separated I–V defect pair is shown Fig.2a.

The small melting cluster containing an I defect just after separation from its V defect often exhibits rapid motion in the lattice, which for the trajectory shown in Fig.1 lasts from 2,150 to 2,730 ps (see Appendix Fig.A4). After this induction period of about 600 ps, the system takes very little time to reach the liquid state (about 0.3 ns in Fig.1). The average time from initial I–V separation to total melting is about 1 ns, with a Poisson-type time distribution peaking around $t = 0.5$ ns. The fast formation of a critical nucleus and subsequent melting is a direct result of separated I (or D) defects enabling facile HB alternations in water ice (see also below).

To quantify the degree of disorder of the HB network of the melting nucleus, I define the topological ‘edit’ distance d_T as the minimum number of all HB additions and deletions (that is, edits) needed to recover from a given disordered HB structure to the network topology of the closest proton-disordered ice structure²⁴ (Appendix section 2.4.2.1 gives the detailed procedure for estimating d_T .) Formation of a 5+7 defect, for example, introduces two off-lattice molecules and increases d_T by 4 (sometimes 6) because it requires the breaking and also the creation of 2 HBs to recover the original ice structure. When 5+7 defect pairs accumulate, d_T increases by $2n$ and I can then define the excess d_T (denoted as d_T^{ex}) as d_T minus $2n$. d_T^{ex} is a measure of the difficulty of resolving HB disorder, and hence of the degree of HB network entanglement.

Fig.2b shows the time evolution of d_T^{ex} for the melting trajectory in Fig.1. d_T^{ex} stays small when melting clusters consisting only of 5+7 defects and/or L+D complexes appear intermittently in the quiescent period, and increases suddenly at 2,150 ps when a separated I–V pair is created. The melting nucleus containing the I defect then changes its position and size for about 600 ps, with d_T^{ex} fluctuating around larger values. The nucleus finally starts growing rapidly at 2,730 ps (Fig.1), and d_T^{ex} increases again. In this growing process, an I defect often couples with its surrounding 5+7 defects to rapidly convert ice HB network structures into liquid structures. The plot of d_T^{ex} against n_{LC} in Fig.2c shows that the melting trajectory is initially characterized by $d_T^{\text{ex}} < 10$ as n_{LC} fluctuates between 0 and 5 (n fluctuates between 0 and 15), and that d_T^{ex} then increases and fluctuates more strongly after a separated I–V pair is created. With the increase of d_T^{ex} , the defect pair becomes harder to remove, because the larger number of HBs rearrangement is required to recover it. The benefit of a quantitative measure of network disorder is also illustrated by the HB

network structures arbitrarily selected from the trajectory of Fig.1 and shown in Fig.2d and e: although the structure in Fig.2d may look more distorted than in Fig.2e, the latter HB network has the larger $d_{T^{ex}}$ value and is thus characterized by larger network disorder.

I calculate the free energy of the system from the melting trajectories, and plot it against n_{LC} and $d_{T^{ex}}$ in Fig.3a. (The procedure used to estimate the free energy is described in section 2.4.2.2). The contours are such that they enforce a strong direction dependence on the minimum free-energy path to melting: $d_{T^{ex}}$ first increases owing to the formation of a separated pair defect, and only then does n_{LC} , the size of the largest cluster that goes on to form the melting nucleus, increase. The solid-state free-energy minimum is found at $n_{LC} = 4$, which is about the average size of the melting cluster in the quiescent period. The critical nucleus, that is, the saddle point between solid and liquid, is located around $n_{LC} = 50$ and is about 18 kJ mol⁻¹ higher in energy than the solid state. Beyond the critical nucleus size, the free energy decreases monotonically to the liquid state. The size of the critical nucleus decreases with increasing temperature, with our calculations predicting that it is around $n_{LC} = 33$ at 280 K.

This contour map can be projected on n_{LC} to obtain a one-dimensional free-energy surface. The resultant free-energy curve F_W (grey line in Fig.3b) is then decomposed into the free energy of the ‘un-separated stage’ (F_U) and that of the ‘separated stage’ (F_S), with $F_W = -k_B T \ln[\exp(-\beta F_U) + \exp(-\beta F_S)]$, where k_B is Boltzmann’s constant, T absolute temperature and $\beta = 1/k_B T$. F_U is calculated by projecting the ‘un-separated’ part of the contour in Fig.3a with $d_{T^{ex}} \leq 10$, and F_S by projecting the other part of the contour. Fig.3b shows that F_U increases monotonically and sharply with n_{LC} . In contrast, the initial slope of F_S is less than half that of F_U , showing the facile growth of the melting nucleus after defect separation. F_S intersects with F_U at around $n_{LC} = 6$.

I next performed separate molecular dynamics simulations for a system with an extra water molecule added to ice. The I defect of this so called ‘doped system’ cannot be annihilated because there is no V-defect counterpart, so it mimics a long-lived I defect. Fig.3c shows that the free energy of the doped system (F_D) is almost identical to F_S , with this strong resemblance between the two free energies illustrating the key role of separated I defect formation in melting.

The functional form of free energy in classical nucleation theory²⁵ is given by $f(n_{LC}) = a*n_{LC}^{2/3} - b*n_{LC} + c$, where n_{LC} and $n_{LC}^{2/3}$ correspond to the volume and the surface area of the melting nucleus, respectively. Fig.3d illustrates that F_W can be fitted well with a function $f(n_{LC})$ (green line), indicating that ice melting can be described within the classical nucleation theory framework when looking at this stochastic part of the melting process. However, homogeneous ice

melting advances under normal super-heated conditions only when a separated I–V (or separated L–D) pair is involved, and this early molecular melting process is much more intricate and elaborate than the simple stochastic process (accumulation of 5+7 defects) considered by classical nucleation theory. Although also stochastic in nature, it requires at least a two-dimensional description of the free energy, as in Fig.3a. But once a defect pair has been separated and entanglement created, melting is seen to proceed in a near-classical stochastic manner and further destroys the HB network of the crystal.

Fig.4 plots the total potential energy changes²⁶ along a typical melting trajectory as either a separated I defect or un-separated defect clusters grow, along with the potential energy changes of individual water molecules involved in defect growth. Energy barriers¹⁵ are particularly low when an I defect is involved in separated cluster growth, as its dangling bonds change the HB coordination with other water molecules in the lattice to reduce its own energy while destabilizing other water molecules. The total potential energy of the system thus slowly increases with the increase of n_{LC} . The average energy needed to create an additional defect, $\Delta U(n_{LC}) = U(n_{LC}+1) - U(n_{LC})$, is only about 4–5 kJ mol⁻¹ per off-lattice molecule for $n_{LC} = 5$ to 20. For comparison, the experimentally obtained value for liquid water (that is, $n_{LC} = \infty$) at 0 °C is 6 kJ mol⁻¹, calculated from the latent heat. But in the absence of separated I (or D) defects, the energies of all individual molecules directly involved in defect formation and growth increase significantly with the breaking of their HBs to form strongly hindered four HB coordinations. The total energy steadily increases with n_{LC} , and $\Delta U(n_{LC})$ ranges from 9 to 13 kJ mol⁻¹ per off-lattice molecule for $n_{LC} = 5$ to 10. This value is much larger than that in the presence of the separated I defect.

The average entropy term $T\Delta S(n_{LC})$ for creating one additional defect in the lattice is about 4 kJ mol⁻¹ for $n_{LC} = 5$ to 10, which is the difference between the free energy and the average potential energy, and then gradually increases with further growth of the liquid fragments²⁷ (Appendix Fig.A6), whereas $\Delta U(n_{LC})$ remains nearly constant at about 5 kJ mol⁻¹ for $n_{LC} > 5$. This results in the difference between ΔU and $T\Delta S$ gradually decreasing, and in $T\Delta S$ surpassing ΔU at the critical nucleus size, so the free-energy surface has a gentle and convex upward slope up to the critical nucleus size (Fig.3b). This feature of the potential energy surface and the fact that it is smooth with small energy barriers ensure that melting proceeds rapidly once separated I defects have formed and the system has gone through the subsequent short induction period.

I note that although I have shown that homogeneous melting under normal super-heating conditions only proceeds when separated I–V (or separated D–L) defect pairs are created, very high

temperatures will simply induce total collapse of the ice network.²⁸ But separated defect pairs may possibly also play a role during the very late stages of water freezing,²⁹⁻³² as separated I–V defects can induce HB reorientations that stabilize the system as a more proton-ordered ice,²² after an overall crystalline structure has been attained.

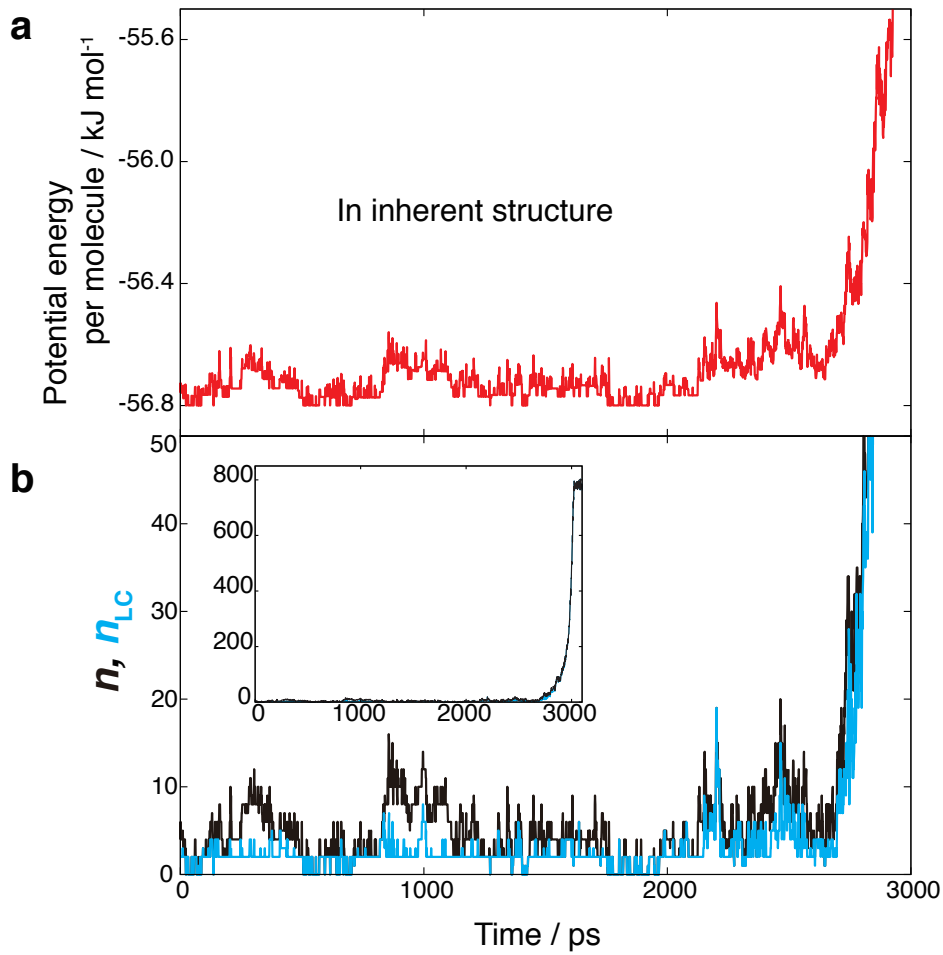


Figure 1 Potential energy per molecule and number of off-lattice molecules in a melting trajectory. (a) Potential energy per molecule of the inherent structures of a typical melting trajectory. At 0 ps, a temperature jump from 270 K to 275 K is imposed. (b) Main panel, total number n of off-lattice water molecules (black line) and size n_{LC} of the largest cluster (blue line) for the trajectory of a. Note that the melting point of the water model (TIP4P) is $T_m = 232$ K, and 275 K is a superheated state. The potential energy fluctuates during the long quiescent period owing to intermittent creation and annihilation of small melting clusters. This lasts until 2,150 ps, when relatively large energy fluctuations appear with the creation of a separated I–V pair. Rapid growth of the melting cluster, and thus rapid energy rise, starts only at 2,730 ps. Inset, data plotted using n up to 800, that is, almost complete melting.

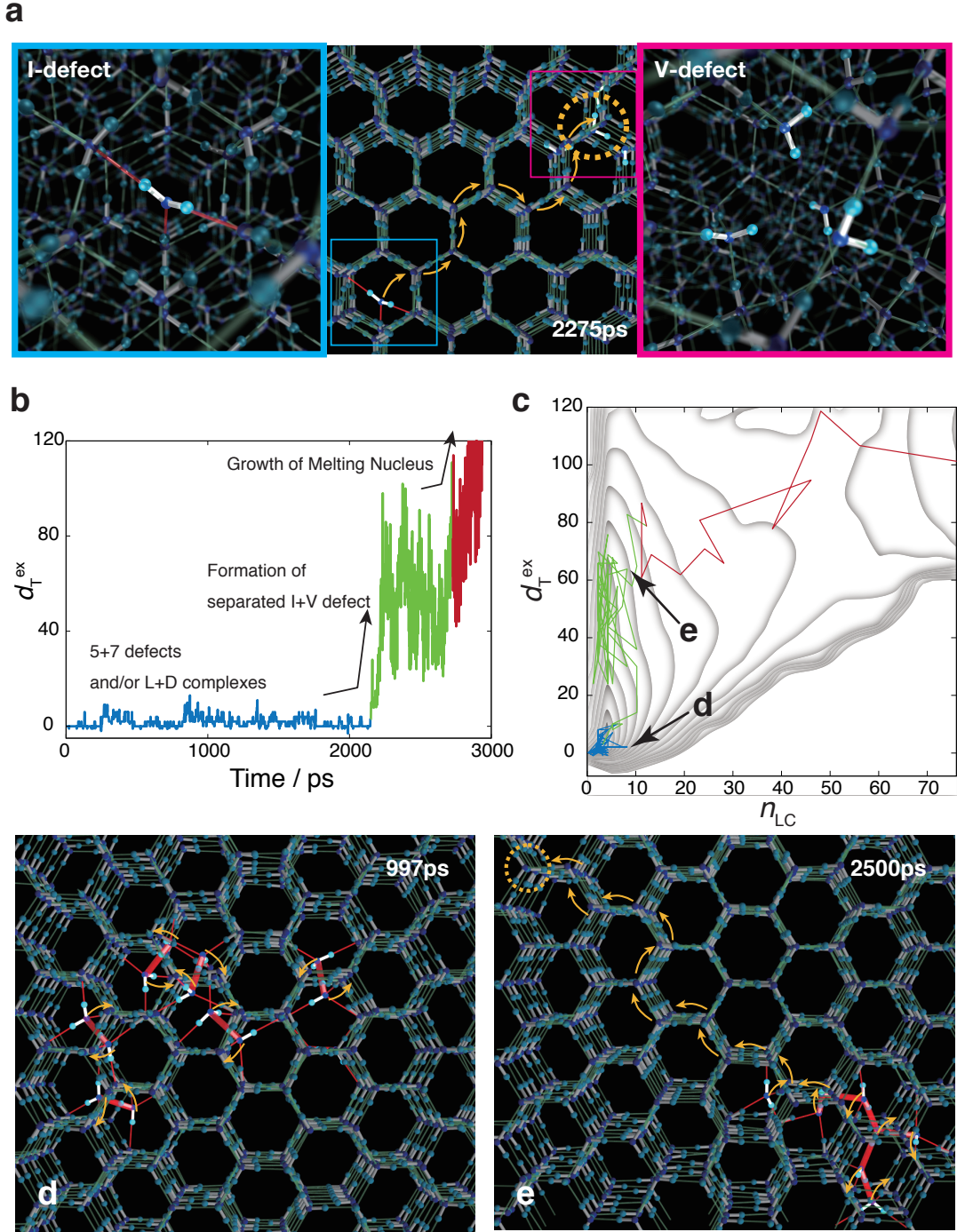


Figure 2 Snapshots of HB structures and the time evolution of excess edit distance, d_T^{ex} . (a) Typical example of a separated I-V defect pair in ice (middle panel). The I defect and the four molecules surrounding the V defect (located at the centre of the dotted circle) are indicated by bright colors and shown magnified in the left and right panels, respectively. The ‘editing’ path to recover a crystalline ice structure (see main text) is indicated by yellow arrows. For this example, with only a pair of I-V defects, d_T^{ex} is 58. (b) Evolution of d_T^{ex} , calculated every 1 ps, for the trajectory of Fig.1. The quiescent period containing 5+7 defects and/or L+D complexes, the formation of the separated I-V defect and the subsequent induction period, and the final period where the melting cluster reaches the size of the critical nucleus and grows rapidly, are colored blue, green and red, respectively. These structurally distinct periods are further emphasized in c, which shows d_T^{ex} against the largest cluster size n_{LC} (calculated every

10 ps). Separation of the defect pair is signalled by the rapid increase in d_{T}^{ex} , while n_{LC} remains small throughout the subsequent induction period lasting from $t = 2,150$ to $2,730$ ps (green lines). In both (b) and (c), dark green shading for $t = 2,210$ – $2,730$ ps indicates when separated defects rapidly dislocate in ice. The thick grey contour lines in the background are the free energy contours shown in Fig.3a. d, e, Snapshots of HB structures from the same melting trajectory, taken at 997 ps and 2,500 ps (indicated by black arrows in (c)). Red lines indicate HBs to off-lattice molecules. The yellow arrows indicate the edit paths that recover a crystalline structure through the formation, cutting and directional inversion of HBs (see Appendix Fig.A5). The structure in (d) contains mostly 5+7 defects, and that in (e) a separated I–V defect pair.

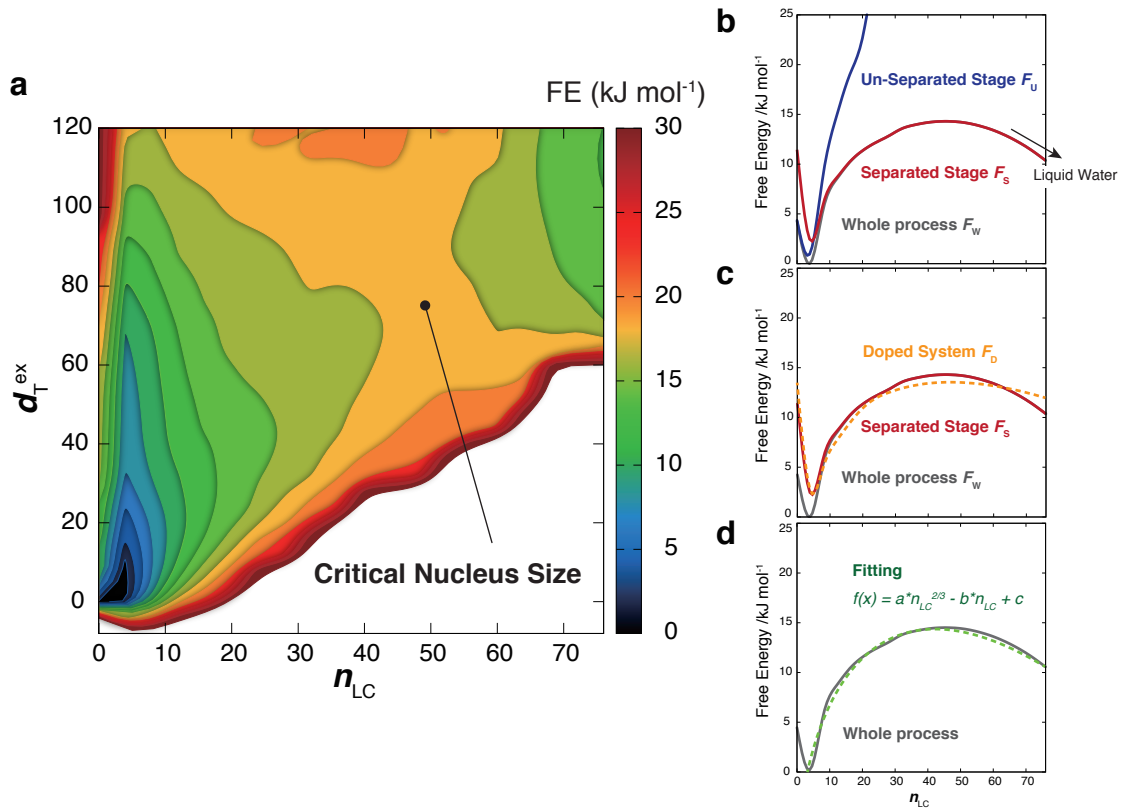


Figure 3 Free energies. (a) The contour map of free energy as a function of largest cluster size (n_{LC}) and excess edit distance (d_{T}^{ex}) at 275 K. The origin ($n_{\text{LC}}, d_{\text{T}}^{\text{ex}}$) = (0,0) represents ice with no defects. The critical nucleus at ($n_{\text{LC}}, d_{\text{T}}^{\text{ex}}$) = (49, 78) is indicated by a black dot. (b) The contour map in a is projected onto n_{LC} to give the ‘whole’ process curve (F_{W} , grey line). The whole process is divided into the unseparated stage (F_{U} : blue line) and the separated stage (F_{S} : red line), see main text. (c) Free energy of the doped system (F_{D} : orange dashed line), which mimics ice containing a permanent I defect (see main text). Note that F_{D} is almost identical to F_{S} (the free energy of the separated stage). (d) The empirical function $f(n) = a \cdot n^{2/3} - b \cdot n + c$ (green dashed line) of classical nucleation theory with the parameters $a=5.93$, $b=1.13$ and $c=-9.71$, optimized by approximating F_{W} . The classical free energy curve overlaps well with F_{W} .

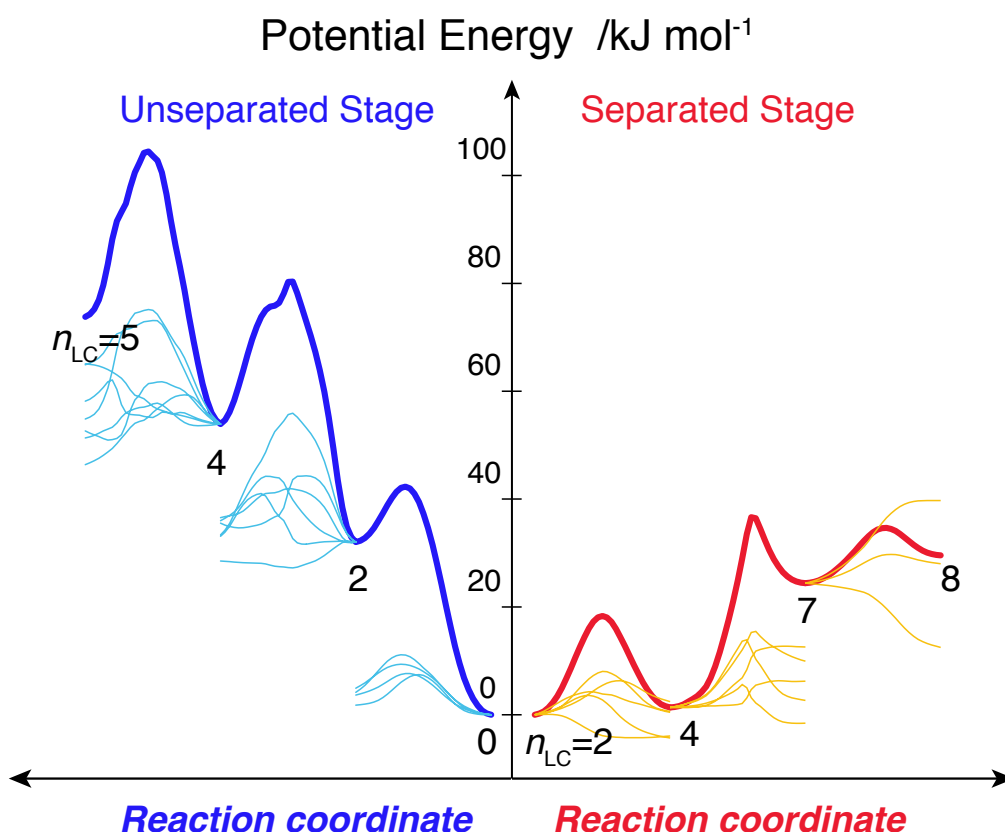


Figure 4 Potential energy surfaces during defect growth. Illustration of how potential energy changes as defect sizes n_{LC} increase, in the unseparated (left: blue lines) and separated (right: red lines) stage. Total potential energy changes are plotted in bold lines, and potential energy changes of individual molecules directly involved in defect growth are plotted with thin blue (left) or thin yellow (right) lines. Potential energies are along reaction coordinates. The nucleus size n_{LC} is indicated by a number at each minimum of an inherent structure. The starting points of individual molecular potential energy changes are shifted to the corresponding minima of the total energy surfaces. For example, five water molecules are directly involved in the defect growth from $n_{LC} = 2$ to 4 in the separated stage (five yellow curves at the bottom of the right side figure); the curve corresponding to the separated I defect exhibits a monotonic decay.

2.4 Appendix

2.4.1 Supplementary figures

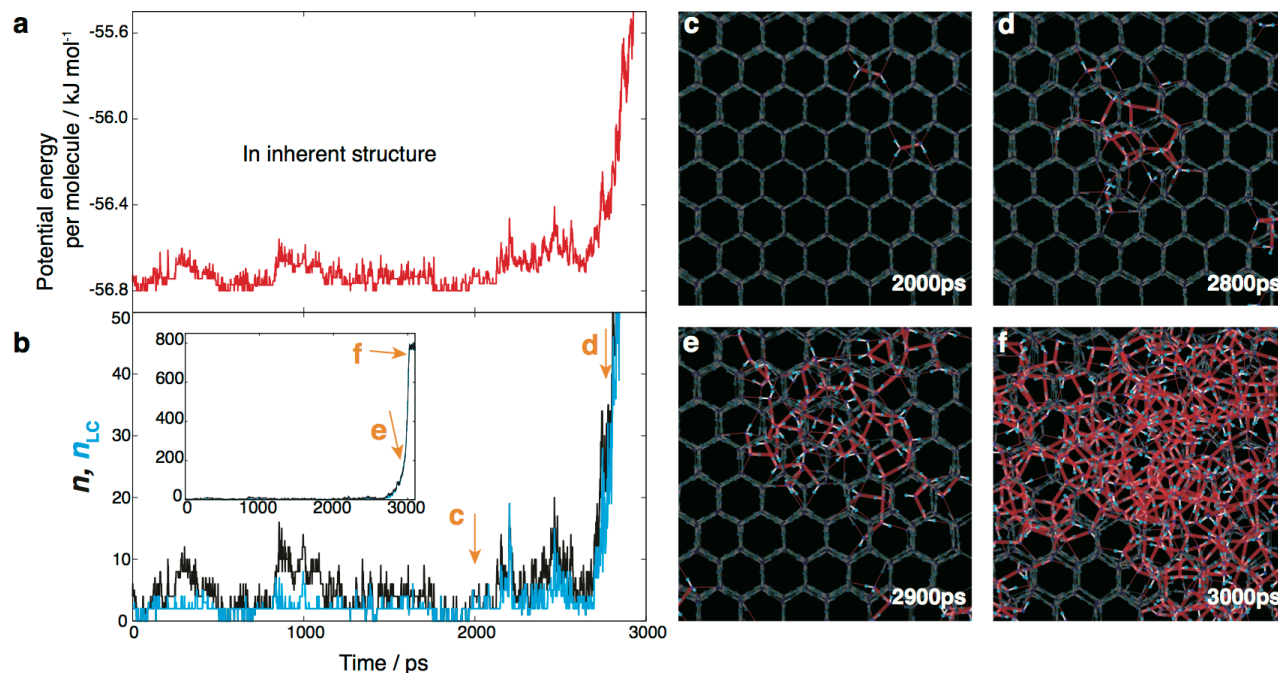


Figure A0 Total potential energy, number of the off-lattice molecules, and HB structures in a melting trajectory. (a) and (b) are same as Fig.1(a) and (b) in the manuscript. (a) Total potential energy of the inherent structures of a typical melting trajectory. At 0 ps, temperature is jumped from 270 K to 275K. (b) Total number of the off-lattice water molecules, n , (black line) and the largest cluster size, n_{LC} , (blue line) for the same trajectory of Fig.1a. (c)-(f) are the snapshots of the HB network structures at 2000, 2800, 2900, and 3000 ps, indicated in (b), respectively. Oxygen and hydrogen atoms of water molecules are represented by small blue and skyblue spheres, respectively. HBs on the off-lattice molecules are represented by red lines. The potential energy fluctuates for a long time (the quiescent period) associated with intermittent creation and the annihilation of small size melting clusters until 2150 ps. Then, the relatively large energy fluctuation appears with the creation of a “separated” I-V pair. The rapid growth of the melting cluster (nucleus), thus the rapid energy rise, starts at 2730 ps.

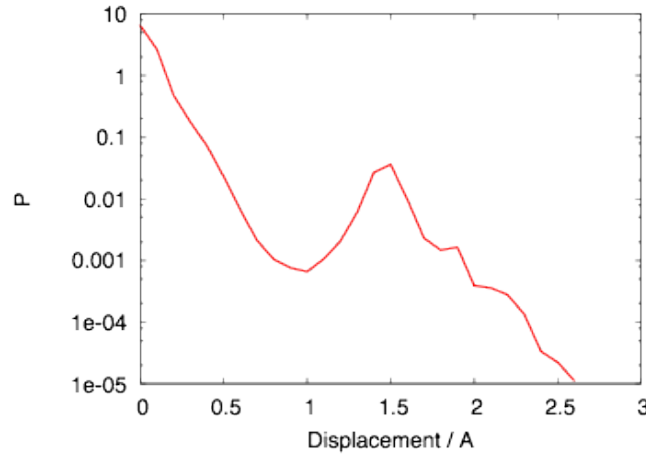


Figure A1 Histogram of water molecular displacements from the nearest lattice points in a melting trajectory (in inherent structures). It clearly shows that the threshold for the off-lattice molecules should be chosen at 1Å.

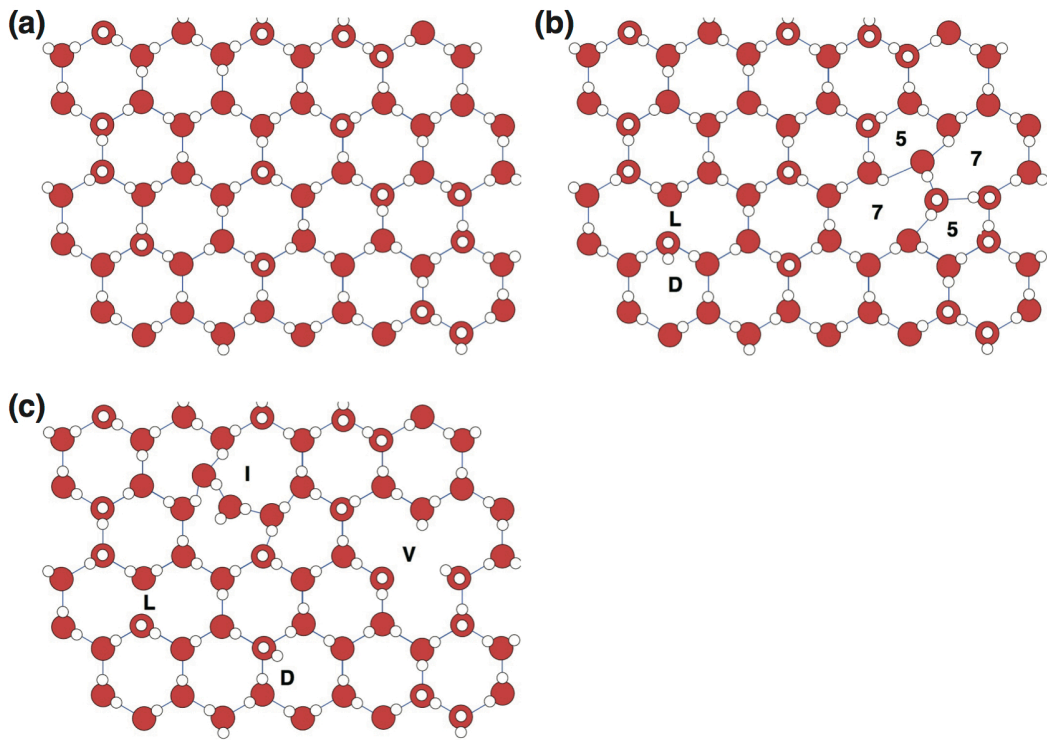


Figure A2 Configurations of the typical defects. (a) Ice Ih with no defect. Each water molecule has four HBs with its four neighbors with obeying Bernal-Fowler ice rule. (b) A 5+7 defect and a L+D complex, directly formed from ice structure with thermal agitation. (c) A separated I+V pair and a separated L+D pair, neither of which can be created directly from ice structure, as discussed in the main text. Other water molecules are in ice structure.

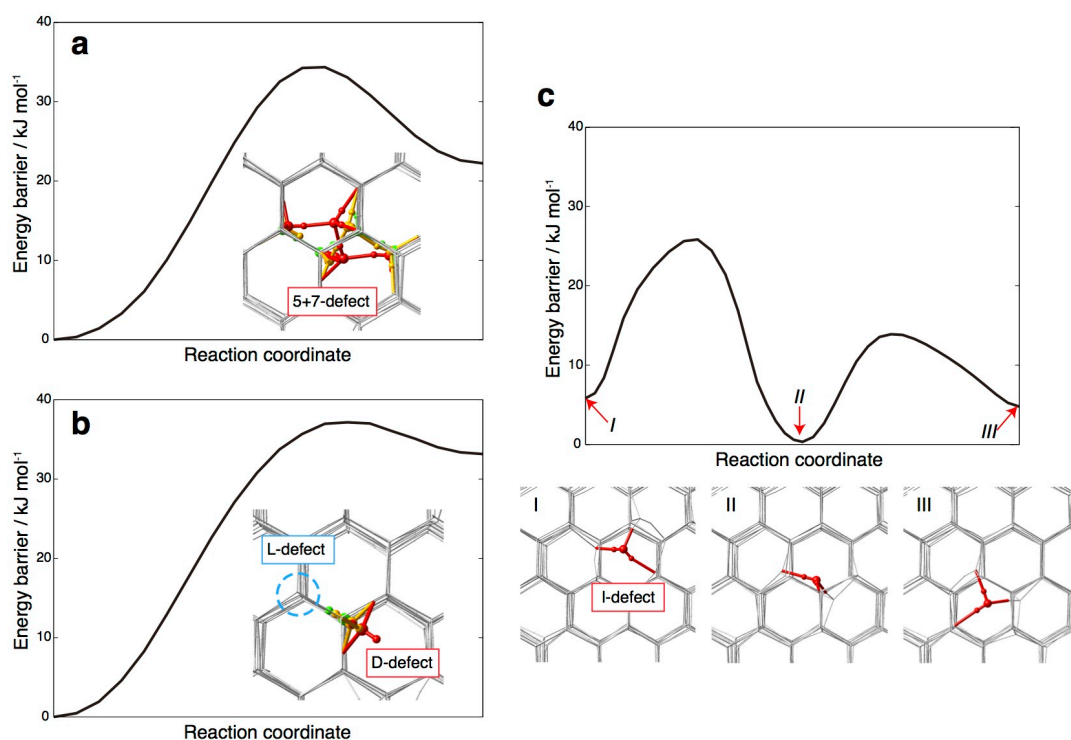


Figure A3 Energy barriers and structure changes along reaction coordinates of HB transformations. (a) Total potential energy change and HB transformation along a reaction coordinate to create a 5+7 defect. Water molecules involved in the HB transformation are colored; green, yellow and red indicates the initial, the intermediate and the final structure of the process, respectively. (b) Total potential energy change and HB transformation along a reaction coordinated to create a L+D complex; colors are used in the same manner as Fig.A3a. Detailed analyses for the reaction coordinates are made by the method of Grishina et al.^{12@} (c) Total potential energy change along reaction coordinates of the sequential dislocation of an I-defect in ice. Configurations of the I-defect in three successive energy minima are displayed in red color. Note that the low energy barriers of backward reactions in (a) and (b) indicate the facile annihilation of these defects. Note also that the barrier of I-defect dislocation in (c) is lower than to create a 5+7 defect (a) or an L+D complex (b). The barrier to create a separated I-V pair is much higher than these energy barriers.

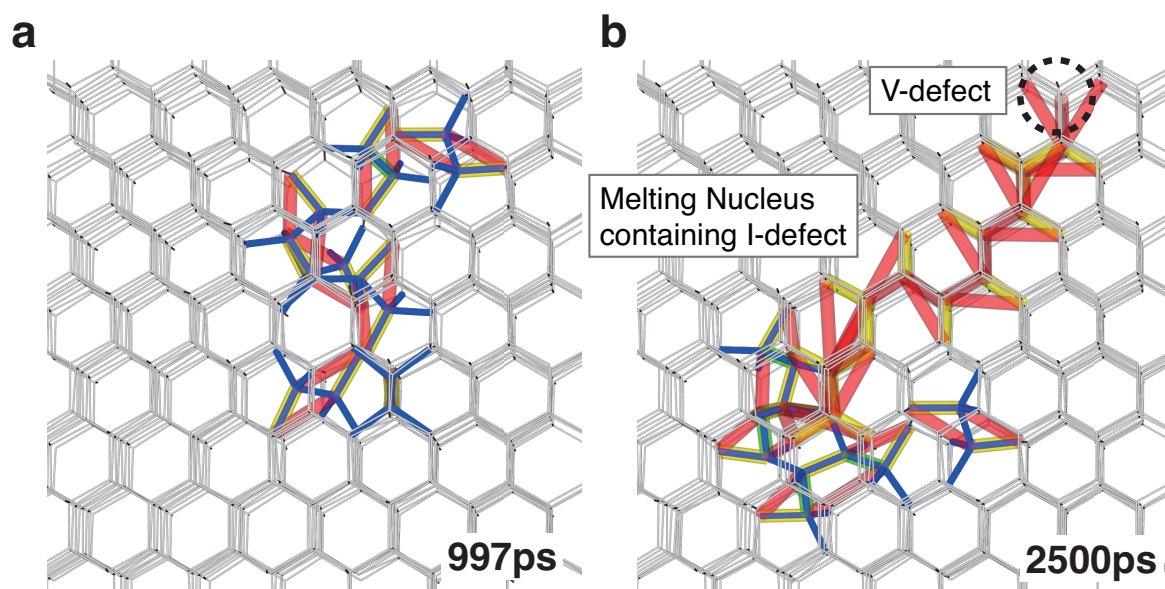


Figure A4 (a) and (b) are the detailed edit paths for the same structures of Fig.3(d) and (e), respectively. Blue lines are HBs on the off-lattice molecules, while gray thin lines are HBs connecting the water molecules on the lattice points of ice. Edit path to recover the crystalline ice consists of forming HBs (red), cutting HBs (yellow) and inverting the HB direction (green).

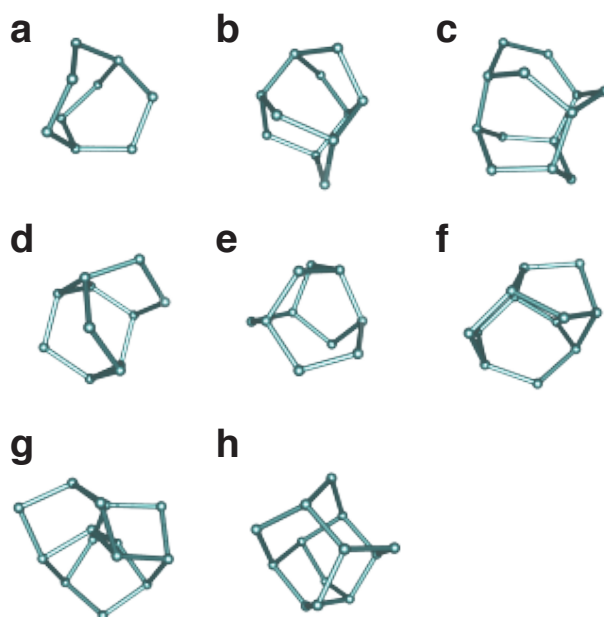


Figure A5 8 typical interfacial fragments lined up from (a) to (h) in the order of their population. First three fragments appear adjacent to a 5+7 defect. Shape and properties of these fragments can be seen at the online fragment database: <http://vitrite.chem.okayama-u.ac.jp> Their ID are #4, #98, #101, #65, #10, #63, #233, #182 in the DB.

(Fragment of the melting nucleus) In the early stage of melting nuclear growth, “interfacial fragments”, consisting of 5-, 6-, and 7-membered rings, appear. Fragments are three-dimensional polyhedral building blocks representing the local connectivity of HB network, in which vertices and edges correspond to water molecules and HBs, respectively. When the melting proceeds, “liquid fragments,” containing distorted rings such as 4- or 8-membered rings, start to appear at the central

domain of the melting nucleus for $n_{LC} > 20$. The “interfacial fragments” then become a wetting boundary of one layer between ice and liquid domain. The “interfacial” molecules, initially highly hindered, become less strained. Eight types of “interfacial” fragments constitute 70% of the interfacial layer, while “liquid fragments” have more varieties.

According to the fragment description of ice melting, one can see there are 4 stages; (1) 0-2150 ps, a considerably long quiescence period, mainly 5+7 defects or L+D complexes intermittently appear and disappear all over the system. Defective structures appearing in this period mainly consist of interfacial fragments (Fig.A5) . (2) 2150-2730 ps. 5+7 defects and L+D complexes are accumulated in a locus and then the system fails to revert back to a crystalline structure, resulting the formation of a I+V defect pair for $t = 2150-2210$ ps,. I-defect undergoes rapid sequential steps of dislocation in ice for $t = 2210-2730$ ps . (3) 2730-3050 ps, when the melting nucleus undergoes the rapid growth. The interfacial fragments exist only in a one-layer boundary between the melting nucleus and the remaining ice structure. (4) after 3050 ps, when the system completes the melting. The interfacial fragments rarely appear in this bulk liquid stage.

2.4.2 Supplementary methods.

2.4.2.1 Edit distance

Recent progress in computer simulation has provided us with new insight into the anomalous properties of liquid water and ice.^{33,34} These discoveries were supported by the development of new analytical methods for computer simulation studies. Most of the anomalous properties originate in the collective aspects of water, which are introduced by the complex HB network. Methods that treat many degrees of freedom are therefore important to explore the properties of water.

How can one measure the difference between structures with many degrees of freedom? In the case of protein folding, for example, the structural difference is often measured by using the root mean square deviation (RMSD) between the atomic positions in two structures.

A similar procedure is applicable when comparing undistorted and thermally fluctuating ice structures. If the latter structure is not quite broken, the structural difference can be measured by summing the displacements between the water positions of the same label in two structures. Counting the dislocated atoms from the lattice point is also a simple way to evaluate how much the given structure is distorted from the crystal structure. The amount of dislocated molecules can be estimated experimentally. Many simulation studies on crystal melting also employ this method as a crude approximation.³⁵

Liquid water and ice have quite a distinct network of hydrogen bonding.³⁶ Ice and liquid water structures can therefore be regarded as the directed graphs. The topology of the graph is expressed by an $N \times N$ adjacency matrix, where N is the number of water molecules in the structure.

The element $a(i, j)$ of the adjacency matrix is 1 when a water molecule i donates an HB to molecule j , and 0 otherwise. Structure change in water can also be measured by the number of HB rearrangements. Hamming distance between two HB networks, say A and B, is defined as the number of different matrix elements in the adjacency matrices of these networks³⁷:

$$D_H(A, B) = \sum_i^N \sum_j^N |a^A(i, j) - a^B(i, j)| \quad (1)$$

where $a^A(i, j)$ is an adjacency matrix element of structure A.

One might think that it is possible to measure the disordered nature of ice at the point of melting by utilizing RMSD, the number of dislocated molecules, or Hamming distance. However, molecular displacements and HB rearrangements in ice do not always introduce disorder. Two typical cases are shown in the following examples.

(I) In ice near the melting point, a water molecule often moves away from its original lattice point by thermal fluctuation and exchange its position with a neighboring water molecule. Repeated exchanges allow the water molecule to diffuse in solid ice. However, the resultant structure is still ice. Thus, exchange of molecules on the lattice points does not introduce disorder in an ice. Mean square displacement cannot be used to figure out the equivalence by molecular exchanges.

(II) In ideal ice, each water molecule donates two HBs and accepts two HBs, the so-called “ice rules”.²⁷ The structure of ice affords many equivalent states satisfying the ice rules. The directed graph of ice contains a certain amount of cyclic paths along which all the edges (i.e. hydrogen bonds) are in the same direction, i.e. homodromic cycles.³⁸ If the original graph obeys the ice rule, inversion of edge directions along a homodromic cycle still conserves the ice rule and the energy difference is very small. Such an inversion of a homodromic cycle, i.e. collective rotation along a cyclic path, can be observed frequently in superheated ice in a molecular dynamics simulation, but it does not introduce disorder in ice. Hamming distance cannot be used to figure out the equivalence by cyclic HB rearrangements along a homodromic cycle.

When I measure the disorder of an ice structure, I should pay attention to essential structural changes but disregard the interconversion between equivalent ice structures. In the present paper, I propose two kinds of edit distances as measures of essential structure change.

2.4.2.1.1 Concept

In information theory, the edit distance is a metric for measuring the amount of difference between two sequences. The edit distance between two strings is defined as the minimal number of edit operations – such as addition or deletion of letters – to transform one string to another. The edit distance is used to measure variation between DNA.

When I apply the concept of edit distance to a water structure, I define the translation of a water molecule or an HB rearrangement, i.e. connection and disconnection, as an edit operation. I call a sequence of edit operations “edit sequence” and a sequence of edit operations in the actual molecular configuration “edit path”. The edit distance between two structures of ice is defined as the minimum number of edit operations to transform one structure to another. Two kinds of edit distances are defined.

2.4.2.1.2 Geometrical edit distance

Geometrical edit distance (GED) is the minimal displacement required to move all water molecules in a given structure to the lattice points of ice. Suppose two ice structures obtained by molecular dynamics simulation. Structures A and B are an initial perfect ice structure and a thermally fluctuated ice structure near the melting point, respectively. In B, most water molecules are on the ice lattice points, while some dislocate, melt locally, or form defective structure such as a 5+7 defect.²³⁻²⁵ I define D as the sum of the square deviations between water molecules of the same label in structures A and B:

$$D(A,B) = \sum_i^N (\mathbf{r}_i^A - \mathbf{r}_i^B)^2 \quad (2)$$

where \mathbf{r}_i^A is the center-of-mass position of the i^{th} water molecule in structure A and N is the number of water molecules in the system. Even when structure B is not molten, water molecules displace from their initial lattice points by exchanging their positions. Therefore, D gradually increases with time.

To eliminate the unessential displacement by point exchanges, I exchange molecular labels in structure A in order to minimize D . If all the water molecules in structure B are on the lattice points, minimal D becomes zero. GED is defined as the minimal D obtained by molecular label permutation. The shortest edit paths to move molecules from a given structure to the lattice points of ice can also be obtained simultaneously. Note that multiple edit paths may give the same GED.

I also define geometrical edit steps (GES) as a discretized expression of GED. GES is the number of times of crossing over the Voronoi bisectors of structure A when water molecules in structure B move along the shortest edit path to the lattice points.

Minimization of D is a combinatorial problem. In practice, D is minimized by exchanging the molecular labels by using a simulated annealing and tempering method. Suppose a water molecule moves away from its lattice point and locates interstitially at a distant place. This kind of pairwise defect is called the separated I-V pair. Though the total structure still looks like almost perfect ice at first glance, both GED and GES become large. In this case, the editing path indicates the shortest route to translate the chain of water molecules from interstitial to vacancy defects (Fig.A6).

If the number of molecules in two structures is different, excess molecules should be removed from the larger system and a smaller number of molecules of both are used to calculate D . Removal of an excess molecule is also regarded as an edit operation with a constant weight w , i.e.,

$$D(A,B) = \sum_i^N (\mathbf{r}_i^A - \mathbf{r}_i^B)^2 + nw \quad (3)$$

where N is the number of molecules in the smaller system and n is the number of excess molecules in the larger system. In calculation of D , the combination of molecular labels i in structure A and i' in structure B is optimized so as to minimize D . The choice of molecules to be removed should be optimized in order to minimize D . Such a situation occurs when extrinsic molecules are injected in the ice lattice or when ice suffers radiation damage.

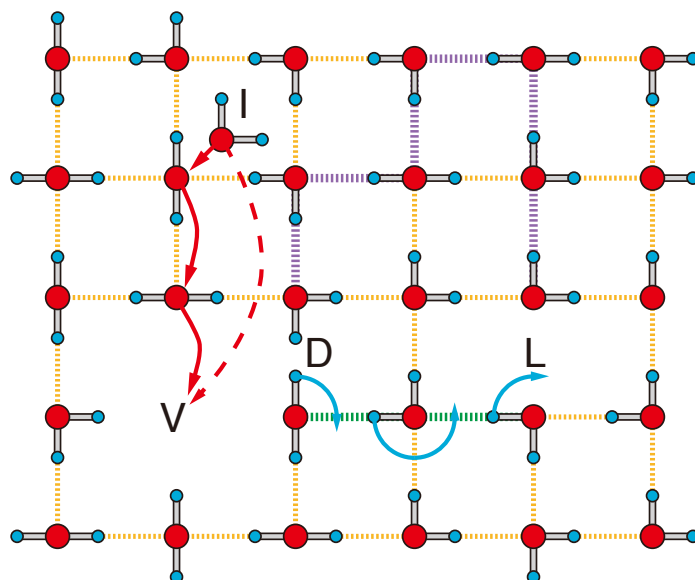


Figure A6 A schematic illustration of the two kinds of edit distances, GED and TED, and the corresponding defect pairs in ice. Dotted lines indicate the HBs. Red solid arrows indicate the geometrical edit path to dissolve the interstitial (I)-vacancy (V) defect pair. A direct leap of an interstitial defect to the vacancy (red dashed arrow) offers a larger D value. Sky-blue arrows indicate the topological edit path to dissolve the D and L defect pair, i.e. the Bjerrum pair. Green lines indicate the HBs to be inverted in the edit path, while purple lines indicate a detour to dissolve the Bjerrum pair.

2.4.2.1.3 Topological edit distance

Topological edit distance (TED) is the minimal number of HB rearrangements (bond creation and removal) required for a given HB network to recover the topology of an ice HB network but which obeys the ice rule. Suppose, there are two ice structures A and B. The former is an initial perfect structure and the latter is a thermally fluctuating structure. Molecular labels are supposed to be optimized by GED calculation in advance.

When structure B is made from structure A by inverting one HB, for example, the Hamming distance between the two structures becomes 2. A perfect crystal structure A obeying the ice rule contains a certain amount of homodromic cycles. One can change the network topology of structure A without violating the ice rule by inverting the cycles, and DH also changes as a consequence.³⁹ TED is defined as the minimal $D_H(A,B)$ obtained by inverting the homodromic cycles in structure A. The shortest edit paths to convert the disordered network structure into a structure of ice obeying the ice rule can also be obtained simultaneously. Note that multiple edit paths may give the same TED.

When A and B are perfect ice structures with a different proton order, TED becomes zero. Minimization of D_H is a combinatorial problem. In practice, D_H is minimized by inverting the bond directions along the homodromic cycles by using a simulated annealing and tempering method.

When one tries to find TED between two structures, one must calculate GED in advance in order to optimize the molecular labeling between them. Otherwise, a geometrically impossible HB rearrangement path might be a candidate for the shortest edit path. Let us exemplify the case when a water molecule moves away from a lattice point and locates interstitially at a distant place. If I try to simply restore the HB network of ice, the shortest edit path would be to remove all the HBs of the interstitial molecule and to create four HBs from that molecule to the molecules around the vacancy. It is, however, impossible because the molecule must leap through the ice lattice.

If the number of molecules in two structures is different, unconnected nodes should be inserted into the smaller system, i.e. empty rows and columns should be added to the adjacency matrix of the smaller system, in order to calculate D_H . TED is finite if GED is finite, since translation of a water molecule in ice always accompanies HB rearrangements. However, the opposite is not true. Therefore TED is a more sensitive index to detect the disorder in ice.

2.4.2.1.4 Pair defect and the edit distance

These edit distances elicit certain kinds of pair defects hidden in disordered ice. In order to dissolve a separated L-D defect pair and relocate the excess protonic defect to the right place, it is necessary to reorient water molecules and invert a chain of HBs by the Bjerrum mechanism.³⁰ TED gives the shortest edit path to dissolve the separated L-D pair. On the contrary, if the shortest edit path in TED includes a chain of HB inversion, there must be D- and L-defects at the ends of the chain. Thus TED elicits the separated L-D pair in a distorted ice.

In order to dissolve a separated I-V pair, a train of water molecules needs to be translated to relocate them to the right places. GED gives the shortest edit path to dissolve the separated I-V pair. On the contrary, if the shortest edit path in GED includes a train of translational motion of water molecules, there must be I- and V-defects at the ends of the train. Thus GED elicits the separated I-V pair in disordered ice.

2.4.2.2 Free energy landscape

I can estimate the topography of the free energy landscape directly from the simulation result with the use of Markov network and the order parameters. Here I introduce the procedure briefly.

2.4.2.2.1 Markov network

Consider there are finite kinds of “states” for a given system and there are the networks of “transition paths” among them with a certain probabilities. At time t , the system takes one of the states with the probability $p(t,i)$. After a small time interval Δt , the system changes its state to j with a transition probability $p(j|i)$. Note that there is also a probability of not changing the state $p(i|i)$ and thus the transition probability satisfies the equation: $\sum_i p(i|j)=1$. Then the probability for the

system to take a state i at time $t+\Delta t$ is given by

$$p(t+\Delta t,i)=\sum_j p(i|j)p(t,j) \quad (4)$$

At the stationary state, the following equation should hold:

$$p(i)=\sum_j p(i|j)p(j), \quad (5)$$

where $p(i)$ is the stationary distribution. When $p(i|j)$ is given, the probabilities $p(i)$ at the stationary state can be calculated by iterating the equation (4).

2.4.2.2.2 Coarse-grained landscape

I can classify structures into finite number of states (i.e. bins), specified with order parameters; for example, a state i is specified with the largest cluster size of defective molecules (n_{LC}), and the excess edit distance (d_T^{ex}), $i = (n_{LC}, d_T^{ex})$. Numbers of structures generated by MD simulation are projected on a state i . Then determine the transition probability $p(j|i)$ between the states i and j by counting the number of trajectories which moves from state i to state j during a time interval Δt . By assuming the local equilibria of whole states, the probabilities $p(i)$ at the stationary state can be calculated by iterating the equation (4).

When the melting nucleus size becomes very large, for example, nuclear growth rate becomes too rapid and I cannot estimate the stationary probabilities for such states precisely any longer. In order to avoid the calculation of transition probabilities between dubious states, I will put a

boundary about the nuclear size and beyond the size is considered to be a single “hidden” state, whose stationary probability is unknowable. Suppose there are only three states, for example, and transition between state 1 and 2 or 2 and 3 are allowed. Here I regard the state 3 as the hidden state, as shown in Fig.A7.

At the stationary distribution, I can assume the detailed balance between the states, that is, $p(3|2)p(t,2) = p(2|3)p(t,3)$. Then the following equation should hold

$$\begin{aligned} p(t+\Delta t, 2) &= p(2|1)p(t,1) + p(2|2)p(t,2) + p(2|3)p(t,3) \\ &= p(2|1)p(t,1) + \{p(2|2) + p(3|2)\}p(t,2) \end{aligned} \quad (6)$$

Thus, the probability for the system to be found in state 2 at the next step, $p(t+\Delta t, 2)$, can be calculated without referencing the transition probability from hidden state 3 to state 2, $p(2|3)$, nor the probability for the system to be found in state 3, $p(t, 3)$. It must be noted that the detailed balance breaks down when the nucleus size is much larger than its critical size. In such case, the gradient of the free energy landscape is so large that the probability of going back to smaller nucleus size becomes too small even with an exhaustive sampling. It is therefore difficult to extend the free energy landscape calculation with the present method to the region of much larger nucleus size.

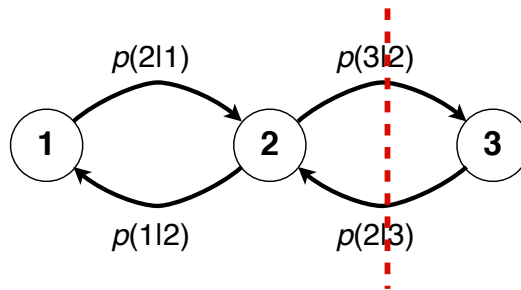


Figure A7 An example system with three states, in which state 3 is the hidden state.

2.4.2.2.3 Effective free energy and the entropy

The effective potential energy of a state i , say $U(i)$, is introduced to represent the coarse-grained potential energy landscape against the order parameter, n_{LC} . $U(i)$ is defined as the average potential energy of the molecular configurations belonging to the largest nucleus. In the stationary state, the probability of choosing a state i , $p(i)$, is considered to be proportional to the Boltzmann distribution ($\exp\{-\beta U(i)\}$) and the degeneracy ($N(i)$) as:

$$p(i) = BN(i) \exp(-\beta U(i)) \quad (7)$$

where B is a proportionality coefficient. From this formula, I can estimate $S(i) = -k_B \ln N(i)$, i.e., the effective entropy of the state i . The effective free energy should be written naturally as:

$$F(i) = -k_B T \ln(p(i)/B) = -k_B T \ln N(i) + U(i) = U(i) - TS(i)$$

It must be noted that $N(i)$ contains all kind of degeneracy, including the number of structures belonging to a state i . I should therefore be careful to deal with the absolute values of $S(i)$ and $F(i)$.

References

- [1] Schmeisser, M., Iglev, H., Laubereau, A., *Chem. Phys. Lett.*, **442**, 171-175 (2007)
- [2] Cahn, R. W., *Nature*, **323**, 668-669 (1986)
- [3] Schmeisser, M., Iglev, H., Laubereau, A., *J. Phys. Chem. B*, **111**, 11271-11275 (2007)
- [4] van der Spoel, D., Maia, F., Caleman, C., *Phys. Chem. Chem. Phys.*, **10**, 6344 (2008)
- [5] Lindemann, F. A., *Phys. Z.*, **11**, 609-612 (1910)
- [6] Born, M. J., *Chem. Phys.*, **7**, 591-601 (1939)
- [7] Jin, Z. H. et al., *Phys. Rev. Lett.*, **87**, 055703 (2001)
- [8] Forsblom, M., Grimvall, G., *Phys. Rev. B*, **72**, 054107 (2005)
- [9] Forsblom, M., Grimvall, G., *Nat. Mater.*, **4**, 388-390 (2005)
- [10] Mizushima, S., *J. Phys. Soc. Jpn.*, **15**, 70-77 (1960)
- [11] Kuhlmann, D., *Phys. Rev.*, **140**, 1599-1610 (1965)
- [12] Grishina, N., Buch, V., *J. Chem. Phys.*, **120**, 5217-5225 (2004)
- [13] Donadio, D., Raiteri, P., Parrinello, M., *J. Phys. Chem. B*, **109**, 5421-5424 (2005)
- [14] Jorgensen, W. L. et al., *J. Chem. Phys.*, **79**, 926-935 (1983)
- [15] Jorgensen, W. L., Madura, J. D., *Mol. Phys.*, **56**, 1381-1392 (1985)
- [16] Vega, C. et al., *J. Chem. Phys.*, **122**, 114507 (2005)
- [17] Jacobson, L. C., Hujo, W., Molinero, V., *J. Phys. Chem. B*, **113**, 10298-10307 (2009)
- [18] Nosé, S., *J. Phys. Condens. Matter*, **2**, SA115-SA119 (1990)
- [19] Henkelman, G., Uberuaga, B. P., Jonsson, H., *J. Chem. Phys.*, **113**, 9901-9904 (2000)
- [20] Frenkel, D., Smit, B., *Understanding Molecular Simulation: From Algorithms to Applications*, Academic (2002)
- [21] McBride, C., *Mol. Phys.*, **103**, 1-5 (2005)
- [22] Tanaka, H., Mohanty, J., *J. Am. Chem. Soc.*, **124**, 8085-8089 (2002)
- [23] Bernal, J. D., Fowler, R. H., *J. Chem. Phys.*, **1**, 515-548 (1933)
- [24] Petrenko, V. F., Whitworth, R. W., *Physics of Ice*, Oxford Univ. Press (1999)
- [25] Abraham, F. F., *Homogeneous Nucleation Theory: The Pretransition Theory of Vapor Condensation*, Academic (1974)
- [26] Wales, D. J., *Energy Landscapes: Applications to Clusters, Biomolecules and Glasses*, Cambridge Univ. Press (2004)
- [27] Matsumoto, M., Baba, A., Ohmine, I., *J. Chem. Phys.*, **127**, 134504 (2007)
- [28] Iglev, H. et al., *Nature*, **439**, 183-186 (2006)

- [29] Matsumoto, M., Saito, S., Ohmine, I., *Nature*, **416**, 409–413 (2002)
- [30] Jacobson, L. C., Hujo, W., Molinero, V., *J. Am. Chem. Soc.*, **132**, 11806–11811 (2010)
- [31] Walsh, M. R. et al., *Science*, **326**, 1095–1098 (2009)
- [32] Moore, E. B., Molinero, V., *Nature*, **479**, 506–508 (2011)
- [33] Mishima, O., Stanley, H. E., *Nature*, **396**, 329–335 (1998)
- [34] Koga, K. et al., *Physica a-Stat. Mech. Its Applications*, **314**, 462–469 (2002)
- [35] Forsblom, M., Grimvall, G., *Phys. Rev. B*, **72**, 053107 (2005)
- [36] Matsumoto, M., *J. Chem. Phys.*, **126**, 054503 (2007)
- [37] Matsumoto, M., Ohmine, I., *J. Chem. Phys.*, **104**, 2705–2712 (1996)
- [38] Bjerrum, N., *Math. -Fys. Medd.*, **27**, 35 (1951)
- [39] Saenger, W., *Nature*, **279**, 343–344 (1979)

Chapter 3

Local structure of methanol-water binary solutions studied by soft X-ray absorption spectroscopy and molecular dynamics simulation

M. Nagasaka, K. Mochizuki, V. Leloup and N. Kosugi

J. Phys. Chem. B (submitted)

3.1 Introduction

It is known that methanol-water binary solutions show smaller entropy than expected in an ideal solution of randomly mixed molecules¹ and show a nonlinear profile in viscosity as changing the mixing ratio.² These characteristics have been discussed by using clathrate-like structure models of methanol molecules with surrounding water molecules and with hydrophobic interactions between methyl groups.¹ However, a consistent picture of the microscopic structure of methanol-water binary solutions is not yet established.

The oxygen atom in a water molecule has two hydrogen-donating ('donor') sites and two hydrogen-accepting ('acceptor') sites, and liquid water forms tetrahedrally-coordinated three-dimensional (3D) hydrogen bond (HB) networks.³ On the other hand, a methanol molecule has one donor and one or two acceptor sites due to the replacement of one donor site by a hydrophobic methyl group, and liquid methanol forms one- and two-dimensional (1D/2D) HB networks, such as chain and ring structures.⁴⁻⁹ In the neutron diffraction experiments of methanol-water binary solutions,¹⁰⁻¹² it is found that 3D HB networks of methanol-water mixtures are formed by hydrophilic and hydrophobic interactions between water and methanol molecules. Dixit et al. measured neutron diffraction at $X=0.7$ in the methanol-water binary solutions $(\text{CH}_3\text{OH})_X(\text{H}_2\text{O})_{1-X}$,¹¹ and revealed that the distance between methyl groups of methanol molecules becomes closer by adding water molecules.

The interaction between methanol and water molecules in the binary solution was studied by nuclear magnetic resonance,¹³ mass spectrometry,¹⁴ Rayleigh scattering,¹⁵ and dielectric relaxation methods.¹⁶ Takamuku et al. measured the number of water molecules per 6 methanol molecules as a function of the methanol molar fraction by the mass spectrometry of methanol-water liquid micro-jets,¹⁴ and found three different dependences with the borders at $X=0.7$ and $X=0.3$. They proposed that the chain structures of methanol clusters are dominant at $X>0.7$, the tetrahedral-like water clusters gradually evolve at $0.7>X>0.3$, and the water cluster is a main species at $0.3>X>0.0$.

The interaction in the binary solution was also investigated by vibrational spectroscopies: infrared spectroscopy¹⁷⁻²⁴ and Raman spectroscopy.^{22,25-28} Dixit et al. found nonlinear profiles of the C-O stretching vibration in the Raman spectroscopy when decreasing the methanol molar fraction,²⁷ in which the behavior of the energy shifts changes at $X=0.70$ and $X=0.25$. They proposed different local structures: in the region $X>0.7$, water molecules connect the terminal of the methanol chain and the chain structure is preserved; in the region

$0.7 > X > 0.25$, the methanol chain is broken by adding water molecules; in the region $0.25 > X > 0.05$, the hydration structure of methanol molecules is formed.

The structure of liquid methanol and methanol-water binary solutions has been investigated theoretically by using molecular dynamics (MD)²⁹⁻⁴³ and Monte Carlo simulations.⁴⁴⁻⁴⁹ In the methanol-rich condition, the 1D/2D HB network structure of methanol clusters is not strongly influenced by water molecules.³² As the mixing ratio of water increases, the HB networks of both water and methanol molecules grow to be mixed with each other.^{38-40,46,48} When the mixing ratio of water is high, the 3D hydration shell is formed around methanol molecules.^{41,43,45}

Although methanol-water binary solutions have been studied experimentally and theoretically as described above, microscopic structures of methanol-water mixtures, such as nearest neighbor interactions, have not yet been known in detail. Soft X-ray absorption spectroscopy (XAS) is an element-selective method to investigate local structures of liquid and aqueous solutions. The structure of liquid water was extensively studied by the O K-edge XAS.⁵⁰⁻⁵² Because the X-ray absorption process occurs within a few femto seconds, XAS gives us the information of the averaged HB structures, which rearrangements occur within the time scale. The hydration structure of cations in aqueous salt solutions was also investigated by the O K-edge XAS.⁵³⁻⁵⁴ Wilson et al. studied the O and C K-edge XAS of liquid methanol in the total electron yield of liquid micro-jet.⁵⁵ Tamenori et al. measured the O and C K-edge XAS of free methanol clusters.⁵⁶ Guo et al. investigated liquid methanol and methanol-water binary solutions at $X=0.5$ by using the O K-edge XAS and X-ray emission spectroscopy.⁵⁷⁻⁵⁸ Guo et al. suggested that liquid methanol shows chains and rings of 6-8 methanol molecules, and proposed that the number of pure methanol chains decreases and the number of mixed methanol-water networks increases when adding water molecules. However, the O K-edge XAS shows contributions of oxygen atoms in both methanol and water molecules, and is difficult to analyze the local structure of methanol-water mixtures. It is necessary to measure the C K-edge XAS to analyze the local structure of the methyl group of methanol molecules in the binary solution.

In the present work, the local structure of methanol-water binary solutions at different concentrations by the O and C K-edge XAS is investigated. The XAS measurement is based on a transmission mode by using a recently developed liquid cell that enables to optimize the absorbance by changing the thickness of liquid layer.⁵⁹ The pre-edge feature in the O K-edge

XAS is found to show almost linear dependence of the concentration. On the other hand, in the C K-edge XAS, the spectral intensity in a characteristic energy region is found to change its behavior at $X=0.7$ and $X=0.3$. With the help of the MD simulation, it has been revealed different local structures of methanol-water mixtures at the three concentration regions.

3.2 Experiments

The experiments were performed at an in-vacuum soft X-ray undulator beam line BL3U at UVSOR-II.⁶⁰ Details of the liquid cell were described previously.^{59,61} The liquid cell consists of four regions, which are separated by 100 nm-thick Si₃N₄ membranes (NTT AT Co., Ltd.). Soft X-rays under vacuum (region I) pass through the buffer region filled with helium gas (region II) and the liquid thin layer (region III), and finally reach a photodiode detector filled with helium gas (region IV). The regions II and IV are connected and can be mixed with other gas molecules for the precise gas-liquid energy shift measurement and photon energy calibration. A liquid sample (region III) is sandwiched between two Si₃N₄ membranes with pressed Teflon spacers, and can be substituted by other samples in combination with a tubing pump system.

The thickness of liquid layer should be optimized in order to transmit soft X-rays with an appropriate absorbance.⁶² In the present liquid cell, the thickness of the liquid layer can be controlled from 2000 nm to 20 nm by increasing the helium pressure in the regions II and IV. The thickness is set to 300 nm in the present O K-edge XAS. In the C K-edge XAS, on the other hand, the thickness of liquid methanol is set to 550 nm, and the thickness is set larger in more dilute methanol aqueous solutions. The energy resolutions of incident soft X-rays at the O and C K-edges are set to 0.40 eV and 0.19 eV, respectively. The XAS spectra are based on the Beer-Lambert law, $\ln(I_0/I)$, where I_0 and I are the detection current through the cell without and with samples, respectively. The liquid flow is stopped during the XAS measurement because the sample liquid has no radiation damage from the long (say, more than one hour) exposure of soft X-rays in the present photon flux. The photon energy in the O K-edge is calibrated by the O 1s - π^* peak (530.80 eV)⁶³ for free O₂ molecules and that in the C K-edge is calibrated by the first peak (287.96 eV)⁵⁶ of free methanol molecules, which are mixed with helium gas in the regions II and IV.

3.3 Results and discussions

3.3.1 Oxygen K-edge XAS

Fig.1 shows O K-edge XAS spectra for methanol-water binary solutions of different concentrations at 25 °C. The absorbance in the O K-edge XAS spectra was normalized by the sample thickness and the concentration of the binary solution considering the soft X-ray absorption coefficients of water and methanol in the O K-edge.⁶² After this normalization, a constant background was subtracted and the resultant the absorption coefficients is shown in Fig. 1. In the previous work,⁵⁷ a small peak was observed around 532 eV in methanol-water mixtures; on the other hand, it is not observed for any concentration in the present measurement. It should not be regarded as an intrinsic peak. The pre-edge peak of liquid water (534.7 eV) corresponds to the O 1s transition to an unoccupied $4a_1^*$ orbital of a water molecule (533.9 eV), which is mainly distributed at the oxygen atom in water molecule and is blue-shifted and broadened by the HB interaction.⁵¹ On the other hand, the pre-edge feature of liquid methanol (534.9 eV) is embedded in the main peak but is similarly blue-shifted from the gas-phase peak (534.0 eV).⁵⁵

Fig.1 shows that the intensity of the pre-edge region around 535.2 eV decreases as the molar fraction of methanol (X) decreases in the binary solution $(\text{CH}_3\text{OH})_x(\text{H}_2\text{O})_{1-x}$. It is known that the pre-edge peak in liquid water reflects the HB interaction, and the intensity of methanol is different from that of water. The pre-edge region shows isosbestic points at 534.8 eV and 535.9 eV. In such a case, one does not have to focus on the peak itself and one can select one of the energy regions showing rather large change without detailed spectral analysis. This is 535 eV or 540 eV. Since the signal-to-noise ratio in 540 eV is worse than in 535 eV, the pre-edge region (535 eV) was used to analyze the change in the HB interaction at different concentrations.

Fig.2 shows the intensity dependence of the energy region between 534.9 and 535.8 eV with a molar fraction step of 0.05. The intensity decreases almost linearly as the molar fraction of methanol decreases. Note that the similar linear dependence is found for the energy region of 540 eV. There may be some information about different local structures on the oxygen atom behind small deviations from the linear dependence, but it is consistent with the result of vibrational spectroscopy,^{22,24} which explains that the ratio of HB interaction of methanol-methanol to that of methanol-water is linearly dependent on the molar fraction of methanol.

Observation of the isosbestic points suggests that the binary solution has two major HB components, though there are four different HB interactions: $O_m^*HO_m$, $O_m^*HO_w$, $O_w^*HO_m$, and $O_w^*HO_w$, where O_m and O_w denote oxygen atoms of methanol and water, respectively, and $*$ denotes the atom with an O 1s hole. There are two possibilities. One is a negligible HB interaction between water and methanol, $O_m^*HO_w$ and $O_w^*HO_m$, indicating that water aggregates are segregated from methanol ones in solution. The other is almost the same HB interaction in $O_w^*HO_m$ as in $O_w^*HO_w$ and that in $O_m^*HO_w$ as in $O_m^*HO_m$. The pre-edge peak of liquid water is sensitive to the HB interaction of liquid water and is dependent on the temperature.⁵² The concentration dependence in the pre-edge region of the methanol-water binary solution is smaller than the temperature dependence of liquid water. Therefore, nearly the same HB interaction between methanol and water could be possible.

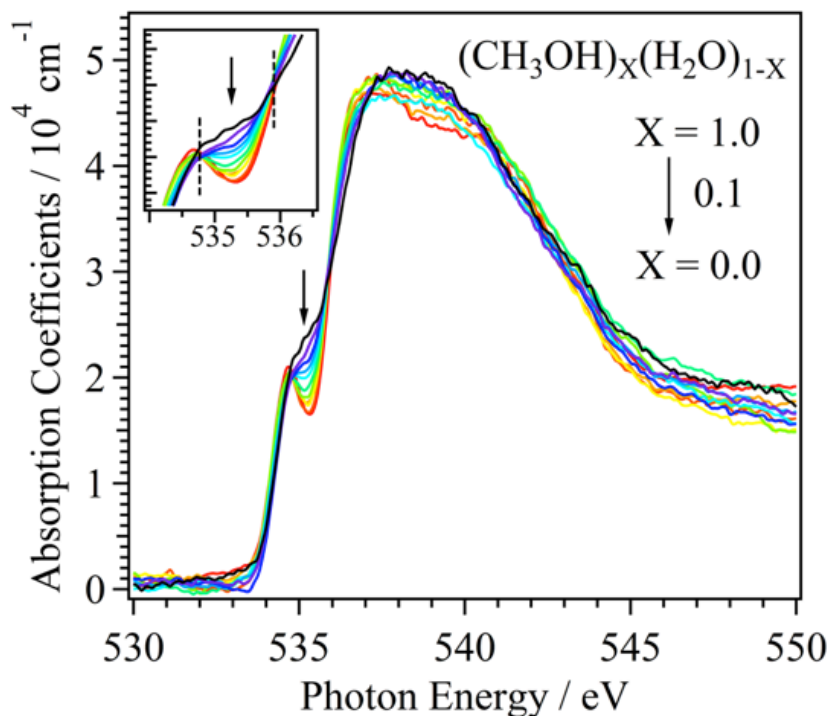


Figure 1 O K-edge XAS spectra of methanol-water binary solutions of different concentrations at 25 °C. The mixing ratio of methanol in the solution decreases with molar fraction steps of 0.1 along indicated arrows. The inset shows isosbestic points (dashed lines) in the pre-edge region.

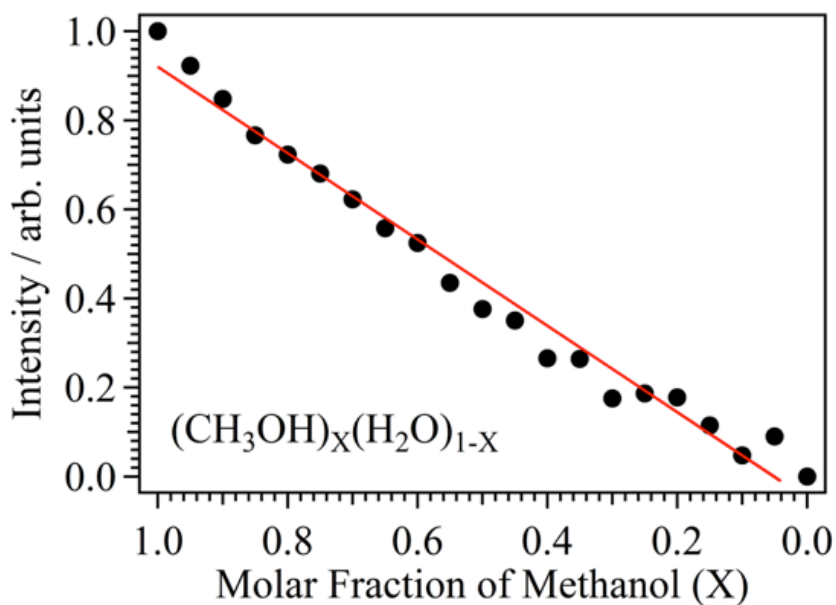


Figure 2 Intensity dependence of the energy region between 534.9 and 535.8 eV in the O K-edge XAS spectra as a function of methanol molar fraction (X) in the methanol-water binary solutions $(\text{CH}_3\text{OH})_X(\text{H}_2\text{O})_{1-X}$. The intensities of methanol ($X=1.0$) and water ($X=0.0$) are normalized to one and zero, respectively.

3.3.2 Carbon K-edge XAS

Fig.3 shows C K-edge XAS spectra of molecular (gas) and liquid methanol at 25 °C by using the same sample cell. The XAS spectrum of methanol gas was measured by mixing methanol vapor into helium buffer gas, and is in agreement with published spectra.^{55-56,64-65} As shown in Fig. 3a, the first peak (287.96 eV) in a gas-phase spectrum arises from a transition of the C 1s electron to the lowest unoccupied orbital (8a') of C-O anti-bonding and O-H bonding characters. The second peak around 289.44 eV arises from a transition of the C 1s electron to the second lowest unoccupied orbital (9a') of pseudo $\text{CH}_3\text{-}\pi^*$ character with a very small OH component. The broad peak around 293 eV arises from a transition of the C 1s electron to the highest unoccupied orbital (11a') within a minimal basis picture, which is of both C-O anti-bonding and O-H anti-bonding character.

Fig.3b shows the present C K-edge XAS spectrum of liquid methanol with a simple structure of three main contributions around 288.4, 289.55, and 293 eV. The published C K-edge XAS spectrum of free methanol clusters⁵⁶ is almost the same as the present liquid spectrum. On the other hand, the published spectrum of liquid micro-jet methanol⁵⁵ is different from the present one and is rather similar to the gas-phase spectrum. It could be

difficult to completely remove the contribution from the molecular methanol in the liquid micro-jet experiment.

In the liquid spectrum, the first peak (8a'-related), which is an excited state with both CH₃ and OH components, shows a 0.44 eV blue shift from that in the gas spectrum. On the other hand, the second peak (9a'-related), which is an excited state with a large CH₃ component, shows a 0.11 eV blue shift from that in the gas spectrum. Considering atomic components in the corresponding molecular orbitals of a methanol molecule, the second excited state has mainly a hydrophobic interaction but the first excited state has not only a hydrophobic interaction but also an HB interaction. The unoccupied orbital level could be destabilized by both the hydrophobic and hydrophilic interactions in liquid, similarly to the case of the blue-shifted O 1s pre-edge peak in liquid water as observed in Fig. 1.

Fig.4 shows C K-edge XAS spectra of methanol-water binary solutions of different concentrations at 25 °C. The C K-edge XAS spectrum is more appropriate than the O K-edge XAS as regards the analysis of the intermolecular interaction of methanol, because the carbon atom is contained only in methanol. The absorbance in the C K-edge XAS spectra was normalized by the sample thickness and the concentration of the binary solution considering the soft X-ray absorption coefficients of methanol in the C K-edge.⁶² After this normalization, the absorbance of water in the C K-edge was subtracted by considering the sample thickness, the concentration of the binary solution, and the soft X-ray absorption coefficient of water in the C K-edge.⁶² Fig.4 shows resultant absorption coefficients. The first peak does not change its energy position so much at different concentrations. This is reasonable if the HB interaction of methanol with water is not so different from that with methanol.

On the other hand, the second peak, which corresponds to the excited state with a large CH₃ component, increases the blue shift as the mixing ratio of water increases. It is reasonable, considering that the blue shift arises from the interaction of methyl group in methanol molecule. Liquid methanol forms 1D/2D network structures and the methyl groups are apart from each other due to its hydrophobic interaction. When water molecules join the 1D/2D HB network of methanol, the 3D HB network might be formed. Then, the interaction of methyl groups can be enhanced in binary solutions and causes the blue shift of the second peak in the C K-edge spectra. The general behavior observed in Fig. 4 is consistent with that of the neutron diffraction¹¹, where mixed methanol-water networks are formed and methyl groups become closer to each other in the methanol-water binary solution.

When the second peak (289.55 eV) related to the methyl group is blue-shifted by increasing the mixing ratio of water. The pre-edge region shows a quasi-isosbestic point around 290 eV. Similarly in the case of the O K-edge region, it is possible to select one of the energy regions showing rather large change without detailed spectral analysis. This is 289 eV or 290.5 eV. The intensity of the valley structure around 289 eV decreases rather sensitively and the signal-to-noise ratio in 290.5 eV is worse than in 289 eV. Therefore, the pre-edge region (289 eV) was used to investigate the change in the hydrophobic interaction at different concentrations.

Fig.5 shows the intensity in the energy region between 288.7 and 289.2 eV at different molar fractions of methanol (X) in the binary solution $(\text{CH}_3\text{OH})_x(\text{H}_2\text{O})_{1-x}$. The intensity is changed nonlinearly, and show three different behaviors with the borders of $X=0.7$ and $X=0.3$. Note that the similar nonlinear dependence is found for the energy region of 290.5 eV. In the methanol-rich region I ($X>0.7$), the intensity is not so much changed as compared to the intensity of liquid methanol ($X=1.0$). The phase transition-like behavior of the intensity change is found at $X=0.7$. The slow decrease in intensity (indicating blue shift of the second pre-edge peak) continues when increasing the mixing ratio of water in the region II ($0.7>X>0.3$). The decrease in intensity becomes faster in the water-rich region III ($0.30>X>0.05$). These results suggest different local interactions of the methyl group at the different concentration regions.

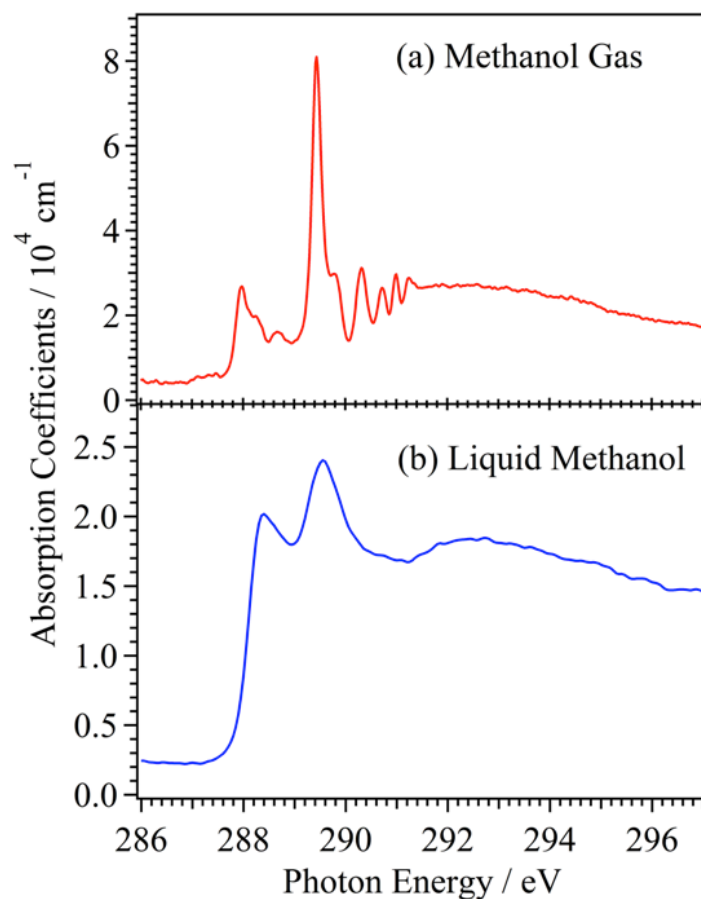


Figure 3 C K-edge XAS spectra of (a) methanol gas and (b) liquid methanol at 25 °C.

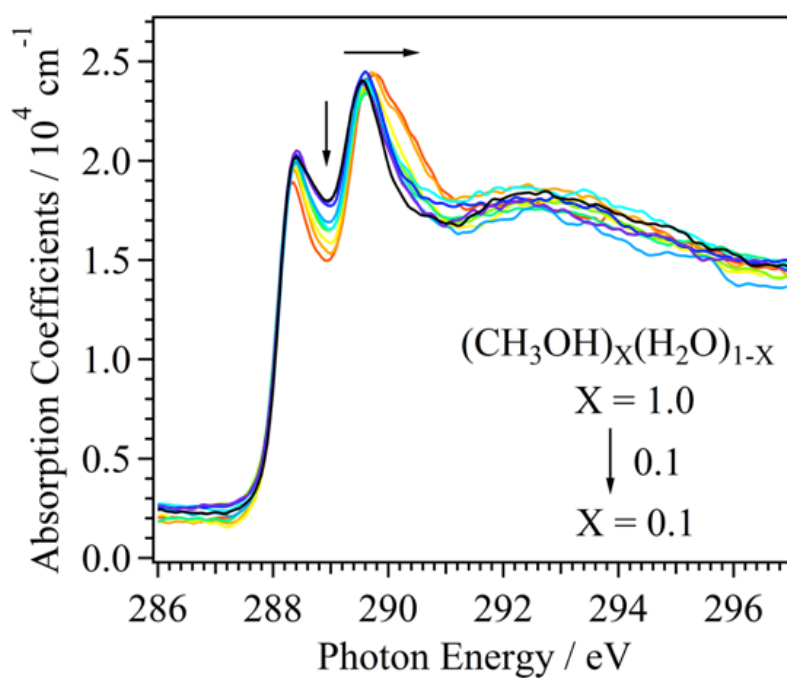


Figure 4 C K-edge XAS spectra of methanol-water binary solutions at different concentrations at 25 °C. The mixing ratio of methanol in the solution decreases with molar fraction steps of 0.1 along indicated arrows.

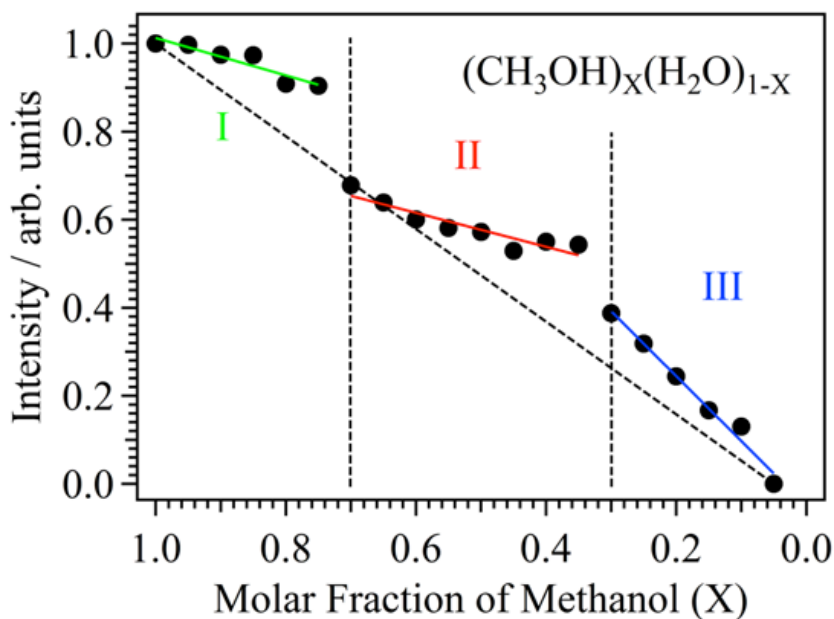


Figure 5 Results of the C K-edge XAS for the intensity dependence of the valley between 288.7 and 289.2 eV as a function of methanol molar fractions (X) in the methanol-water binary solutions $(\text{CH}_3\text{OH})_X(\text{H}_2\text{O})_{1-X}$. The intensities of methanol ($X=1.0$) and the dilute solution at $X=0.05$ are normalized to one and zero, respectively. Three characteristic regions are found with the borders of $X=0.7$ and $X=0.3$.

3.3.3 MD simulation

Without any spectral calculation based on time-consuming density functional theory or ab initio approaches, it would be simply understood that the C K-edge region is sensitive to the hydrophobic interaction around the methyl group of methanol. In order to get such information from the radial distribution function (RDF) of intermolecular interaction in the solutions, I have carried out the MD simulation by using GROMACS 4.5.5.⁶⁶ The potential of methanol molecule is described by OPLSAA,⁶⁷⁻⁶⁸ and that of water molecule is TIP5P.⁶⁹ When compared with earlier models (TIP3P and TIP4P), TIP5P model forms a more "tetrahedral" water structure that better reproduces the experimental radial distribution functions from neutron diffraction. The temperature is controlled by the Nosé-Hoover thermostat method.⁷⁰ The pressure is adjusted by the Parrinello-Rahman method.⁷¹ The simulation was performed at a time step of 1 fs with a periodic boundary condition and the particle-mesh Ewald method.⁷² The unit cell consists of 500 molecules, and the molar fraction of methanol (X) is changed from $X=0.0$ to $X=1.0$. Randomly distributed structures were

optimized by the simulations, which run during 50 ps at 100 K in the NVT condition, 50 ps at 200 K and 1 atm in the NPT condition, and 400 ps at 298.15 K and 1 atm in the NPT condition. The equilibrium structures were obtained by sampling the structures every 1 ps during a simulation time of 2 ns.

First, I calculated RDF of four different HB: O_m-HO_m , O_m-HO_w , O_w-HO_m , and O_w-HO_w . The distances of both the first peak and the first minimum point in RDF are not changed even at different molar fractions. It means that the HB interaction of water is nearly the same as that of methanol as already discussed in the O K-edge XAS.

The number of HBs was counted at different molar fractions, based on the criterion of the distance between HO and O within the first minimum point (2.5 Å) in RDF.⁷³ By increasing the molar fraction of water, the total number of HBs increases linearly. The average number of HBs around water molecules is between 3.2 and 3.8, and that around methanol is between 1.8 and 2.5. This is consistent, considering the water molecule has two donor and two acceptor sites and the methanol molecule has one donor and one or two acceptor sites. The total average number of HBs increases at the higher mixing ratio of water molecules.

Fig.6 shows the RDF from C in the methyl group of methanol to surrounding atoms at different molar fractions of methanol (X) in the binary solution $(CH_3OH)_X(H_2O)_{1-X}$. Fig.6a shows the RDF of C with C and hydrogen H (HC) atoms in the CH_3 group of neighboring methanol molecules. The RDF distances of both C-C and C-HC are slightly reduced by increasing the mixing ratio of water. This result is consistent with the results of neutron diffraction.¹¹ Fig.6b shows the RDF of C with O_m and HO_m in neighboring methanol molecules. The intensities of both HO_m and O_m in the first coordination peak decrease as increasing the mixing ratio of water. On the other hand, as shown in Fig. 6c, the intensities of both HO_w and O_w in the first coordination peak increase as increasing the mixing ratio of water.

Fig.7 shows the coordination number by nearest neighbors HO_m and HO_w to the C atom in the methyl group of methanol at different binary solutions. The coordination is defined within the RDF distance of 3.2 Å, which is the first minimum point of HO_m and HO_w as shown in Fig. 6. The methyl group in liquid methanol ($X=1.0$) is surrounded by HO_m (blue). By increasing the molar fraction of water, the number of HO_m coordination decreases and instead that of HO_w (red) increases. When the methanol molar fraction is below $X=0.7$, the number of HO_w coordination becomes larger than the HO_m coordination. It is reasonable

considering that the molar ratio of methanol and water is 2:1 at $X=0.67$ and the ratio of the H donating site is 1:1. However, the rate of increase in the number of HO_w coordination is larger than the rate of decrease in the number of HO_m coordination. The rate of increase in the number of HO_w coordination is accelerated in the region III ($0.3 > X$). On the other hand, the number of HO_m coordination is nearly zero in the region III. Note that $X=0.7$ and 0.3 are the same borders in the spectral change of the C K-edge XAS as shown in Fig. 5.

Next, I investigate the mesoscopic scale HB network in the binary solution. The average size of methanol clusters in pure liquid methanol ($X=1.0$) is 40 in the present MD simulation. The hydrophobic interaction of the methyl group prevents a large HB network formation, and permit only formation of small methanol clusters with an average size of 40. This size is larger than the previously predicted size, 6-8 molecules.⁸

On the other hand, water molecules like to meet (bond) together to form a large HB network in solution. Fig.8 shows results of the MD simulation for the average HB network size and the average size of water-only clusters embedded in total HB networks at different molar fractions of methanol (X) in the binary solutions $(\text{CH}_3\text{OH})_X(\text{H}_2\text{O})_{1-X}$. A unit cell in the present MD simulation contains totally 500 molecules.

In the methanol-rich region I ($X > 0.7$), the average size of water-only clusters is rather small considering that a water molecule is difficult to meet another water molecule in this region. Water clusters start to grow at $X=0.7$. As the molar fraction of methanol is down to $X=0.3$ in the region II ($0.7 > X > 0.3$), the rate of growing in size of water clusters is accelerated, and finally the average size of water clusters is equal to the total number of water molecules in the region III. At $X=0.7$, all the methanol and water molecules join a large HB network, though a large water-only cluster is not yet formed in the network. The ratio of the total number of H donating (accepting) sites is 1:1 for water and methanol at $X=0.67$; therefore, all the water and methanol molecules can meet (bond) together to form a large HB network at around $X=0.7$.

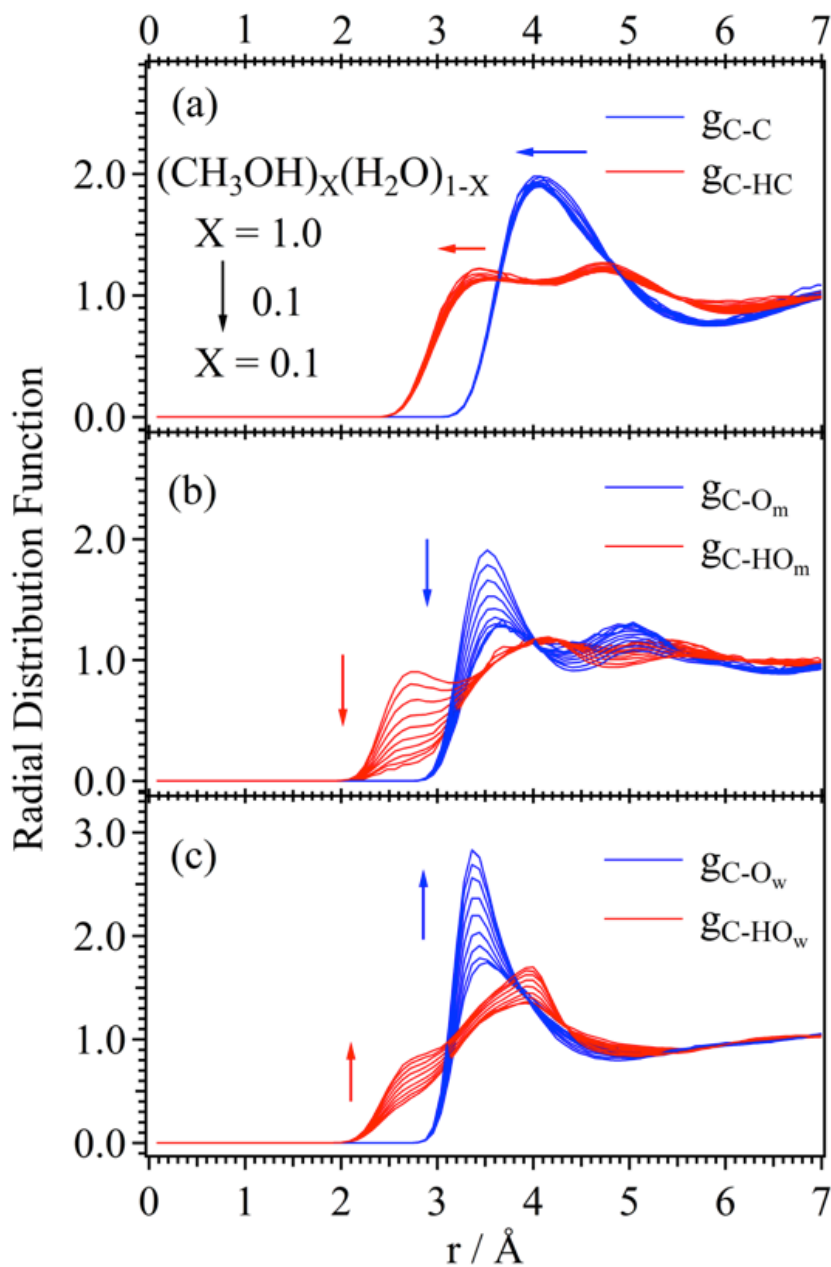


Figure 6 Results of the MD simulation for the RDF of the C atom with surrounding atoms: (a) C and HC in methanol, (b) Om and HOm in methanol, and (c) Ow and HOw in water. The mixing ratio of methanol in the solution decreases from $X=1.0$ to $X=0.1$ with molar fraction steps of 0.1 along indicated arrows. Note that RDF in (c) is changed from $X=0.9$ to $X=0.1$.

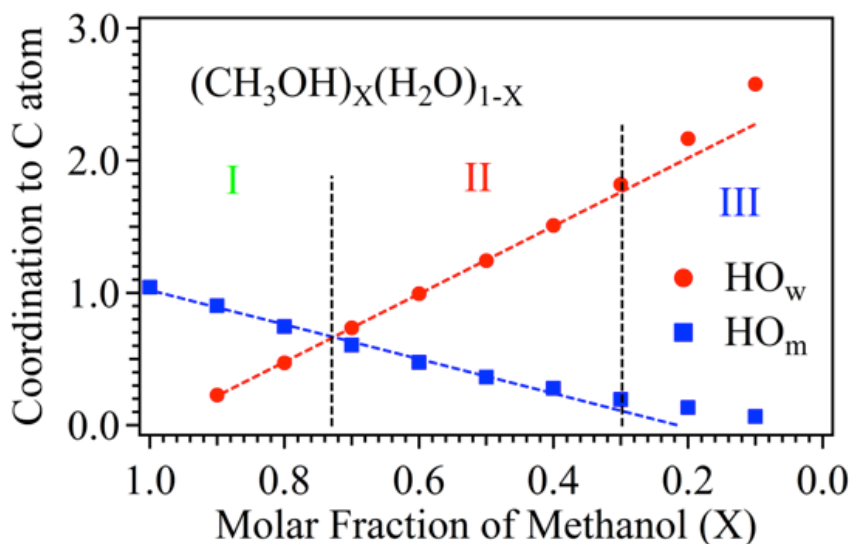


Figure 7 Results of the MD simulation for the coordination numbers of HO_m (blue) and HO_w (red) to the C atoms of the methyl group of methanol at different molar fractions of methanol (X) in the binary solutions $(\text{CH}_3\text{OH})_X(\text{H}_2\text{O})_{1-X}$. Below $X=0.73$ (region II), the coordination number of HO_w becomes larger than that of HO_m . In the region III ($0.3 > X$), the linear dependence of HO_m and HO_w is not valid as shown by dashed lines.

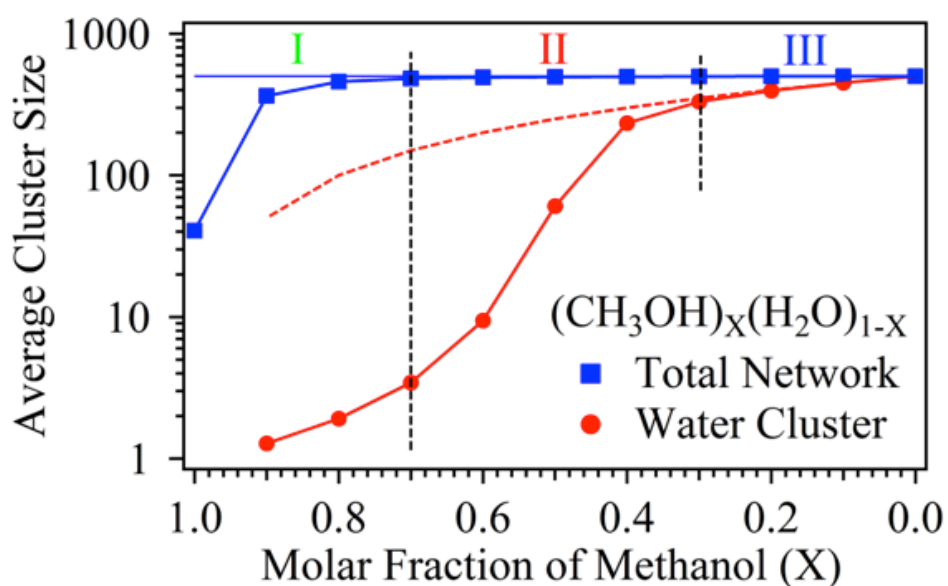


Figure 8 Results of the MD simulation for the average HB network size (blue) and the average size of water-only clusters (red) embedded in total HB networks in a unit cell (500 molecules) of the binary solution $(\text{CH}_3\text{OH})_X(\text{H}_2\text{O})_{1-X}$ as a function of molar fraction of methanol (X). The dashed line (red) is a total number of water molecules in the solution.

3.3.4 Structures of methanol-water mixtures

From the result of the C K-edge XAS shown in Fig. 5, the interaction around the methyl group of methanol molecule shows characteristic changes at the three concentration regions. The MD simulations also show similar three concentration regions from the coordination number around the methyl group shown in Fig. 6 and the average cluster size shown in Fig. 7.

Fig.9a shows a typical structure in the binary solution at $X=0.9$ in the methanol-rich region I. As shown in Fig.8, the average size of water-only clusters is much smaller than the total number of water molecules. Water molecules form the HB network with methanol clusters and stabilize the total energy of the binary solution. However, the interaction around the methyl group of methanol is not so much influenced by water molecules because of a small amount of isolated water molecules. It is consistent with the previous work,^{16, 27} where water molecules are coordinated to the terminal of methanol chains in the methanol-rich region.

Fig.5 shows a phase transition-like intensity change at $X=0.7$ in the C K-edge XAS. Fig. 7 shows that the number of HO_w coordination to the methyl group becomes larger than that of HO_m coordination below $X=0.7$. In addition, Fig.8 shows that the HB network of water clusters grows rapidly below $X=0.7$. Fig.9b shows a typical structure in the binary solution at $X=0.5$ in the region II. Water molecules form a large cluster and have the 3D HB network with methanol molecules, resulting in the increase of the interaction of the methyl group in methanol with water molecules. The phase transition-like behavior at $X=0.7$ in the C K-edge XAS (Fig. 5) indicates that the 3D HB network involving water clusters is dominant over the 1D/2D HB network of methanol in the binary solutions. This result is consistent with the previous MD simulation, in which the 1D chain structure of methanol molecules is changed to 3D mixed clusters by adding water molecules.⁴⁰

Fig.9c shows a typical structure in the binary solution at $X=0.1$ in the water-rich region III. The HB networks between methanol molecules are mostly diminished, and methanol molecules are isolated in the 3D HB network of water. The hydration structures of methanol molecules are dominated by the 3D HB network of water, and the numbers of water coordination to the methyl group increase. As a result, the hydrophobic interaction around the methyl group is enhanced in this region, increasing a blue-shift in the C K-edge XAS. The previous theoretical studies suggested that hydration structures of methanol molecules are formed in this concentration region,^{41, 43, 45} consistent with the present result.

It is known that the thermodynamic parameter such as entropy and viscosity shows an extreme value at the molar fraction of $X=0.30$.¹⁻² Dougan et al. studied neutron diffraction experiments and MD simulations, and suggested that both methanol and water molecules are percolated in this region, and the thermodynamic parameters show extreme values at the molar fraction of $X=0.27$.¹² It means that the structure and abundance of large mixed methanol-water HB networks in the binary solution affect macroscopic thermodynamic properties.

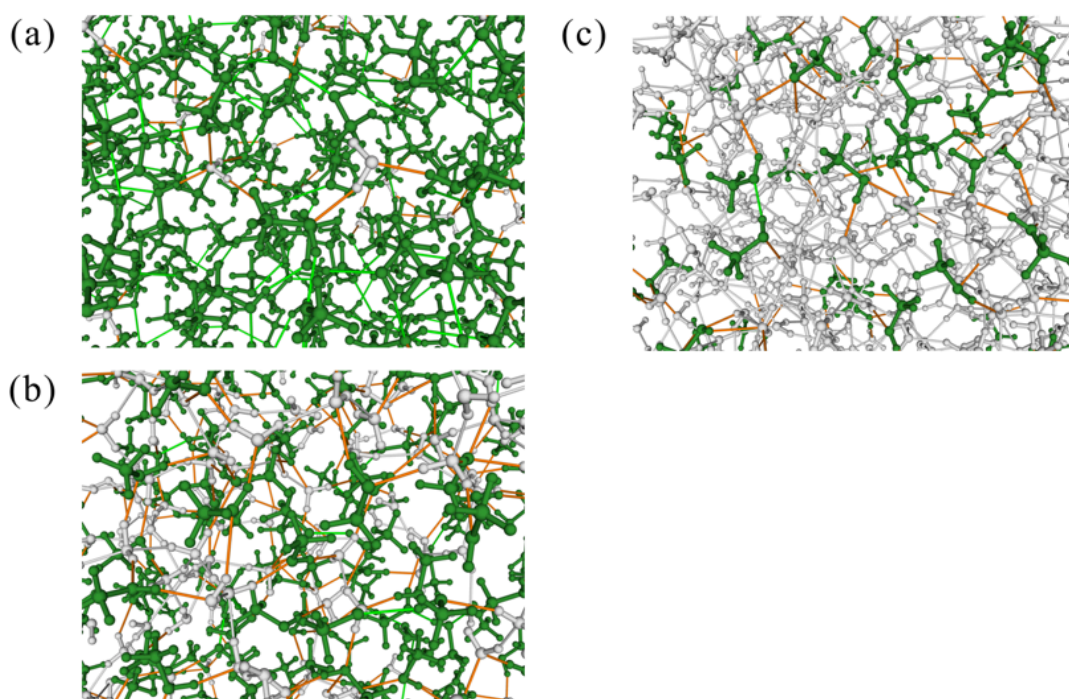


Figure 9 Typical structures of methanol-water binary solutions at different concentrations: (a) $X=0.9$, (b) $X=0.5$, and (c) $X=0.1$. Methanol and water molecules are marked as green and white, respectively. HB between methanol and water is marked as an orange line.

3.4 Conclusions

The local structure of methanol-water binary solutions was studied by the O and C K-edge XAS. The pre-edge peak in the O K-edge XAS reflects the HB interaction of oxygen atoms and shows almost linear concentration dependence in intensity. It indicates that the HB interaction of methanol with surrounding water molecules is nearly the same as in pure liquid methanol and the HB interaction of water with surrounding methanol molecules is nearly the same as in pure liquid water.

The C K-edge XAS enables us to investigate local structures around the methyl group of methanol molecules in the binary solution. The peak around 290 eV in the C K-edge XAS corresponds to a transition of the C 1s electron to the unoccupied orbital around the methyl group, and shows the higher photon energy (blue shift) as the mixing ratio of water increases. It predicts enhancement in the interaction between the hydrophobic methyl groups by large water clusters in mixed methanol-water networks. The intensity change shows a nonlinear profile with three characteristic concentration regions in the binary solution $(\text{CH}_3\text{OH})_x(\text{H}_2\text{O})_{1-x}$.

The three regions are consistently interpreted with the help of the MD simulation. Liquid methanol is known to have the 1D/2D HB network structure.⁸ In the methanol-rich region I ($X > 0.7$), the XAS spectra show only weak concentration dependence in intensity. A small amount of water molecules exists separately around the methanol clusters. Because the structure of methanol clusters is not so much influenced with water molecules, the interaction around the methyl group in methanol is not changed in this region. The phase transition-like decrease in the pre-edge intensity occurs at the molar fraction of $X = 0.7$. The 3D HB network of water start to grow rapidly and the HO_w coordination to the methyl group becomes dominant over the HO_m coordination when the molar fraction is below $X = 0.7$. In the region II ($0.7 > X > 0.3$), methanol molecules form a large HB network with water molecules. As a result, the hydrophobic interaction of the methyl group is enhanced in this region. This behavior is reasonable considering that the molar ratio of methanol and water is 2:1 at $X = 0.67$ and the ratio of the H donating site is 1:1. The thermodynamic parameters such as entropy and viscosity are closer to extreme values as the number of mixed HB networks increases at $X = 0.3$. In the water-rich region III ($0.3 > X > 0.05$), the decrease of intensity in the 289 eV region in the C K-edge XAS is accelerated, indicating methanol molecules are embedded in the 3D HB network of water molecules.

The methyl group in the methanol-water binary solution shows three characteristic local structures: methanol-dominant 1D/2D HB network structure, methanol-water mixed 3D HB network structure, and water-dominant 3D HB network structure. These features are successfully revealed by the pre-edge analysis in the C K-edge XAS, which is sensitive to the hydrophobic interaction of the methyl group, and the MD simulation.

References

- [1] Frank, H. S., Evans, M. W., *J. Chem. Phys.*, **13**, 507-532 (1945)
- [2] Mikhail, S. Z., Kimel, W. R., *J. Chem. Eng. Data*, **6**, 533-537 (1961)
- [3] Ludwig, R., *Angew. Chem.-Int. Edit.*, **40**, 1808-1827 (2001)
- [4] Magini, M., Paschina, G., Piccaluga, G., *J. Chem. Phys.*, **77**, 2051-2056 (1982)
- [5] Narten, A. H., Habenschuss, A., *J. Chem. Phys.*, **80**, 3387-3391 (1984)
- [6] Tanaka, Y., Ohtomo, N., Arakawa, K., *Bull. Chem. Soc. Jpn.*, **57**, 644-647 (1984)
- [7] Tanaka, Y., Ohtomo, N., *Bull. Chem. Soc. Jpn.*, **58**, 270-276 (1985)
- [8] Sarkar, S., Joarder, R. N., *J. Chem. Phys.*, **99**, 2032-2039 (1993)
- [9] Yamaguchi, T., Hidaka, K., Soper, A. K., *Mol. Phys.*, **96**, 1159-1168 (1999)
- [10] Soper, A. K., Finney, J. L., *Phys. Rev. Lett.*, **71**, 4346-4349 (1993)
- [11] Dixit, S. et al., *Nature*, **416**, 829-832 (2002)
- [12] Dougan, L. et al., *J. Chem. Phys.*, **121**, 6456-6462 (2004)
- [13] Corsaro, C. et al., *J. Phys. Chem. B*, **112**, 10449-10454 (2008)
- [14] Takamuku, T. et al., *Z. Naturforsch.*, **55**, 513-525 (2000)
- [15] Micali, N. et al., *Phys. Rev. E*, **54**, 1720-1724 (1996)
- [16] Sato, T., Chiba, A., Nozaki, R., *J. Chem. Phys.*, **112**, 2924-2932 (2000)
- [17] Falk, M., Whalley, E., *J. Chem. Phys.*, **34**, 1554-1568 (1961)
- [18] Passchier, W. F., Klompmaker, E. R., Mandel, M., *Chem. Phys. Lett.*, **4**, 485-488 (1970)
- [19] Bertie, J. E. et al., *Appl. Spectrosc.*, **47**, 1100-1114 (1993)
- [20] Venables, D. S., Schmuttenmaer, C. A., *J. Chem. Phys.*, **113**, 11222-11236 (2000)
- [21] Adachi, D. et al., *Appl. Spectrosc.*, **56**, 357-361 (2002)
- [22] Ma, G., Allen, H. C., *J. Phys. Chem. B*, **107**, 6343-6349 (2003)
- [23] Woods, K. N., *J. Chem. Phys.*, **123**, 134507 (2005)
- [24] Ahmed, M. K., Ali, S., Wojcik, E., *Spectr. Lett.*, **45**, 420-423 (2012)
- [25] Kabisch, G., Pollmer, K., *J. Mol. Struct.*, **81**, 35-50 (1982)
- [26] Schwartz, M., Moradi-Araghi, A., Koehler, W. H., *J. Mol. Struct.*, **81**, 245-252 (1982)
- [27] Dixit, S., Poon, W. C. K., Crain, J., *J. Phys. Condens. Matter*, **12**, L323-L328 (2000)
- [28] Lin, K. et al., *J. Phys. Chem. B*, **114**, 3567-3573 (2010)
- [29] Haughney, M., Ferrario, M., McDonald, I. R., *J. Phys. Chem.*, **91**, 4934-4940 (1987)
- [30] Ferrario, M. et al., *J. Chem. Phys.*, **93**, 5156-5166 (1990)
- [31] Skaf, M. S., Ladanyi, B. M., *J. Phys. Chem.*, **100**, 18258-18268 (1996)
- [32] Laaksonen, A. et al., *J. Phys. Chem. A*, **101**, 5910-5918 (1997)

- [33] Tsuchida, E., Kanada, Y., Tsukada, M., *Chem. Phys. Lett.*, **311**, 236-240 (1999)
- [34] van Erp, T. S., Meijer, E. J., *Chem. Phys. Lett.*, **333**, 290-296 (2001)
- [35] Morrone, J. A., Tuckerman, M. E., *J. Chem. Phys.*, **117**, 4403-4412 (2002)
- [36] Pagliai, M. et al., *J. Chem. Phys.*, **119**, 6655-6662 (2003)
- [37] Handgraaf, J. W. et al., *J. Chem. Phys.*, **121**, 10111-10119 (2004)
- [38] Yu, H. et al., *J. Comput. Chem.*, **27**, 1494-1504 (2006)
- [39] Zhong, Y., Warren, G. L., Patel, S., *J. Comput. Chem.*, **29**, 1142-1152 (2007)
- [40] Bakó, I. et al., *Phys. Chem. Chem. Phys.*, **10**, 5004-5011 (2008)
- [41] Silvestrelli, P. L., *J. Phys. Chem. B*, **113**, 10728-10731 (2009)
- [42] Ishiyama, T., Sokolov, V. V., Morita, A., *J. Chem. Phys.*, **134**, 024509 (2011)
- [43] Moin, S. T. et al., *J. Comput. Chem.*, **32**, 886-892 (2011)
- [44] Jorgensen, W. L., *J. Am. Chem. Soc.*, **102**, 543-549 (1980)
- [45] Okazaki, S., Touhara, H., Nakanishi, K., *J. Chem. Phys.*, **81**, 890-894 (1984)
- [46] Adamovic, I., Gordon, M. S., *J. Phys. Chem. A*, **110**, 10267-10273 (2006)
- [47] Valdéz-González, M. et al., *J. Chem. Phys.*, **127**, 224507 (2007)
- [48] da Silva, J. A. B. et al., *Phys. Chem. Chem. Phys.*, **13**, 6452-6461 (2011)
- [49] da Silva, J. A. B. et al., *Phys. Chem. Chem. Phys.*, **13**, 593-603 (2011)
- [50] Smith, J. D. et al., *Science*, **306**, 851-853 (2004)
- [51] Wernet, P. et al., *Science*, **304**, 995-999 (2004)
- [52] Huang, C. et al., *Proc. Natl. Acad. Sci. U. S. A.*, **106**, 15214-15218 (2009)
- [53] Naslund, L.A. et al., *J. Phys. Chem. A*, **109**, 5995-6002 (2005)
- [54] Waluyo, I. et al., *J. Chem. Phys.*, **134**, 224507 (2011)
- [55] Wilson, K. R. et al., *J. Phys. Chem. B*, **109**, 10194-10203 (2005)
- [56] Tamenori, Y. et al., *J. Chem. Phys.*, **128**, 124321 (2008)
- [57] Guo, J.H. et al., *Phys. Rev. Lett.*, **91**, 157401 (2003)
- [58] Kashtanov, S. et al., *Phys. Rev. B*, **71**, 104205 (2005)
- [59] Nagasaka, M. et al., *J. Electron Spectrosc. Relat. Phenom.*, **177**, 130-134 (2010)
- [60] Hatsui, T., Shigemasa, E., Kosugi, N., *AIP Conf. Proc.*, **705**, 921-924 (2004)
- [61] Nagasaka, M. et al., *J. Phys. Chem. C*, **117**, 16343-16348 (2013)
- [62] Chantler, C. T., *J. Phys. Chem. Ref. Data*, **29**, 597-1048 (2000)
- [63] Coreno, M. et al., *Chem. Phys. Lett.*, **306**, 269-274 (1999)
- [64] Hempelmann, A. et al., *J. Phys. B: At. Mol. Opt. Phys.*, **32**, 2677-2689 (1999)
- [65] Prince, K. C. et al., *J. Phys. Chem. A*, **107**, 1955-1963 (2003)

- [66] Hess, B. et al., *J. Chem. Theory Comput.*, **4**, 435-447 (2008)
- [67] Jorgensen, W.L., Tirado-Rives, J., *Proc. Natl. Acad. Sci. U.S.A.*, **102**, 6665-6670 (2005)
- [68] Caleman, C. et al., *J. Chem. Theory Comput.*, **8**, 61-74 (2012)
- [69] Mahoney, M. W., Jorgensen, W. L., *J. Chem. Phys.*, **112**, 8910-8922 (2000)
- [70] Nosé, S., *J. Phys. Condes. Matter*, **2**, SA115-SA119 (1990)
- [71] Parrinello, M., Rahman, A., *J. Appl. Phys.*, **52**, 7182-7190 (1981)
- [72] Darden, T., York, D., Pedersen, L., *J. Chem. Phys.*, **98**, 10089-10092 (1993)
- [73] Matsumoto, M., *J. Chem. Phys.*, **126**, 054503 (2007)

Chapter 4

A conformational factorisation approach for estimating the binding free energies of macromolecules

K. Mochizuki, C.S. Whittieston, S. Somani, H. Kusumaatmaja and D.J. Wales

Phys. Chem. Chem. Phys., **16**, 2842-2853 (2014)

4.1 Introduction

Predicting binding affinity between two non-covalently bound molecules is a challenging problem in molecular science. Calculating binding affinities using atomistic simulations can provide detailed molecular level insight into molecular recognition, and help inform fields such as structure-based drug design^{1–3} and self-assembly.^{4,5} For instance, an accurate and efficient method for predicting protein–ligand binding free energy can help screen a library of candidate compounds against a protein target, or assist in lead optimisation, by predicting the impact of chemical modifications. Hence this is an active field for the computational drug design community.^{1,6,7}

A broad class of methods for computing protein–ligand binding docks the ligand into the binding pocket and uses a scoring function to estimate the binding affinity.^{6,7} The scoring functions have explicit terms to model various contributions to the binding free energy, such as the hydrophobic effect, hydrogen-bonding, and further entropic contributions, which are usually fitted to experimental binding data. This docking and scoring approach is fast, but may not be accurate, due to training set bias and an approximate treatment of conformational entropy.

An alternative class of methods employs atomistic force fields to model the interatomic and intermolecular interactions. To describe protein–ligand binding the energy function is typically taken to be an empirical form, either with explicit water molecules, or an implicit solvent model; I employ AMBER⁸ in the present study.

A range of simulation methods have previously been developed to compute binding free energies using force field energy models and molecular dynamics (MD) or Monte Carlo (MC) simulations. Alchemical methods, where atoms of one ligand are transformed to those for another ligand, are used to compute relative binding affinity. Thermodynamic integration⁹ and free energy perturbation^{10–12} have been employed for alchemical free energy simulations. In another approach, the absolute free energy of binding is computed by equilibrium^{13,14} or non-equilibrium simulations^{15,16} along a physical pathway between the free and the bound ligand. These methods are formally rigorous, but are computationally expensive due to sampling limitations in MD or MC simulations of proteins. Another class of methods, including Molecular Mechanics/Poisson–Boltzmann Surface Area (MM/PBSA),^{17–20} Linear Interaction Energy (LIE),^{21–24} and others,²⁵ rely on MD simulation of only the free and bound states. These endpoint methods are relatively less expensive than pathway methods, since intermediate states are not considered, but still require adequate MD sampling of the end states, which can be challenging for protein sized systems. On the whole, current physics-based methods are

computationally much more expensive than docking and scoring based methods and therefore have found limited utility in applications such as virtual screening.

The force field-based methods are potentially more accurate than docking and scoring approaches, as they have been developed to account for explicit intermolecular interactions and are typically fitted using statistical mechanical theories. However, these methods can be computationally expensive, since MD or MC simulations are easily trapped in local minima of the potential or free energy surface for relatively long time scales.

The superposition approach provides an alternative formulation for global thermodynamics within the energy landscape framework.^{26–28} Here, the partition function is written as the sum of contributions from the catchment basins²⁹ of local potential energy minima.^{30–32} The contribution of each minimum can be estimated using the harmonic approximation, possibly with anharmonic³³ or quantum³⁴ corrections. To apply this procedure to calculate a binding free energy one can evaluate the free energy of the complex and the free molecules separately from databases of local minima for each species. This approach to binding free energy calculations is employed in the mining minima algorithm,³⁵ which has been successfully applied to various biomacromolecular systems, especially small host–guest systems. Benchmark superposition calculations for atomic and molecular clusters show that the energy landscape approach can be much faster than MD or MC based methods, especially for cases of broken ergodicity,^{36–39} since the superposition partition function is explicitly ergodic.

To apply the superposition method for large systems requires appropriate sampling, because the number of local minima increases exponentially with system size.^{40,41} A new method has recently been described to implement such sampling systematically, and was applied successfully to atomic cluster.⁴² Alternatively, the mining minima method has been extended to larger protein–ligand systems⁴³ by focusing the calculation on regions around the binding pocket. For example, in ref. 43, protein atoms were partitioned into three layers of different thickness with respect to the distance from ligand atoms. Atoms in the 7 Å layer closest to the ligand were free to move, while those in the middle layer of thickness 5 Å were fixed. Atoms in the outermost layer were deleted.

In the present contribution, I present a method conceptually similar to mining minima, but with key differences in the implementation, which aim to improve the accuracy and sampling efficiency. I again partition the protein atoms into three layers according to distance from the bound ligand. Atoms in the ‘inner’ region, adjacent to the ligand, are unconstrained, while those in the ‘intermediate’ region are treated using the local rigid body framework.⁴⁴ All atoms in the ‘outer’ layer were grouped as one

rigid body, but their contributions to the potential energy of the system are retained. The local rigid body framework is used to reduce the number of degrees of freedom, both in sampling minima and in the calculation of normal mode frequencies. All ligand atoms are fully flexible.

To benchmark the procedure I systematically increase the radius defining the innermost unconstrained region until the binding free energy converges. The key idea is that contributions from minima corresponding to alternative conformations of groups that are sufficiently distant from the binding site are expected to cancel between the free protein and the complex. Hence, I only need to sample consistent conformations for these degrees of freedom. The theory, described in Section 4.2.1, therefore corresponds to a factorisation of the partition functions for the protein, ligand, and complex. I therefore refer to the method as *a factorised superposition approach* (FSA).

I apply the FSA procedure to compute the binding free energy for human aldose reductase (5113 atoms) and one of its inhibitors, phenyl acetic acid (PAC). Human aldose reductase is an NADPH-dependent oxidoreductase, which catalyses the reduction of a variety of aldehydes and carbonyls, including monosaccharides. It is primarily known for catalysing the reduction of glucose to sorbitol, the first step in the polyol pathway of glucose metabolism.⁴⁵

The next section describes the theory underlying the factorisation procedure, the calculation of approximate free energies, local rigidification, and the sampling of local minima. I then describe the system setup for aldose reductase in Section 4.3, and discuss the conditions for convergence. I find that a flexible region of 14 Å, corresponding here to rigidification of about 80% of the protein, is required to obtain a converged binding free energy.

4.2 Methodology

4.2.1 Factorised superposition approach

I wish to estimate the binding free energy or the free energy change, ΔF , involved in forming a complex AB from non-covalent association of two molecules A and B. The standard free energy difference of this reaction is given by ref. 46–48

$$e^{-\beta\Delta F^0} = \frac{C^0}{8\pi^2} \frac{Z_{AB}}{Z_A Z_B}, \quad (1)$$

where Z_X is the configurational part of the single-molecule partition function of species $X \in \{A, B, AB\}$, C^0 is the standard concentration, and $\beta = 1/kT$, with k the Boltzmann constant and T the temperature. The above expression is derived using classical thermodynamics so that the momentum factors in the partition function of the free molecules and the bound complex cancel. The translational and rotational degrees of freedom have been integrated out from the configurational integral; see ref. 47 and 48 for a detailed derivation of the above expression.

I compute the partition functions [henceforth referring to the configurational integrals in Eq. (1)] using the superposition approach,^{26,28,30} where each Z_X is written as a sum of contributions from local minima of the potential energy surface. In this section, I describe the FSA framework, an extension of the superposition approach, which facilitates calculation of an approximate binding free energy from a subset of local minima.

The FSA framework was developed in order to provide a route to protein–ligand binding energies, where the number of relevant local minima becomes problematic for the standard superposition approach. To limit the number of minima, I assume that the contributions of analogous alternative conformations of functional groups that are sufficiently distant from the binding region cancel out. This scheme can be formalised by thinking in terms of the possible local conformations of distinct parts of the protein, such as backbone and side chain geometries. As a further simplification, I consider molecules that are not rotating or translating, and focus on the vibrational partition function for each local minimum.

To index the local minima I consider the possible conformations for each part of the molecule, and assume that I can identify them independently of the conformations adopted by the rest of the system (a factorisation). Each minimum can then be represented by a vector, $x = (x_1, x_2, \dots)$, where each component x_i for $i = 1, 2, 3, \dots$ identifies the local conformation of a region i . Some conformations of one region will preclude conformations of other regions, so the permitted combinations of x_i are restricted, which prevents further factorisation in general. I now identify local minima corresponding

to AB as $x_{AB} = (x_A, x_B)$. If certain local conformations are only possible in the complex, then the corresponding geometries in the separate A and B molecules are presumably high in energy, but can still be included in the possible conformations enumerated by x_A and x_B . To analyse the most plausible cancellation of contributions from alternative conformations that lie outside the binding region (the factorisation) I now assume that x_{AB} can be partitioned into two sets as $x_{AB} = (\mathbf{u}_{AB}, \mathbf{v}_{AB})$, as shown in Fig. 1. This formalism is designed to reflect our intuition that some local conformations are common to each molecule, while others associated with the binding region are not. The conformations collected in the \mathbf{v}_{AB} set therefore correspond to local structure that is identifiable in each of A, B and AB for all the conformations specified by the vector \mathbf{u}_{AB} . The corresponding regions in the separate A and B molecules are written as $\mathbf{u}_A, \mathbf{v}_A, \mathbf{u}_B$, and \mathbf{v}_B , and I assume that all possible conformations specified by \mathbf{v}_A and \mathbf{v}_B are also available in \mathbf{v}_{AB} for any \mathbf{u}_{AB} .

A significant simplification is possible if I need only consider a consistent reference conformation, \mathbf{v}_A^0 and \mathbf{v}_B^0 , respectively, for each group collected in \mathbf{v}_A and \mathbf{v}_B . In fact, this choice produces a combinatorial reduction in the number of minima that may need to be sampled. The analysis that follows defines the conditions under which this simplification will be valid. Furthermore, our local rigidification procedure⁴⁴ provides an ideal framework for implementing this approach, and enables us to determine a minimal set of states for estimating free energies of binding.

The partition function for separate A and B molecules factorises and I therefore consider

$$\begin{aligned} Z_A &= \sum_{\mathbf{u}_A} \sum_{\mathbf{v}_A} z_A(\mathbf{u}_A, \mathbf{v}_A) e^{-\beta V_A(\mathbf{u}_A, \mathbf{v}_A)} \\ &= \sum_{\mathbf{u}_A} z_A(\mathbf{u}_A, \mathbf{v}_A^0) e^{-\beta V_A(\mathbf{u}_A, \mathbf{v}_A^0)} \times \sum_{\mathbf{v}_A} \frac{z_A(\mathbf{u}_A, \mathbf{v}_A)}{z_A(\mathbf{u}_A, \mathbf{v}_A^0)} e^{-\beta [V_A(\mathbf{u}_A, \mathbf{v}_A) - V_A(\mathbf{u}_A, \mathbf{v}_A^0)]} \end{aligned} \quad (2)$$

and similarly for B. In Eq. (2), $z_A(\mathbf{u}_A, \mathbf{v}_A)$ and $V_A(\mathbf{u}_A, \mathbf{v}_A)$ are the vibrational partition function and potential energy for minimum $(\mathbf{u}_A, \mathbf{v}_A)$. The sum is over all local minima of A, identified via their \mathbf{u}_A and \mathbf{v}_A conformational assignment. Next, I define $f(\mathbf{u}_A, \mathbf{v}_A)$ as the free energy of a minimum

$$f(\mathbf{u}_A, \mathbf{v}_A) \equiv -\frac{1}{\beta} \ln Z_A(\mathbf{u}_A, \mathbf{v}_A) + V_A(\mathbf{u}_A, \mathbf{v}_A) \quad (3)$$

and rewrite Eq. (2) as

$$Z_A = \sum_{\mathbf{u}_A} \sum_{\mathbf{v}_A} e^{-\beta f(\mathbf{u}_A, \mathbf{v}_A)}. \quad (4)$$

I define a free energy shift, $\Delta f_A(\mathbf{u}_A, \mathbf{v}_A; \mathbf{v}_A^0)$, as the free energy difference between a given minimum $(\mathbf{u}_A, \mathbf{v}_A)$ and the corresponding reference $(\mathbf{u}_A, \mathbf{v}_A^0)$,

$$\Delta f_A(\mathbf{u}_A, \mathbf{v}_A; \mathbf{v}_A^0) \equiv f_A(\mathbf{u}_A, \mathbf{v}_A) - f_A(\mathbf{u}_A, \mathbf{v}_A^0). \quad (5)$$

The partition function for each molecule can then be written as

$$Z_X = \sum_{\mathbf{u}_X} e^{-\beta f_X(\mathbf{u}_X, \mathbf{v}_X^0)} \sum_{\mathbf{v}_X} e^{-\beta \Delta f_X(\mathbf{u}_X, \mathbf{v}_X; \mathbf{v}_X^0)}, \quad (6)$$

and the ratio of partition functions in Eq. (1) becomes

$$\begin{aligned} \frac{Z_{AB}}{Z_A Z_B} &= e^{-\beta \Delta F} \\ &= \frac{\sum_{\mathbf{u}_{AB}} e^{-\beta f_{AB}(\mathbf{u}_{AB}, \mathbf{v}_{AB}^0)} \sum_{\mathbf{v}_{AB}} e^{-\beta \Delta f_{AB}(\mathbf{u}_{AB}, \mathbf{v}_{AB}; \mathbf{v}_{AB}^0)}}{\sum_{\mathbf{u}_A} e^{-\beta f_A(\mathbf{u}_A, \mathbf{v}_A^0)} \sum_{\mathbf{v}_A} e^{-\beta \Delta f_A(\mathbf{u}_A, \mathbf{v}_A; \mathbf{v}_A^0)} \times \sum_{\mathbf{u}_B} e^{-\beta f_B(\mathbf{u}_B, \mathbf{v}_B^0)} \sum_{\mathbf{v}_B} e^{-\beta \Delta f_B(\mathbf{u}_B, \mathbf{v}_B; \mathbf{v}_B^0)}} \end{aligned} \quad (7)$$

As noted above, I require the $\mathbf{v}_{AB} = (\mathbf{v}_A, \mathbf{v}_B)$ conformations to appear in both the separate molecules and in the complex, and they must be identifiable for each minimum specified by different conformations in $\mathbf{u}_{AB} = (\mathbf{u}_A, \mathbf{u}_B)$. Next I introduce two assumptions, schematically described in Fig. 1, to simplify Eq. (7).

First, I assume that the shifts with respect to the reference conformation in the free energies Δf_X , are independent of \mathbf{u}_X when the conformations of the complex are chosen appropriately:

$$\Delta f_X(\mathbf{u}_X, \mathbf{v}_X; \mathbf{v}_X^0) \approx \Delta f_X(\mathbf{v}_X; \mathbf{v}_X^0) \quad \forall \mathbf{u}_X. \quad (8)$$

The summations over the common regions then factorise giving:

$$e^{-\beta \Delta F} = \frac{\left[\sum_{\mathbf{v}_{AB}} e^{-\beta \Delta f_{AB}(\mathbf{v}_{AB}; \mathbf{v}_{AB}^0)} \right] \left[\sum_{\mathbf{u}_{AB}} e^{-\beta f_{AB}(\mathbf{u}_{AB}, \mathbf{v}_{AB}^0)} \right]}{\left[\sum_{\mathbf{v}_A} e^{-\beta \Delta f_A(\mathbf{v}_A; \mathbf{v}_A^0)} \right] \left[\sum_{\mathbf{u}_A} e^{-\beta f_A(\mathbf{u}_A, \mathbf{v}_A^0)} \right] \left[\sum_{\mathbf{v}_B} e^{-\beta \Delta f_B(\mathbf{v}_B; \mathbf{v}_B^0)} \right] \left[\sum_{\mathbf{u}_B} e^{-\beta f_B(\mathbf{u}_B, \mathbf{v}_B^0)} \right]}. \quad (9)$$

Second, I assume that, for a given minimum, the shifts in the free energy relative to the reference conformation, $\mathbf{v}_{AB}^0 = (\mathbf{v}_A^0, \mathbf{v}_B^0)$, are the same in the complex and the separated molecules for all \mathbf{u}_{AB} , that is,

$$\Delta f_{AB}(\mathbf{v}_{AB}; \mathbf{v}_{AB}^0) \approx \Delta f_A(\mathbf{v}_A; \mathbf{v}_A^0) + \Delta f_B(\mathbf{v}_B; \mathbf{v}_B^0) \quad \forall \mathbf{u}_{AB}. \quad (10)$$

Note that, by construction, every \mathbf{v}_{AB} conformation in the numerator of Eq. (7). can be associated with a product of terms from the \mathbf{v}_A and \mathbf{v}_B sums in the denominator. Therefore, using Eq. (10), the factors with summations over the common region in Eq. (10). cancel, giving the final result

$$e^{-\beta\Delta F} \approx \frac{\sum_{\mathbf{u}_{AB}} e^{-\beta f_{AB}(\mathbf{u}_{AB}, \mathbf{v}_{AB}^0)}}{\sum_{\mathbf{u}_A} e^{-\beta f_A(\mathbf{u}_A, \mathbf{v}_A^0)} \sum_{\mathbf{u}_B} e^{-\beta f_B(\mathbf{u}_B, \mathbf{v}_B^0)}}. \quad (11)$$

I must therefore sum over members of the \mathbf{u}_{AB} minima and over local minima corresponding to all conformations of A and B in the same regions, with a Boltzmann weighting. A consistent set of local reference conformations \mathbf{v}_{AB}^0 must be used for the other regions corresponding to \mathbf{v}_{AB} .

In the present work, the common regions included only the protein atoms (molecule A) since the ligand (molecule B) was treated as fully flexible, reducing Eq. (11) to

$$e^{-\beta\Delta F} \approx \frac{\sum_{\mathbf{u}_{AB}} e^{-\beta F_{AB}(\mathbf{v}_{AB}^0)}}{e^{-\beta F_A(\mathbf{v}_A^0)} \times e^{-\beta F_B}}, \quad (12)$$

where

$$e^{-\beta F_X(\mathbf{v}_X^0)} \equiv \sum_{\mathbf{u}_X} e^{-\beta f_X(\mathbf{u}_X, \mathbf{v}_X^0)}, \quad X \in A, AB,$$

and

$$e^{-\beta F_B} \equiv \sum_{\mathbf{u}_B} e^{-\beta f_B(\mathbf{u}_B)}$$

are the free energies for the free molecules and the complex. Note that the free energy of the protein and complex depend on the reference configuration of the common regions. Eq. (12) is the working equation for the applications considered below. Since I am primarily interested in the convergence of the binding free energy with respect to the FSA framework, I do not include the $8\pi^2/C^0$ prefactor from Eq. (1), and I treat all molecules in vacuum for this initial benchmarking.

Eq. (12) is quite intuitive, with consistent reference conformations selected for regions of the protein that interact only weakly with the binding site. The derivation defines the validity of this approximation. In particular, it is clear that sampling over a small number of local minima where the conformations in the weakly interacting region are not consistent would introduce systematic errors. For large systems the number of possible conformations will be combinatorial, and randomly chosen conformations are unlikely to be in correspondence.

A straightforward method for implementing Eq. (12) is to sample local minima with the common region constrained in the reference conformation. This sampling is accomplished here using the local rigidification framework.⁴⁴ Our strategy for testing Eq. (12) is to check the convergence of the binding free energy as I expand the unconstrained region specified by \mathbf{u} . As a cross-validation, the result should be independent of the reference conformations specified by \mathbf{v}^0 .

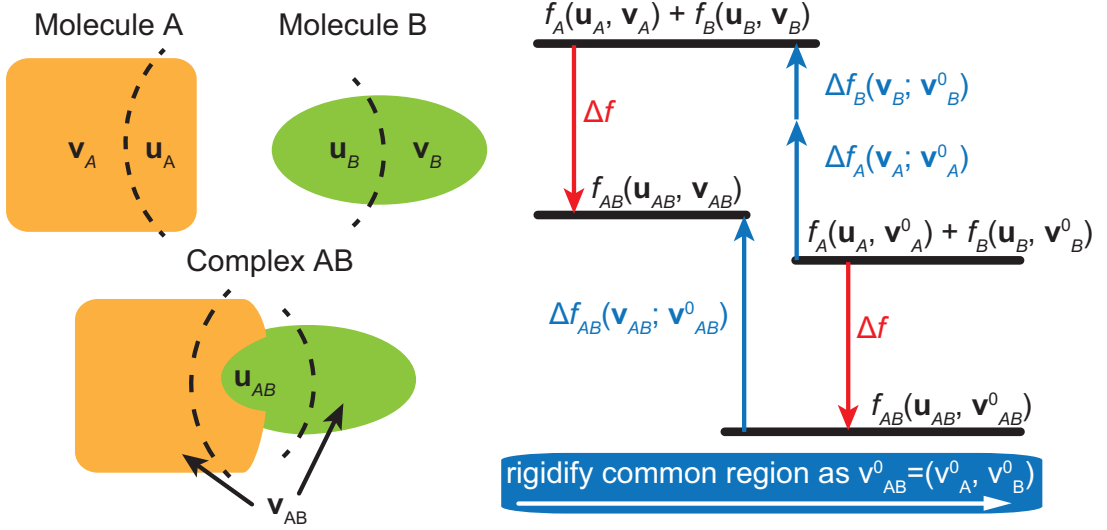


Figure 1 (left) Schematic representation of the conformational indexing vectors for molecules A and B, and for the complex AB. (right) Schematic representation for the free energies of one local minimum of A, B and AB, representing $f_A(\mathbf{u}_A, \mathbf{v}_A)$, $f_B(\mathbf{u}_B, \mathbf{v}_B)$ and $f_{AB}(\mathbf{u}_{AB}, \mathbf{v}_{AB})$. The difference, Δf , does not change if the shift corresponding to different \mathbf{v}_{AB} , $\Delta f_{AB}(\mathbf{v}_{AB}; \mathbf{v}_{AB}^0)$, is independent of \mathbf{u}_{AB} [Eq. (8)] and is, additive for $\Delta f_A(\mathbf{v}_A; \mathbf{v}_A^0)$ and $\Delta f_B(\mathbf{v}_B; \mathbf{v}_B^0)$ [Eq. (10)].

4.2.2 Free energy of local minima

In the harmonic approximation, the free energy, f , of a minimum is given by

$$e^{-\beta f} = \frac{e^{-\beta V_{\min i}}}{(\beta h \bar{\nu})^\kappa}, \quad \text{with} \quad \bar{\nu} = \left(\prod_i \nu_i \right)^{\frac{1}{\kappa}}, \quad (13)$$

where h is Planck's constant, V_{\min} is the potential energy of the minimum, and $\bar{\nu}$ is the geometric mean of the $\kappa = 3N - 6$ vibrational normal mode frequencies. For a fully flexible molecule, $\kappa = 3N - 6$ where N is the number of atoms. The number of vibrational modes is reduced if parts of the molecule are rigidified, as described in the next section. When applying the superposition formula I collect together the identical contributions for all permutation-inversion isomers of a given minimum, which

corresponds to weighting $e^{-\beta f}$ by $1/o$, with o the order of the corresponding point group.^{26,30,49} An additional factor that depends on the atomic composition of the system is needed to enumerate the distinct local minima precisely, but cancels from all thermodynamic quantities. Since the point group is C_1 for all the minima considered in the present work, $o = 1$. Eq. (13) also ignores overall translational and rotational contributions (Section 4.2.3.1 in ref. 47), which were found to make a negligible contribution to the free energy differences of interest in the present study.

4.2.2.1 Normal mode analysis in the local rigid body framework.

The cost of diagonalisation of the $3N \times 3N$ dimensional Hessian matrix required for calculating the normal mode frequencies for each minimum scales as $O(N^3)$. The computational expense is reduced when I consider the Hessian corresponding to local rigidification. Since the ligand is treated as fully flexible, its normal mode frequencies are computed by diagonalising the standard all-atom Hessian.

I need to address two issues in order to perform a normal mode calculation with local rigidification. First, the Hessian matrix of second derivatives required for the normal mode analysis has dimension $3N$ for N atoms. Rigidification reduces the dimensionality, and corresponds to a projection of the degrees of freedom of the constrained atoms onto the rotational and translational degrees of freedom of the rigid bodies. Second, the moment of inertia tensor is generally not diagonal for the kinetic energy expressed in the local rigid body coordinates. Hence I need two steps to calculate the corresponding normal modes, as detailed below.

First I establish our notation, denoting the number of rigid bodies by N_{RB} and the number of unconstrained atoms by N_{A} . In the angle-axis representation^{27,50} each rigid body I has six degrees of freedom: three representing the position of the centre of mass (translational degrees of freedom) $\mathbf{r}^I = \{r^I_1, r^I_2, r^I_3\}$, and three representing its orientation (rotational degrees of freedom) $\mathbf{p}^I = \{p^I_1, p^I_2, p^I_3\}$. For clarity, I employ capital letters for rigid bodies, and lower case for the sites in the rigid bodies. The coordinates of the sites, i , for rigid body I are denoted by $\mathbf{r}^I(i) = \{r^I_1(i), r^I_2(i), r^I_3(i)\}$, where

$$\mathbf{r}^I(i) = \mathbf{r}^I + \mathbf{S}^I \mathbf{x}^I(i); \quad i \in I. \quad (14)$$

I define $\mathbf{x}^I(i) = \{x^I_1(i), x^I_2(i), x^I_3(i)\}$ as the reference coordinates of the sites relative to the centre of mass of rigid body I , and \mathbf{S}^I as the rotation matrix constructed from the rotational degrees of freedom $\{\mathbf{p}^I\}$ (in the angle-axis representation) that rotates rigid body I from its reference frame to its current orientation,

$$\mathbf{S}' = \mathbf{I} + (1 - \cos \theta') \tilde{\mathbf{p}}^I \tilde{\mathbf{p}}^I + \sin \theta' \tilde{\mathbf{p}}^I, \quad (15)$$

with \mathbf{I} the identity matrix, $\theta' = \{(p_1^I)^2 + (p_2^I)^2 + (p_3^I)^2\}^{1/2}$ and $\tilde{\mathbf{p}}^I$ the skew-symmetric matrix obtained from the rotation vector \mathbf{p}^I :

$$\tilde{\mathbf{p}}^I = \frac{1}{\theta'} \begin{pmatrix} 0 & -p_3^I & p_2^I \\ p_3^I & 0 & -p_1^I \\ -p_2^I & p_1^I & 0 \end{pmatrix}. \quad (16)$$

Using the above notation, the Hessian corresponding to local rigid body coordinates is given by

$$\begin{aligned} \frac{\partial^2 V}{\partial r_\alpha^I \partial r_\beta^J} &= \sum_{i \in I} \sum_{j \in J} \frac{\partial^2 V}{\partial r_\alpha^I(i) \partial r_\beta^J(j)}, \\ \frac{\partial^2 V}{\partial r_\alpha^I \partial p_\beta^J} &= \sum_{i \in I} \sum_{j \in J} \sum_{a=1}^3 \frac{\partial^2 V}{\partial r_\alpha^I(i) \partial r_a^J(j)} \left[\frac{\partial \mathbf{S}^J}{\partial p_\beta^J} \mathbf{x}^J(j) \right]_a, \\ \frac{\partial^2 V}{\partial p_\alpha^I \partial p_\beta^J} &= \sum_{i \in I} \sum_{j \in J} \sum_{a=1}^3 \sum_{b=1}^3 \frac{\partial^2 V}{\partial r_b^I(i) \partial r_a^J(j)} \left[\frac{\partial \mathbf{S}^I}{\partial p_\alpha^I} \mathbf{x}^I(i) \right]_b \left[\frac{\partial \mathbf{S}^J}{\partial p_\beta^J} \mathbf{x}^J(j) \right]_a, \quad \text{for } I \neq J \\ \frac{\partial^2 V}{\partial p_\alpha^I \partial p_\beta^J} &= \sum_{i_1 \in I} \sum_{i_2 \in I} \sum_{a=1}^3 \sum_{b=1}^3 \frac{\partial^2 V}{\partial r_b^I(i_1) \partial r_a^I(i_2)} \left[\frac{\partial \mathbf{S}^I}{\partial p_\alpha^I} \mathbf{x}^I(i_1) \right]_b \left[\frac{\partial \mathbf{S}^I}{\partial p_\beta^I} \mathbf{x}^I(i_2) \right]_a + \sum_{i \in I} \sum_{a=1}^3 \frac{\partial V}{\partial r_a^I(i)} \left[\frac{\partial^2 \mathbf{S}^I}{\partial p_\alpha^I \partial p_\beta^I} \mathbf{x}^I(i) \right]_a, \end{aligned} \quad (17)$$

where I have used

$$\frac{\partial r_a^I(i)}{\partial p_\alpha^I} = \left[\frac{\partial \mathbf{S}^I}{\partial p_\alpha^I} \mathbf{x}^I(i) \right]_a, \quad i \in I. \quad (18)$$

The notation $[\dots]_a$ corresponds to the a -th component of the vector given inside the bracket. Further details of the derivations are given in Section 4.5.1.

To illustrate the computation of the normal modes, I first focus on the kinetic energy terms for the rigid bodies:

$$K_{RB} = \sum_I^{N_{RB}} \frac{1}{2} M^I (\dot{r}^I)^2 + \sum_I^{N_{RB}} \sum_{\alpha=1}^3 \sum_{\beta=1}^3 \frac{1}{2} J_{\alpha\beta}^I \dot{p}_\alpha^I \dot{p}_\beta^I, \quad (19)$$

where the mass of rigid body I is $M^I = \sum_{i \in I} m^i$ and $J_{\alpha\beta}^I$ is the corresponding moment of inertia tensor. I

choose to work in the moving frame of reference, where $\mathbf{S}^I = \mathbf{I}$, as I find diagonalisation of the inertia matrix the most straightforward procedure. Here the moment of inertia has the usual definition.

I now wish to transform to coordinates where the kinetic energy is diagonal, with $\mathbf{Q}^{I,T}$ and $\mathbf{Q}^{I,R}$ for the translational and rotational degrees of freedom of rigid body I , so that

$$K_{RB} = \sum_{I=1}^{N_{RB}} \left[\frac{1}{2} (\dot{\mathbf{Q}}^{I,T})^2 + \frac{1}{2} (\dot{\mathbf{Q}}^{I,R})^2 \right]. \quad (20)$$

For the translational degrees of freedom, the required coordinate transformation is a simple rescaling: $\mathbf{Q}^{I,T} = \sqrt{M^I} \mathbf{r}^I$. However, for the rotational degrees of freedom, I must first apply a coordinate transformation $\mathbf{w}^I = \mathbf{A}^I \mathbf{p}^I$, so that the moment of inertia becomes a diagonal matrix with diagonal elements $\Omega_\alpha^I (\alpha = 1, 2, 3)$.^{51,52} Then I can simply rescale the orientational coordinates by the moment of inertia $Q_\alpha^{I,R} = \sqrt{\Omega_\alpha^I} w_\alpha^I$.

More generally, the total kinetic energy of the system consists of contributions from the rigid bodies and free atoms, and I can write it as

$$K = \sum_{i=1}^{\eta} \frac{1}{2} \dot{Q}_i^2, \quad (21)$$

where $\eta = 3N_A + 6N_{RB}$ is the total number degrees of freedom. For the unconstrained atoms,

$Q_i = X_i \sqrt{m^i}$, where m^i is the mass of the atom corresponding to atomic Cartesian coordinate X_i .

The next step in computing the normal modes is to expand the potential energy, V , in a Taylor series around a local minimum configuration with potential energy V_{\min} up to second order in the Q coordinates:

$$V = V_{\min} + \frac{1}{2} \sum_{i,j=1}^{\eta} \frac{\partial^2 V}{\partial Q_i \partial Q_j} Q_i Q_j. \quad (22)$$

Here, \mathbf{Q} is understood as the deviation from the local minimum configuration, which is defined as the local origin of coordinates. The Hessian matrix $H_{ij} = \partial^2 V / \partial Q_i \partial Q_j$ can be diagonalised using a matrix \mathbf{B} , whose columns are the eigenvectors of \mathbf{H} with associated eigenvalues $\omega_i^2 = 4\pi^2 \nu_i^2$:

$$\sum_{j=1}^{\eta} H_{ij} B_{jk} = \omega_k^2 B_{ik}; \quad q_i = \sum_{j=1}^{\eta} B_{ij} Q_j, \quad (23)$$

where q_i are the normal mode coordinates. In this coordinate system the Hamiltonian H can be written as

$$H = V_{\min} + \frac{1}{2} \sum_{i=1}^{\eta} (\dot{q}_i^2 + \omega_i^2 q_i^2). \quad (24)$$

Due to the overall translational and rotational symmetries, there are six zero normal mode eigenvalues. The total number of vibrational degrees of freedom in local rigid body coordinates is therefore $\kappa = \eta - 6$, which is used in Eq. (13) to define the harmonic free energy of an individual minimum.

4.2.3 Basin-hopping parallel tempering

The minima used in Eq. (12) to compute the binding free energy were sampled using basin-hopping global optimisation.^{53–55} The basin-hopping method steps between local minima of the potential energy surface, proposing moves by perturbing the current minimum, and accepting or rejecting the new structure obtained after minimisation using criteria such as the energy difference.^{53–55} I used the group rotation moves⁵⁶ described in Section 4.2.3.1 for perturbing the conformation of the current minimum in both the unconstrained inner and locally rigid intermediate regions. The perturbed conformation was minimised using a modified L-BFGS algorithm⁵⁷ with a tolerance of 0.001 kcal mol⁻¹ Å⁻¹ on the root mean square force. The new minimum was accepted or rejected using a Metropolis criterion based on the potential energy difference with respect to the previous minimum. Since the Metropolis criterion is based on the energy difference between local minima, all downhill barriers on the potential energy surface are removed. Uphill barriers are reduced to the difference in energy of the two minima. The minimisation and reduced barriers permit large perturbations of geometry, leading to effective sampling of the low energy regions of the potential energy surface of interest.

To enhance the sampling I employed the basin-hopping parallel tempering (BHPT) approach.⁵⁸ Conventional parallel tempering involves carrying out a parallel set of canonical Monte Carlo simulations at a range of temperatures, with periodic exchange attempts between the runs.^{59,60} In the BHPT approach the replicas evolving at different temperatures are all basin-hopping runs⁵⁸ and the exchanges are between the current minima in adjacent replicas.

4.2.3.1 Group rotations.

To propose perturbed conformations within each basin-hopping replica, generalised rotation moves were developed. This scheme allows arbitrary groups of atoms to be rotated about an axis defined by a bond vector, maintaining maximum flexibility without introducing reliance on standard topologies. Each group i has an associated user specified selection probability, $P(\text{select})_i$, and maximum rotation angle, θ_i^{\max} , to allow for further fine tuning of the conformational sampling. These perturbations are referred to as group rotation moves.⁵⁶ During each basin-hopping step:

1. For each group i , a random number ρ_1 is drawn between zero and one. If $P(\text{select})_i > \rho_1$ then the group is rotated in this step.

2. A second random number ρ_2 in the range $[-0.5, 0.5]$ is drawn and the rotation angle to be applied to the group is calculated as $\theta_i = 2\pi\rho_2\theta_i^{\max}$, where θ_i^{\max} is the maximum desired rotation angle for group i as a fraction of 2π .

3. The bond vector that connects the group to the rest of molecule is calculated and normalised before being scaled by θ_i .

4. For an atom with position vector \mathbf{r} , the rotation matrix \mathbf{S} is generated using an implementation of Rodrigues' rotation formula,^{62,63}

$$\mathbf{S}\mathbf{r} = \left[(\mathbf{I}\cos\theta) + \hat{\mathbf{k}}_x \sin\theta + \hat{\mathbf{k}}\hat{\mathbf{k}}^T (1 - \cos\theta) \right] \mathbf{r}, \quad (25)$$

where \mathbf{I} is the identity matrix, $\hat{\mathbf{k}}$ is the rotation axis, and $\hat{\mathbf{k}}_x$ is the 'cross-product matrix':

$$\hat{\mathbf{k}} = \begin{pmatrix} k_1 \\ k_2 \\ k_3 \end{pmatrix}, \quad \hat{\mathbf{k}}_x = \begin{pmatrix} 0 & -k_3 & k_2 \\ k_3 & 0 & -k_1 \\ -k_2 & k_1 & 0 \end{pmatrix}. \quad (26)$$

This scheme was initially developed to allow for comprehensive sampling of small ligands, but in the current work it has been adapted to sample the rotameric states of protein side chains. Fig. 2 shows the rotatable groups used to explore the conformations of the LYS side chain as an illustration. I define up to three such rotatable groups for each amino acid side chain, where atoms are rotated about the C_α - C_β , C_β - C_γ and C_γ - C_δ bonds. For simplicity, I set $P(\text{select})_i = 0.025$ for all groups, giving an average of 5.5 rotations per basin-hopping step for the 220 groups present when $R = 14 \text{ \AA}$ (see Section 4.3.4). The maximum rotation amplitude θ_i^{\max} for each group was chosen based on the group's size and spatial extent, in an effort to achieve the largest possible step size while minimising possible atom clashes

following a rotation. The values used in the current work can be found in Table A1 along with associated input files.

While the conformational changes during sampling are mainly determined by the group rotation of side chains and ligand, I also included small (0.1 Å) random Cartesian perturbations for all atoms, including the backbone, at every basin-hopping step. In addition, the backbone was free to move during minimization in the free and locally rigid regions to accommodate side chain/ ligand movement. Thus, the backbone moves during the sampling. To estimate the contribution of the backbone movement, I looked at eight aldose reductase crystal structures with different ligands bound, which were obtained from the Protein Data Bank. Among these complexes, the smallest ligand has 18 atoms and the largest has 49. The highest C_{α} -RMSD between the one I used as a starting point and any other is 0.723 Å for the whole protein and 0.609 Å for the residues within 16 Å of the ligand (Table A2). These small differences in backbone conformation reflect the fact that the backbone conformation is quite well defined for the species considered in the FSA procedure.

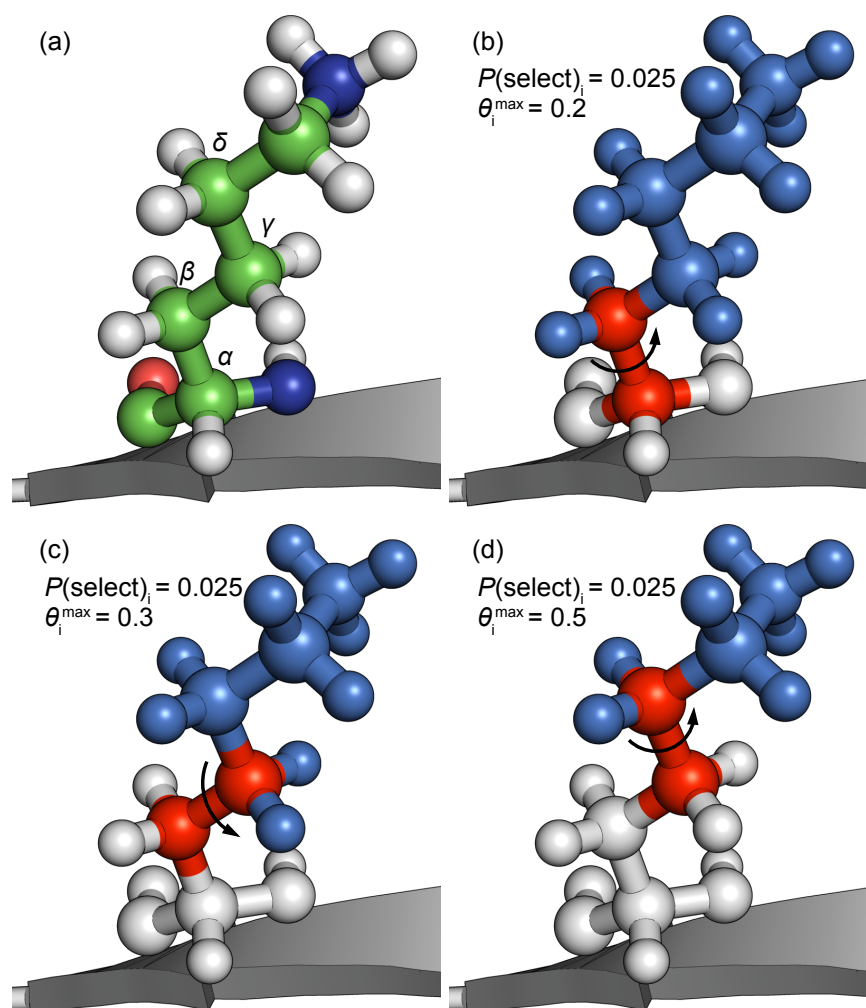


Figure 2 The amino acid lysine (LYS) (a) coloured by element with carbon atoms labelled. (b)–(d) show the $\alpha\beta$, $\beta\gamma$ and $\gamma\delta$ groups that can be rotated during basin-hopping, with their associated selection probabilities $P(\text{select})_i$ and maximum max rotation amplitudes θ_i^{\max} . The axis of rotation is shown in red, while the atoms to be rotated are shown in blue. The graphical representations in Fig. 2–5 were prepared using the Pymol program.⁶¹

4.3 Application to human aldose reductase

I employ the binding of human aldose reductase⁶⁴ with phenyl acetic acid (PAC) as a model system to test the factorisation superposition approach (see Fig. 3). For the purposes of this study, since I am not focussing on the catalytic activity of the enzyme, the NADP⁺ cofactor of the enzyme is considered to be part of the protein. The details of the simulation and local rigidification are described in Sections 4.3.1 and 4.3.2, respectively. The goals of the present work are to test the following two hypotheses. First, that the binding free energy should converge if the active binding site region is sufficiently large. Second, that the binding free energy should then be independent of the configuration of the inactive region. These hypotheses are tested by computing the binding free energy for systematic rigidification with three different reference conformations and examining the convergence to identify the maximum rigidification (factorisation) for which Eq. (12) holds.

4.3.1 Simulation set up

The simulations were performed using the AMBER ff99SB force field⁸ for the protein. Parameters for NADP⁺ were obtained from the AMBER parameter database.⁶⁵ The PAC ligand was parametrised using the General Amber Force Field^{66,67} with RESP^{68,69} charges generated iteratively using GAMESS-US.⁷⁰ A cutoff radius of 999.99 Å was used for non-bonded interactions, effectively including all pair-wise interactions. To evaluate the influence of the reference conformation, corresponding to \mathbf{v}^0_A in Eq. (12), I prepared three initial conformations with different geometries for the rigid region. One conformation, named ‘St-1’, was obtained from the Protein Data Bank (PDB code 2INE).⁶⁴ The other conformations, named ‘St-2’ and ‘St-3’, were prepared using a small number of basin-hopping steps starting from St-1 without any rigidification. St-2 and St-3 are arbitrarily chosen local minima. Similar convergence of the binding free energy with respect to R illustrates the point that any reference conformation \mathbf{v}^0 may be used in the Factorised Superposition Approach (FSA), as long as it is used consistently throughout the minima sampling and normal mode calculations. Fig. 4(a) shows St-1 and St-2 aligned on all atoms (RMSD 1.5 Å), while Fig. 4(b) shows the alignment for St-1 and St-3 (RMSD 1.9 Å). The main differences are the partial unfolding of a helix in St-2 and St-3, respectively. Most of the calculations were performed in vacuo to reduce the computational cost and facilitate more thorough benchmarking. An accurate solvation model is not required for the present study since the objective is to test the factorisation approach, rather than compare directly with experiment. Calculations involving explicit solvent will be the focus of future work, as will be

discussed in Section 4.3.6. In the present contribution I have simply relaxed the key local minima using an implicit solvent model to check that the convergence criteria are robust.

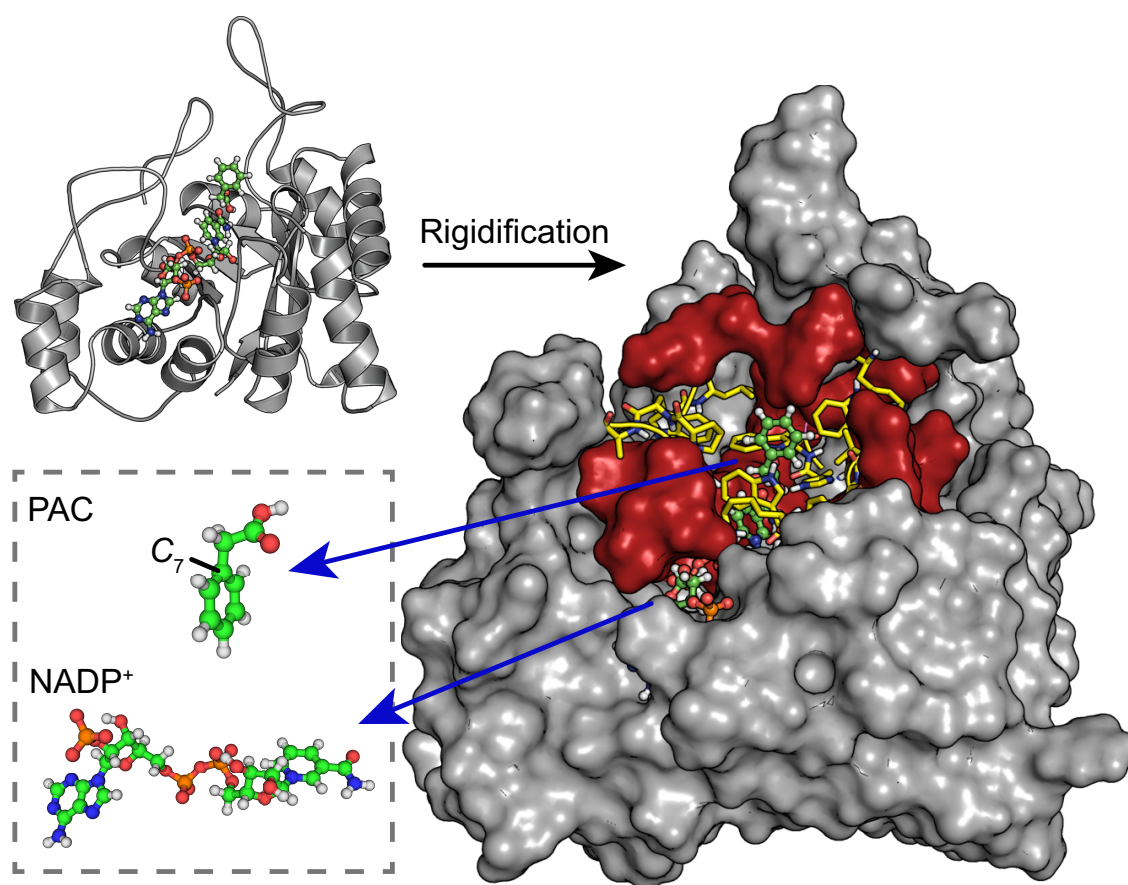


Figure 3 Cartoon (left top) and rigidified representations (right) of St1-Comp-R12. In the rigidified structure, the yellow lines represent the unconstrained inner region, the red surface shows the locally rigidified intermediate region, and the gray part is the outer region, rigidified as one group. The ligand PAC with the atom labels used in the text and the cofactor NADP⁺ are described in the insert.

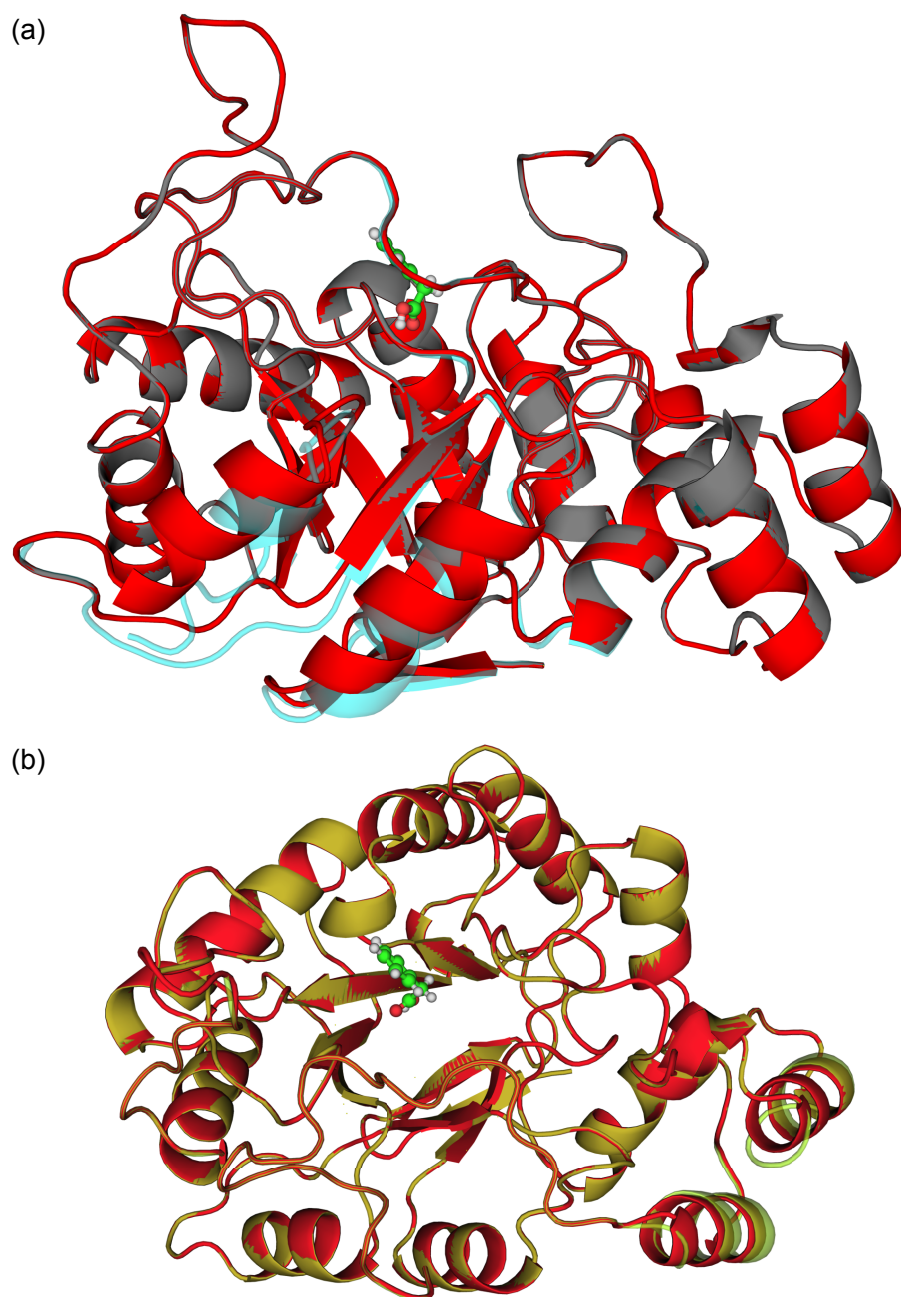


Figure 4 (a) Cartoon descriptions of St-1 (red) and St-2 (sky blue) after alignment. The blue color of St-2 is translucent, thus the overlapped region looks gray. (b) St-1 (red) and St-3 (yellow).

4.3.2 Systematic rigidification

For each structure, the free energy calculations were performed on multiple rigidified versions of the protein. The rigidification was applied systematically to fewer protein atoms, with the corresponding complex initially defined from identical protein and ligand coordinates. I determined the rigidified regions using the distance, R (in Angstroms), from the C_7 atom of the PAC ligand, labelled in Fig. 3. The unconstrained inner layer consisted of all atoms of amino acid residues having any atom within a radius R of the C_7 reference atom. If any atoms of a residue lay between radii R and $R + 1$ then I rigidified their peptide bonds, sp^2 centres, and aromatic rings. This set formed the intermediate layer with local rigid bodies. Atoms in the outer layer were rigidified as a single group. Fig. 3 shows the resulting rigidification scheme for the complex with a threshold value of $R = 12$ Å defined for St-1 (denoted St1-Comp-R12) and the details of the local rigidification are shown in Fig. 5. For St-1, six different rigidified versions were used, corresponding to $R = 6, 8, 10, 12, 14$ and 16 Å. For St-2, four versions ($R = 8, 10, 12, 14$ Å) were prepared, and for St-3, two versions ($R = 12, 14$ Å) were prepared. In each case the cofactor $NADP^+$ was part of the rigidified region. The number of atoms in each group is summarised in Table 1 as a function of R .

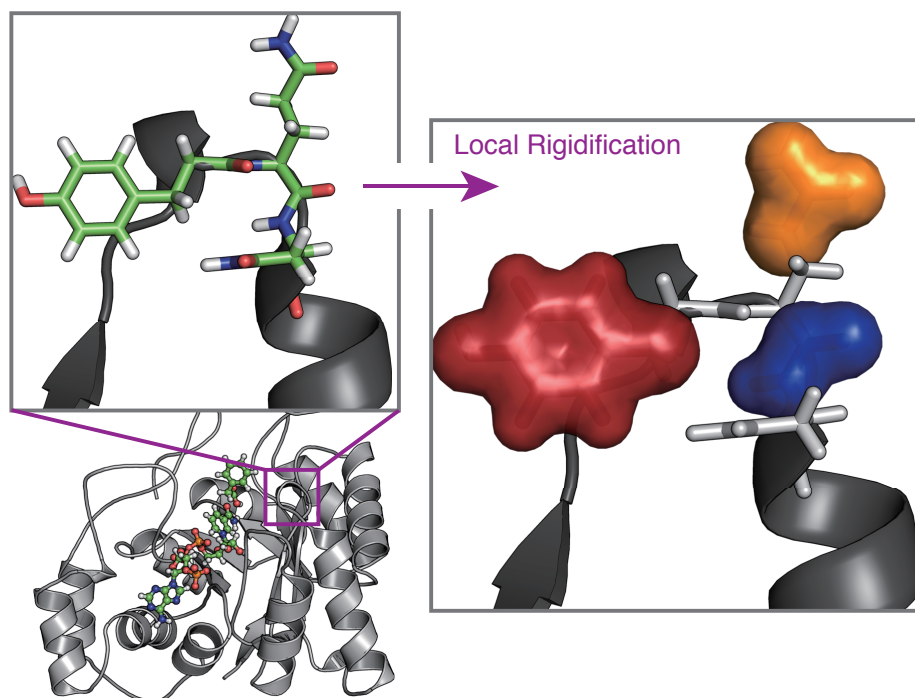


Figure 5 Examples of the local rigidifications corresponding to the intermediate region of Fig. 3. The image at the top left shows residues 47–49 in human aldose reductase. In the local rigidification (right), the red group is an aromatic ring corresponding to the TYR47 residue, the blue group corresponds to a peptide bond between GLN48 and ASN49, and the orange group is a trigonal centre (an amide group in this case) in the side chain of GLN48.

St	Radius R (Å)	% Rigid (protein)	Number of rigidified atoms			
			Intermediate (# LRB)	Outer	κ_A	κ_{AB}
1,2	6	99	0 (0)	5091	66	120
	8	97	36 (7)	4903	564	618
	10	92	92 (17)	4640	1245	1299
	12	87	114 (25)	4338	2133	2187
	14	80	135 (29)	3972	3192	3246
1	16	78	100 (21)	3886	3507	3561
3	12	88	137 (27)	4378	1956	2010
	14	81	146 (32)	3992	3117	3171

Table 1 The binding free energy calculations were performed with the protein (molecule A) atoms separated into three different regions. The inner region is fully flexible, the intermediate region consists of local rigid bodies (LRB), and the outer region is treated as a single rigid body. The total number of atoms in the ligand, protein (including the cofactor NADP⁺) and complex (molecule AB) are 18, 5113 and 5131, respectively. The table gives the number of degrees of freedom for protein (κ_A) and complex (κ_{AB}). The number of degrees of freedom for the ligand is $\kappa_B = 48$

4.3.3 Sampling local minima

The BHPT method (Section 4.2.3)⁵⁸ implemented in our GMIN⁷¹ program was used to sample local minima for the protein and complex for the different R values with both reference structures. All BHPT simulations were performed with 12 replicas and temperatures exponentially spaced between 97 K and 2435 K. Minimisation was performed using a modified version of the L-BFGS⁵⁷ algorithm with a tolerance of 0.001 kcal mol⁻¹ Å⁻¹ for the root mean square force. Minima with energies within 0.01 kcal mol⁻¹ were considered duplicates and excluded from the set used for computing the free energy. For the BHPT run for the complex of St-1, the probability of escape from the previous minimum is 37% at the lowest temperature and 62% at the highest temperature, which corresponds to an efficient choice of parameters.

For each minimum, normal mode frequencies were computed using our OPTIM⁷² program and harmonic free energies were obtained from Eq. (13) as the database of minima expanded (Fig. 6). Sampling was terminated when the change within the previous 2400 basin-hopping steps was less than 0.01 kcal mol⁻¹. Table 2 gives the total number of basin-hopping steps and the number of distinct

minima sampled for the different simulations. Among these minima, only a few hundred contribute significantly to the super- position sums in Eq. (12), and this number increases with R , as expected. As an example, I show how F_{AB} varies for St-1 in Table 3.

A single basin-hopping run of 1000 steps was performed for the ligand using a temperature of 300 K in the accept/reject step. Only the lowest two local minima contribute significantly to the partition function of the noninteracting ligand.

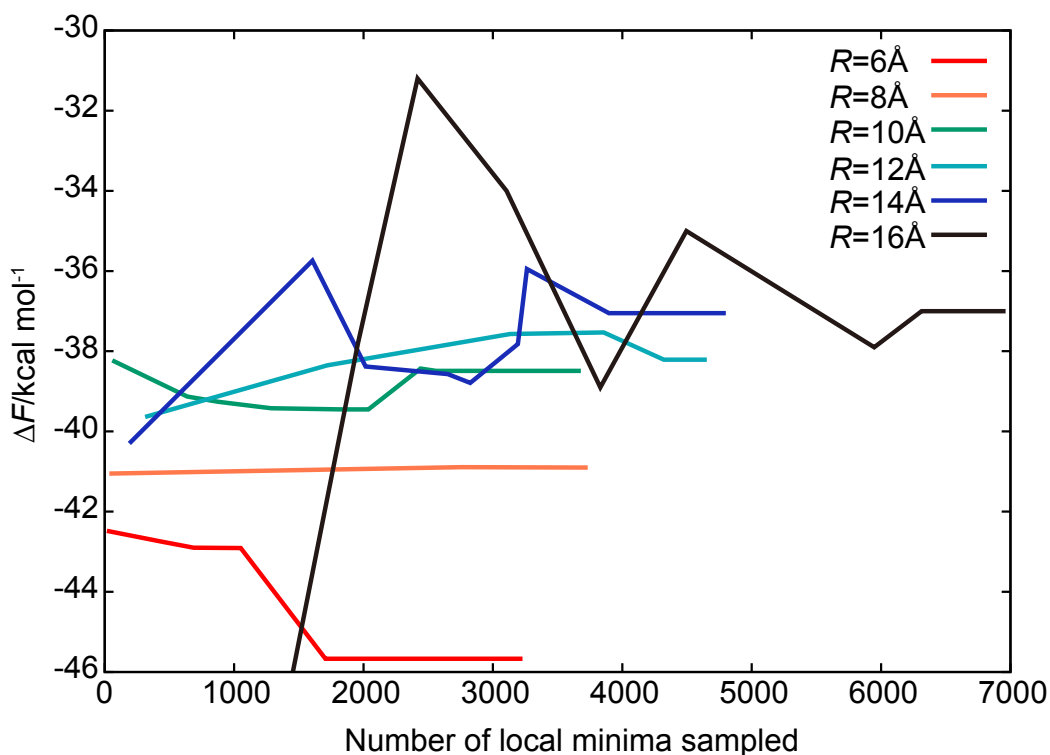


Figure 6 Binding free energies as a function of the number of distinct local minima sampled, corresponding to the progress of the BHPT run. Results are shown for six different values of the radius R , which defines the unconstrained region.

St	R	Total BH steps		# Minima obtained		Free energies (kcal mol ⁻¹)		
		Complex	Protein	Complex	Protein	Complex	Protein	ΔF
1	6	50 723	46 602	3229	2298	-9920.3	-9886.1	-45.6
	8	25 002	24 106	3733	3595	-9464.1	-9434.6	-40.9
	10	24 072	20 859	3680	2861	-8921.4	-8894.4	-38.4
	12	28 467	29 130	4873	4606	-8221.7	-8194.9	-38.2
	14	26 351	23 930	4800	4056	-7348.8	-7323.2	-37.0
	16	47 880	41 172	6962	6200	-4052.7	-4027.1	-37.0
2	8	18 303	19 473	2117	2259	799.3	822.7	-34.8
	10	15 610	31 317	1593	2091	1337.2	1362.1	-36.3
	12	17 183	16 474	2400	2351	2021.2	2046.6	-36.8
	14	19 235	15 030	3932	2899	2876.9	2901.9	-36.4
3	12	38 676	31 608	4388	3394	-5104.1	-5077.9	-37.5
	14	33 816	26 340	3382	3405	-4398.7	-4372.9	-37.2

Table 2 Total basin-hopping (BH) steps for 12 temperatures and the number of distinct local minima obtained for the complex (AB) and the protein (A). The binding free energies (ΔF) are calculated from the free energies of the complex (F_{AB}), protein (F_A) and ligand (F_B). A converged value of $F_B = 11.4$ kcal mol⁻¹ is obtained from the two lowest minima characterised in a BH run of 1000 steps

N_{\min}	$R = 6 \text{ \AA}$	$R = 8 \text{ \AA}$	$R = 10 \text{ \AA}$	$R = 12 \text{ \AA}$	$R = 14 \text{ \AA}$	$R = 16 \text{ \AA}$
1	-9916.821	-9463.359	-8920.421	-8219.631	-7346.759	-4049.182
10	-9917.676	-9464.074	-8921.024	-8220.497	-7347.315	-4050.390
30	-9920.302	-9464.075	-8921.306	-8220.884	-7348.228	-4050.837
50			-8921.314	-8221.049	-7348.355	-4051.481
70			-8921.416	-8221.126	-7348.515	-4052.084
90				-8221.161	-7348.649	-4052.097
110				-8221.589	-7348.732	-4052.173
130				-8221.667	-7348.777	-4052.181
150				-8221.669	-7348.796	-4052.634
170					-7348.801	-4052.636
190					-7348.804	-4052.683

Table 3 Free energies of the complex (kcal mol⁻¹) for reference St-1 using the N_{\min} lowest minima. The free energies changing by more than 0.001 kcal mol⁻¹ from the previous value are summarised below. The final values correspond to F_{AB} in Table 2

4.3.4 Convergence of the free energy with the size of the unconstrained region

I calculated binding free energies, ΔF , using Eq. (13) for St-1, St-2 and St-3 as a function of R , as shown in Fig. 7. For St-1, ΔF increases from $R=6\text{\AA}$ to $R=14\text{\AA}$ and appears to have converged at $R = 14\text{\AA}$. The ΔF values obtained for St-2 deviate significantly from that of St-1 at $R = 8\text{\AA}$, but at $R = 10\text{\AA}$ ΔF approaches the value obtained at $R = 14$ and 16\AA for St-1. Similar ΔF values are also obtained for St-3. ΔF at $R = 14, 16\text{\AA}$ for St-1, $R=12, 14\text{\AA}$ for St-2 and $R=12, 14\text{\AA}$ for St-3 are within 1.1 kcal mol^{-1} , even though the number of degrees of freedom (κ_X) and absolute free energies (F_X) are quite different for each R , as detailed in Tables 1 and 2. Thus, I conclude that the factorisation superposition approach seems to be applicable for this system with $R \geq 14\text{\AA}$, for appropriate reference conformations.

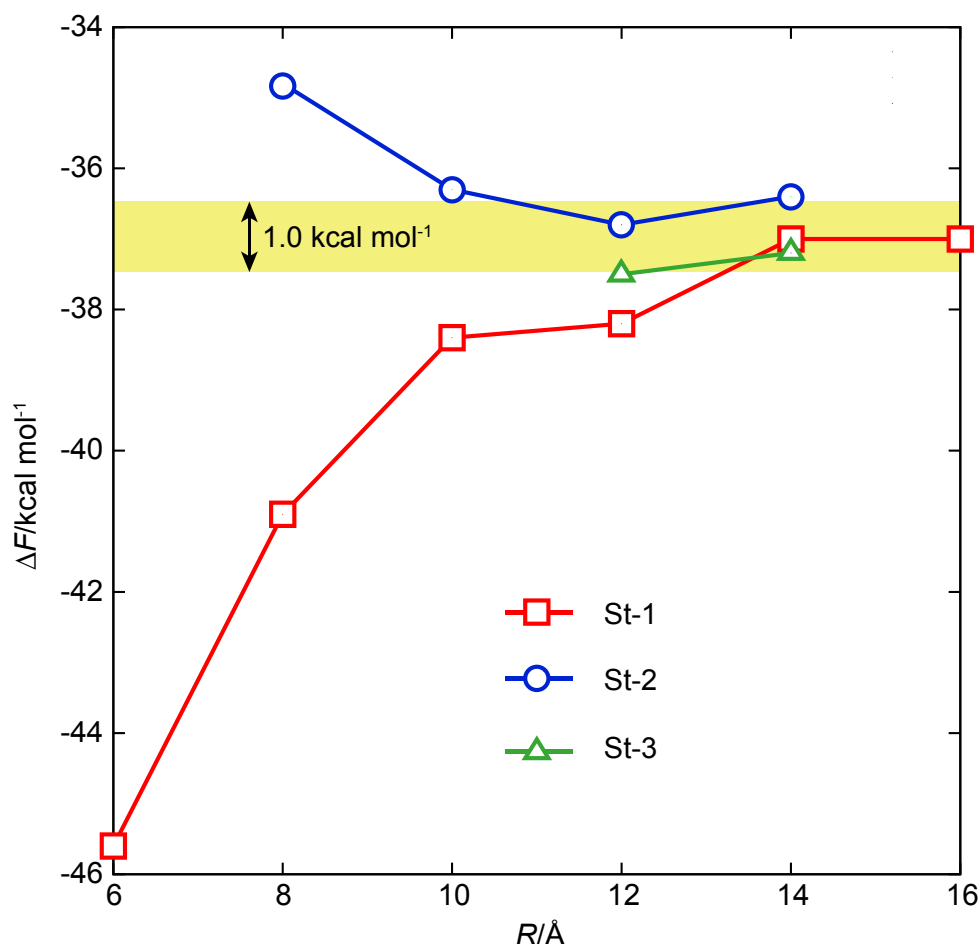


Figure 7 Binding free energies as a function of the rigidification radius, R . Results for St-1, St-2 and St-3 are shown in red, blue and green, respectively. The shaded region represents 1.0 kcal mol^{-1} around the average converged value.

4.3.5 Computational cost

In the BHPT sampling using GMIN, each basin-hopping step for St1-Comp-R14 takes about 3.1 times longer than for St1-Comp-R6 on average, because the coordinate space is larger for St1-Comp-R14. For the normal mode analysis using OPTIM, the diagonalisation of the Hessian matrix for one minimum with $\kappa = 3246$ ($R = 14$) and $\kappa = 15387$ (without any rigidification) took 6 minutes and 46 minutes of cpu time on average, respectively, the computational time scales roughly as $\kappa^{1.5}$, as suggested by the data in Fig. A1.

In spite of the speedup achieved using the rigid body framework, normal mode calculations for the protein and complex minima are still relatively expensive. It is therefore desirable to use as few minima as possible in the superposition sums. Due to the Boltzmann weight in Eq. (13), the low-energy minima dominate these sums. Table 4 shows the binding free energy computed using minima whose energies lie within a cutoff (E_c) of the global minimum energy. I find that the binding free energy is determined by minima with energies within $10 kT$ of the global minimum at $T = 298$ K. This cutoff corresponds to a small fraction of the total number of minima sampled for the protein and complex. For example, a cut-off of $10 kT$ applied to the database of minima for the $R = 14$ Å simulations with St-1 drastically reduces the number of minima of the complex from 4452 to 149. The number of relevant minima for smaller R is even less. A substantial reduction in the total computational cost can therefore be achieved by restricting the normal mode calculations to the low-energy minima.

E_{cut}	$R = 6$ Å	$R = 8$ Å	$R = 10$ Å	$R = 12$ Å	$R = 14$ Å	$R = 16$ Å
2	−42.8	−40.9	−40.6	−39.0	−36.7	−35.2
5	−45.6	−40.8	−38.4	−39.5	−37.7	−36.5
10	−45.6	−40.8	−38.4	−38.2	−37.0	−37.0
20	−45.6	−40.8	−38.4	−38.2	−37.0	−37.0

Table 4 Binding free energy (kcal mol^{−1}) for reference St-1 using protein and complex minima with energies within E_{cut} of the global minimum. The binding free energy computed using all the minima is given in Table 2. E_{cut} is in units of kT for $T = 298$ K

4.3.6 Incorporating solvent effects

The converged binding free energy was found to be approximately $-36.8 \text{ kcal mol}^{-1}$, corresponding to a standard binding free energy of $-29.8 \text{ kcal mol}^{-1}$, which is significantly lower than experimental binding affinity of $-5.5 \text{ kcal mol}^{-1}$.⁶⁴ I suspected that this discrepancy is primarily due to the absence of solvent effects. To test this hypothesis, I repeated the calculation for $R = 14 \text{ \AA}$ with St-1, using the generalized Born implicit solvent model, as implemented in AMBER.⁷³ I relaxed the lowest 500 minima identified in vacuum and recomputed the normal mode frequencies for both the protein and the complex. I did not resample minima for these tests, since the vacuum and the corresponding recomputed potential energies in the implicit solvent were found to be highly correlated (Fig. A2). Both the ligand minima were also relaxed using for this case implicit solvent. The resulting binding free energy was $\Delta F^0 = -8.4 \text{ kcal mol}^{-1}$, much closer to the experimental value. I expect that sampling with a more accurate implicit solvent model, such as linearized Poisson–Boltzmann,⁷⁴ would further improve the agreement with experiment.

I note that, in principle, the FSA framework can also be applied for explicit solvent. However, a large number of explicit water molecules would significantly increase the number of minima, and further work would be needed to sample these structures efficiently, probably by including solvent degrees of freedom in the factorisation. Including a few explicit water molecules might be also desirable, for example, in situations where the crystal structure contains bridging water molecules between the ligand and the protein.

4.4 Conclusions

I have presented a new method based on potential energy landscape theory,²⁶ the factorisation superposition approach (FSA), for computing the binding free energy of a protein–ligand complex. In this scheme the free energy of the free and bound molecules are computed using the superposition approach from a database of local potential energy minima. Due to the exponential increase in the number of minima with system size, exhaustive sampling is not feasible for a protein sized system. The FSA approach addresses this problem by focusing the calculation on protein atoms that interact strongly with the ligand. In Section 4.2.1 I presented the theory for factorising the conformational space of the protein and complex into two regions based on the size of the binding pocket. The factorisation facilitates estimation of the binding free energy using minima corresponding to fluctuations of the binding region, thereby reducing the number of degrees of freedom significantly. I describe the approximations under which such a factorisation is valid, employing a local rigid-body framework⁴⁴ to implement the FSA by treating atoms further from the binding site as collections of local rigid bodies. This procedure reduces the number of active degrees of freedom, but retains all the terms in the force field.

I applied the FSA method to calculate the free energy change for ligand binding with human aldose reductase protein while varying the size of the binding region. I performed the calculations for three different conformations of the rigid part of the protein and for different sizes of the binding pocket. For a given conformation of the rigidified region, I found that the binding free energy converged to within 1 kcal mol⁻¹ as the size of the binding pocket was increased to about 14 Å, corresponding to an 80% reduction in the number of protein degrees of freedom. The converged binding free energy for all three conformations were found to be within 1.1 kcal mol⁻¹, suggesting weak interactions between the ligand and protein atoms beyond 14 Å.

Several further improvements in the accuracy and speed of the FSA method as presented here can be envisioned. Larger systems are likely to derive a greater benefit from the factorisation scheme, because the whole region unrelated to ligand binding can be rigidified into a single unit, with only six rigid-body degrees of freedom. A surprising result of this study is that, even though the number of minima increased rapidly with the size of the unconstrained region around the binding pocket, the number of thermally relevant minima remained small, of the order of few hundred conformations. Anharmonicity corrections^{33,75,76} could improve the accuracy of the method, and the computationally intensive minima sampling and normal mode calculations should be highly amenable to distributed computing.

One key aspect of the FSA approach is the rigidification of large protein regions distant from the binding site. This approach assumes that the configurations of such regions change relatively little upon ligand binding. For proteins with significant allosteric effects,^{77,78} the regions should be rigidified in smaller domains, to avoid freezing out the protein allostery. The local rigid body approach, used in the ‘intermediate region’, and group rotations for sampling should still be applicable for any ligand binding system.

The converged radius for the flexible region obtained in the present work, $R = 14 \text{ \AA}$, is not expected to be universal, and other protein–ligand combinations will require analogous convergence checks. However, for a given protein, the value of R is likely to be transferable for different ligands of comparable size.

In future work I will consider solvent effects in more detail, and present comparisons with alternative approaches for calculating the free energy difference. Our main purpose in the present work was to demonstrate the convergence of the FSA scheme. I hope that, with further benchmarking and computational optimisation, the FSA method could facilitate screening calculations associated with drug design.

4.5 Appendix

4.5.1 Hessian in the local rigid body coordinates

In this section, I detail the derivations of the Hessian in the local rigid body coordinates. Our starting point is the first derivatives of the potential energy. For the translational degrees of freedom, \mathbf{r}^I , this is given by

$$\frac{\partial V}{\partial r_\alpha^I} = \sum_{i \in I} \sum_{a=1}^3 \frac{\partial V}{\partial r_a^I(i)} \frac{\partial r_a^I(i)}{\partial r_\alpha^I} = \sum_{i \in I} \frac{\partial V}{\partial r_\alpha^I(i)}. \quad (\text{A.1})$$

Additionally, the first derivative of the potential energy with respect to the rotational degrees of freedom, \mathbf{p}^I , gives

$$\frac{\partial V}{\partial p_\alpha^I} = \sum_{i \in I} \sum_{a=1}^3 \frac{\partial V}{\partial r_a^I(i)} \frac{\partial r_a^I(i)}{\partial p_\alpha^I} = \sum_{i \in I} \sum_{a=1}^3 \frac{\partial V}{\partial r_a^I(i)} \left[\frac{\partial \mathbf{S}^I}{\partial p_\alpha^I} \mathbf{x}^I(i) \right]_a. \quad (\text{A.2})$$

I have employed the following relations in the above partial derivatives

$$\mathbf{r}^i = \mathbf{r}^I + \mathbf{S}^I \mathbf{x}^I(i); \quad i \in I, \quad (\text{A.3})$$

$$\frac{\partial r_a^I(i)}{\partial r_\alpha^I} = \delta_{a\alpha}, \quad (\text{A.4})$$

$$\frac{\partial r_a^I(i)}{\partial p_\alpha^I} = \left[\frac{\partial \mathbf{S}^I}{\partial p_\alpha^I} \mathbf{x}^I(i) \right]_a. \quad (\text{A.5})$$

The second derivatives then follow in a similar manner. There are four separate cases to consider, and I derive them below for each case

$$\frac{\partial^2 V}{\partial r_\alpha^I \partial r_\beta^J} = \frac{\partial}{\partial r_\alpha^I} \left(\sum_{j \in J} \frac{\partial V}{\partial r_\beta^J(j)} \right) = \sum_{j \in J} \frac{\partial}{\partial r_\beta^J(j)} \left(\frac{\partial V}{\partial r_\alpha^I} \right) = \sum_{i \in I} \sum_{j \in J} \frac{\partial^2 V}{\partial r_\alpha^I(i) \partial r_\beta^J(j)}. \quad (\text{A.6})$$

$$\begin{aligned}
\frac{\partial^2 V}{\partial r_\alpha^I \partial p_\beta^J} &= \frac{\partial}{\partial r_\alpha^I} \left(\sum_{j \in J} \sum_{a=1}^3 \frac{\partial V}{\partial r_a^J(j)} \left[\frac{\partial \mathbf{S}^J}{\partial p_\beta^J} \mathbf{x}^J(j) \right]_a \right) \\
&= \sum_{j \in J} \sum_{a=1}^3 \frac{\partial}{\partial r_a^J(j)} \left(\frac{\partial V}{\partial r_\alpha^I} \right) \left[\frac{\partial \mathbf{S}^J}{\partial p_\beta^J} \mathbf{x}^J(j) \right]_a \\
&= \sum_{i \in I} \sum_{j \in J} \sum_{a=1}^3 \frac{\partial^2 V}{\partial r_\alpha^I(i) \partial r_a^J(j)} \left[\frac{\partial \mathbf{S}^J}{\partial p_\beta^J} \mathbf{x}^J(j) \right]_a.
\end{aligned} \tag{A.7}$$

$$\begin{aligned}
\frac{\partial^2 V}{\partial p_\alpha^I \partial p_\beta^J} &= \frac{\partial}{\partial p_\alpha^I} \left(\sum_{j \in J} \sum_{a=1}^3 \frac{\partial V}{\partial r_a^J(j)} \left[\frac{\partial \mathbf{S}^J}{\partial p_\beta^J} \mathbf{x}^J(j) \right]_a \right), \quad \text{for } I \neq J \\
&= \sum_{j \in J} \sum_{a=1}^3 \frac{\partial}{\partial r_a^J(j)} \left(\frac{\partial V}{\partial p_\alpha^I} \right) \left[\frac{\partial \mathbf{S}^J}{\partial p_\beta^J} \mathbf{x}^J(j) \right]_a \\
&= \sum_{i \in I} \sum_{j \in J} \sum_{a=1}^3 \sum_{b=1}^3 \frac{\partial^2 V}{\partial r_b^I(i) \partial r_a^J(j)} \left[\frac{\partial \mathbf{S}^I}{\partial p_\alpha^I} \mathbf{x}^I(i) \right]_b \left[\frac{\partial \mathbf{S}^J}{\partial p_\beta^J} \mathbf{x}^J(j) \right]_a.
\end{aligned} \tag{A.8}$$

$$\begin{aligned}
\frac{\partial^2 V}{\partial p_\alpha^I \partial p_\beta^I} &= \frac{\partial}{\partial p_\alpha^I} \left(\sum_{i \in I} \sum_{a=1}^3 \frac{\partial V}{\partial r_a^I(i)} \left[\frac{\partial \mathbf{S}^I}{\partial p_\beta^I} \mathbf{x}^I(i) \right]_a \right) \\
&= \sum_{i \in I} \sum_{a=1}^3 \frac{\partial}{\partial r_a^I(i)} \left(\frac{\partial V}{\partial p_\alpha^I} \right) \left[\frac{\partial \mathbf{S}^I}{\partial p_\beta^I} \mathbf{x}^I(i) \right]_a + \sum_{i \in I} \sum_{a=1}^3 \frac{\partial}{\partial r_a^I(i)} \left[\frac{\partial}{\partial p_\alpha^I} \left(\frac{\partial \mathbf{S}^I}{\partial p_\beta^I} \right) \mathbf{x}^I(i) \right]_a \\
&= \sum_{i_1 \in I} \sum_{i_2 \in I} \sum_{a=1}^3 \sum_{b=1}^3 \frac{\partial^2 V}{\partial r_b^I(i_1) \partial r_a^I(i_2)} \left[\frac{\partial \mathbf{S}^I}{\partial p_\alpha^I} \mathbf{x}^I(i_1) \right]_b \left[\frac{\partial \mathbf{S}^I}{\partial p_\beta^I} \mathbf{x}^I(i_2) \right]_a + \sum_{i \in I} \sum_{a=1}^3 \frac{\partial V}{\partial r_a^I(i)} \left[\frac{\partial^2 \mathbf{S}^I}{\partial p_\alpha^I \partial p_\beta^I} \mathbf{x}^I(i) \right]_a
\end{aligned} \tag{A.9}$$

4.5.2 Supporting tables and figures

Residue name	Three letter code	$\theta_i^{\max}/2\pi$		
		$\alpha\beta$	$\beta\gamma$	$\gamma\delta$
Alanine	ALA	1.0	-	-
Arginine	ARG	0.2	0.3	0.5
Asparagine	ASN	0.5	1.0	-
Aspartic acid	ASP	0.5	1.0	-
Cysteine	CYS	1.0	-	-
Glutamic acid	GLU	0.3	0.5	1.0
Glutamine	GLN	0.3	0.5	1.0
Glycine	GLY	-	-	-
Histidine	HIS	0.3	0.5	-
Isoleucine	ILE	0.5	1.0	-
Leucine	LEU	0.5	1.0	-
Lysine	LYS	0.2	0.3	0.5
Methionine	MET	0.5	0.6	-
Phenylalanine	PHE	0.3	0.5	-
Proline	PRO	-	-	-
Serine	SER	1.0	-	-
Threonine	THR	1.0	-	-
Tryptophan	TRP	0.3	0.4	-
Tyrosine	TYR	0.3	0.5	-
Valine	VAL	1.0	-	-

Table A1 The maximum rotation amplitudes θ_i^{\max} for the amino acid side chain groups used in the current work. $\alpha\beta$, $\beta\gamma$ and $\gamma\delta$ refer to rotation about the C_α - C_β , C_β - C_γ and C_γ - C_δ bonds respectively.

PDB ID	C_{α} -RMSD from 2INE (Å)		Note
	All	in 16 Å from ligand	
2INE	0.000	0.000	Complexed with Phenylacetic Acid, RMSD reference
2IQ0	0.094	0.085	Complexed with Hexanoic Acid
2IS7	0.137	0.145	Complexed with Dichlorophenylacetic Acid
2INZ	0.155	0.134	Complexed with 2-Hydroxyphenylacetic Acid
1AH0	0.497	0.248	Pig aldose reductase complexed with Sorbinil
1EL3	0.301	0.156	Complexed with IDD384 inhibitor
1EKO	0.474	0.290	Pig aldose reductase complexed with IDD384 inhibitor
1IEI	0.723	0.609	Complexed with Zenarestat
1MAR	0.482	0.495	Complexed with Zopolrestat

Table A2 C_{α} -RMSD between 2INE which I used as “St-1” and 8 aldose reductase crystal structures with different ligands bound. C_{α} -RMSD for each pair is calculated for the whole system and the residues within 16 Å of the ligand, respectively.

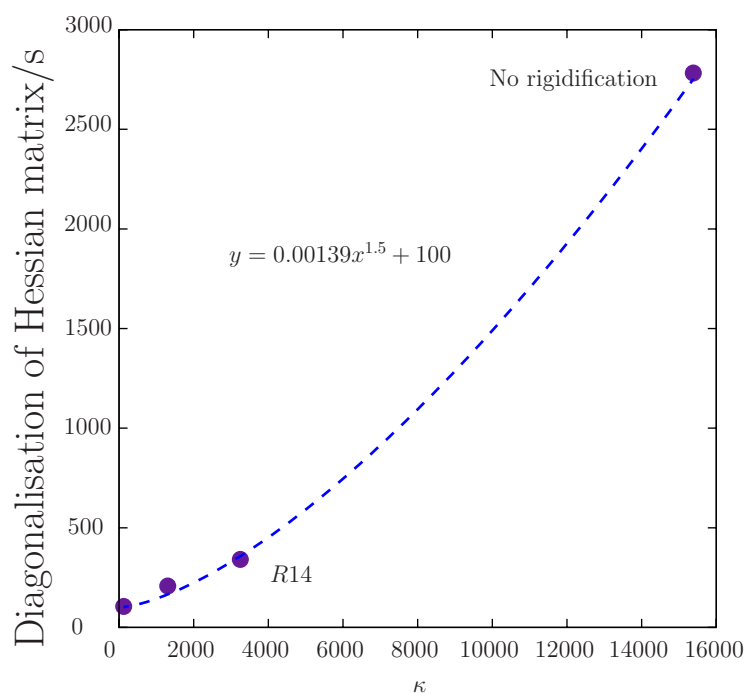


Figure A1 Average calculation time for diagonalisation of the Hessian matrix. $\kappa = 120$, 1299, 3246, corresponding to complexes with $R = 6$, 10, 14, and $\kappa = 15387$ (without any rigidification) are plotted. The fitting line scales as $\kappa^{1.5}$, described as blue dashed line. The calculations were performed on a 2.6GHz Xeon X5650 machine.

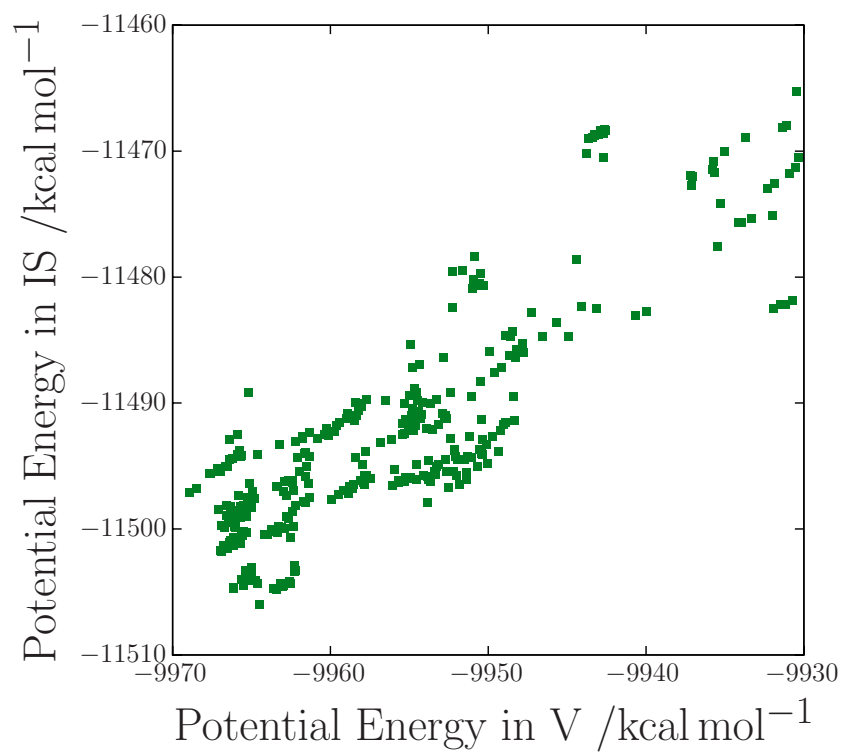


Figure A2 The potential energies in vacuum (V) and the corresponding recomputed energies in the implicit solvent (IS) are plotted. The Pearson correlation coefficient is 0.86.

References

- [1] Jorgensen, W. L., *Science*, **303**, 1813–1818 (2004)
- [2] Michel, J., Foloppe, N. and Essex, J. W., *Mol. Inf.*, **29**, 570–578 (2010)
- [3] Guitierrez-de-Teran, H., Aqvist, J., *Methods Mol. Biol.*, **819**, 305–323 (2012)
- [4] Johnson, R. R. et al., *Nano Lett.*, **9**, 537–541 (2009)
- [5] Ercolani, G., *J. Am. Chem. Soc.*, **125**, 16097–16103 (2003)
- [6] Cheng, T. et al., *AAPS J.*, **14**, 133–141 (2012)
- [7] Kitchen, D. B. et al., *Nat. Rev. Drug Discovery*, **3**, 935–949 (2004)
- [8] Pearlman, D. A. et al., *Comput. Phys. Commun.*, **91**, 1–41 (1995)
- [9] Lybrand, T., McCammon, J. A., Wipff, G., *Proc. Natl. Acad. Sci. U.S.A.*, **83**, 833–835 (1986)
- [10] Kollman, P. A., *Chem. Rev.*, **93**, 2395–2417 (1993)
- [11] Beveridge, D. L., Dicapua, F. M., *Annu. Rev. Biophys. Biophys. Chem.*, **18**, 431–492 (1989)
- [12] Reddy, M. R., Erion, M. D., *Curr. Pharm. Des.*, **11**, 283–294 (2005)
- [13] Deng, Y., Roux, B., *J. Phys. Chem. B*, **113**, 2234–2246 (2009)
- [14] Chen, P. C., Kuyucak, S., *Biophys. J.*, **100**, 2466–2474 (2011)
- [15] Park, S., Schulten, K., *J. Chem. Phys.*, **120**, 5946–5961 (2004)
- [16] Ytreberg, F. M., *J. Chem. Phys.*, **130**, 164906 (2009)
- [17] Brown, S. P., Muchmore, S. W., *J. Chem. Inf. Model.*, **46**, 999–1005 (2006)
- [18] Gouda, H. et al., *Biopolymers*, **68**, 16–34 (2003)
- [19] Pearlman, D. A., *J. Med. Chem.*, **48**, 7796–7807 (2005)
- [20] Srinivasan, J. et al., *J. Am. Chem. Soc.*, **120**, 9401–9409 (1998)
- [21] Aqvist, J., Medina, C., Samuelsson, J. E., *J. Protein Eng.*, **7**, 385–391 (1994)
- [22] Carlsson, J. et al., *J. Phys. Chem. B*, **110**, 12034–12041 (2006)
- [23] Jones-Hertzog, D., Jorgensen, W., *J. Med. Chem.*, **40**, 1539–1549 (1997)
- [24] Zhou, R. et al., *J. Phys. Chem. B*, **105**, 10388–10397 (2001)
- [25] Oostenbrink, C., van Gunsteren, W. F., *Proteins*, **54**, 237–246 (2004)
- [26] Wales, D. J., *Energy Landscapes*, Cambridge University Press, Cambridge, 364–433 (2003)
- [27] Wales, D. J., *Philos. Trans. R. Soc. London, Ser. A*, **363**, 357–377 (2005)
- [28] Strodel, B., Wales, D. J., *Chem. Phys. Lett.*, **466**, 105–115 (2008)
- [29] Mezey, P. G., *Potential Energy Hypersurfaces*, Elsevier, Amsterdam, 198–368 (1987)
- [30] Wales, D. J., *Mol. Phys.*, **78**, 151–171 (1993)
- [31] Wales, D. J., Doye, J. P. K., *J. Chem. Phys.*, **103**, 3061–3070 (1995)
- [32] Doye, J. P. K., Wales, D. J., *J. Chem. Phys.*, **102**, 9673–9688 (1995)

- [33] Calvo, F., Doye, J. P. K., Wales, D. J., *J. Chem. Phys.*, **115**, 9627–9636 (2001)
- [34] Calvo, F., Doye, J. P. K., Wales, D. J., *J. Chem. Phys.*, **114**, 7312–7329 (2001)
- [35] Chen, W., Chang, C. E., Gilson, M. K., *Biophys. J.*, **87**, 3035–3049 (2004)
- [36] Doye, J. P. K., Miller, M. A., Wales, D. J., *J. Chem. Phys.*, **110**, 6896–6906 (1999)
- [37] Wales, D. J., Bogdan, T. V., *J. Phys. Chem. B*, **110**, 20765–20776 (2006)
- [38] Sharapov, V. A., Meluzzi, D., Mandelshtam, V. A., *Phys. Rev. Lett.*, **98**, 105701 (2007)
- [39] Sharapov, V. A., Mandelshtam, V. A., *J. Phys. Chem. A*, **111**, 10284–10291 (2007)
- [40] Stillinger, F. H., Weber, T. A., *Science*, **225**, 983–989 (1984)
- [41] Wales, D. J., Doye, J. P. K., *J. Chem. Phys.*, **119**, 12409–12416 (2003)
- [42] Bogdan, T. V., Wales, D. J., Calvo, F., *J. Chem. Phys.*, **124**, 044102 (2006)
- [43] Chen, W. et al., *J. Chem. Theory Comput.*, **6**, 3540–3557 (2010)
- [44] Kusumaatmaja, H., Whittleston, C. S., Wales, D. J., *J. Chem. Theory Comput.*, **8**, 5159–5165 (2012)
- [45] Petrash, J. M., *Cell. Mol. Life Sci.*, **61**, 737–749 (2004)
- [46] Gilson, M. K., *Biophys. J.*, **72**, 1047–1069 (1997)
- [47] Zhou, H., Gilson, M. K., *Chem. Rev.*, **109**, 4092–4107 (2009)
- [48] Emilio, G., Ronald, M. L., *Recent theoretical and computational advances for modeling protein-ligand binding affinities*, Academic Press, **85**, 27–80 (2011)
- [49] Amar, F. G., Berry, R. S., *J. Chem. Phys.*, **85**, 5943–5954 (1986)
- [50] Chakrabarti, D., Wales, D. J., *Phys. Chem. Chem. Phys.*, **11**, 1970–1976 (2009)
- [51] Pohorille, A., *J. Chem. Phys.*, **87**, 6070–6077 (1987)
- [52] Wales, D. J., Ohmine, I., *J. Chem. Phys.*, **98**, 7257–7268 (1993)
- [53] Li, Z. Q., Scheraga, H. A., *Proc. Natl. Acad. Sci. U. S. A.*, **84**, 6611–6615 (1987)
- [54] Wales, D. J., Doye, J. P. K., *J. Phys. Chem. A*, **101**, 5111–5116 (1997)
- [55] Wales, D. J., Scheraga, H. A., *Science*, **285**, 1368–1372 (1999)
- [56] Whittleston, C., *PhD thesis*, University of Cambridge (2011)
- [57] Press, W. et al., *Numerical Recipes*, Cambridge University Press, Cambridge, 487–555 (1986)
- [58] Strodel, B. et al., *J. Am. Chem. Soc.*, **132**, 13300–13312 (2010)
- [59] Geyer, G. J., *Stat. Sci.*, **7**, 437 (1992)
- [60] Hukushima, K., Nemoto, K., *J. Phys. Soc. Jpn.*, **65**, 1604–1608 (1996)
- [61] *The PyMOL Molecular Graphics System*, Version 1.2r3pre, Schrodinger, LLC.
- [62] Goldstein, H., *Classical mechanics*, Addison-Wesley, Reading, Massachusetts, 128–187 (1980)
- [63] Chakrabarti, D., Wales, D. J., *Phys. Chem. Chem. Phys.*, **11**, 1970–1976 (2009)

- [64] Brownlee, J. M. et al., *Bioorg. Chem.*, **34**, 424–444 (2006)
- [65] Holmberg, N., Ryde, U., Bulow, L., *Protein Eng.*, **12**, 851–856 (1999)
- [66] Wang, J., Cieplak, P., Kollman, P. A., *J. Comput. Chem.*, **21**, 10491074 (2000)
- [67] Wang, J. et al., *J. Mol. Graphics Modell.*, **25**, 247–260 (2006)
- [68] Bayly, C., *J. Phys. Chem.*, **97**, 10269–10280 (1993)
- [69] Cornell, W., *J. Am. Chem. Soc.*, **115**, 9620–9631 (1993)
- [70] Schmidt, M. W., *J. Comput. Chem.*, **14**, 1347–1363 (1993)
- [71] Wales, D. J., *GMIN*: a program for basin-hopping global optimisation, basin-sampling, and parallel tempering.
- [72] Wales, D. J., *OPTIM*: a program for optimizing geometries and calculating reaction pathways.
- [73] Onufriev, A., Bashford, D., Case, D. A., *Proteins*, **55**, 383–394 (2004)
- [74] Sigalov, G., Fenley, A., Onufriev, A., *J. Chem. Phys.*, **124**, 124902 (2006)
- [75] Chang, M. J. P. C. E., Gilson, M. K., *J. Phys. Chem. B*, **107**, 1048–1055 (2003)
- [76] Temelso, K. A. A. B., Shields, G. C., *J. Phys. Chem. A*, **115**, 12034–12046 (2011)
- [77] Hawkins, R. J., McLeish, T. C. B., *J. R. Soc., Interface*, **3**, 125–138 (2006)
- [78] Hawkins, R. J., McLeish, T. C. B., *Phys. Rev. Lett.*, **93**, 098104 (2004)

Chapter 5

Summary

The theoretical studies on complex chemical systems were presented in three chapters of this thesis. First, the basic molecular mechanism of the homogeneous ice melting was analyzed. It was shown that a separated defect pair (I-V), which is very small in size, is the key initial step of ice melting. This defect pair is accidentally created in the process of a thermally distorted HB network segment (accumulation of 5+7 defects) trying to go back to the crystalline order. The defect pair is entangled, i.e., is long-lived and persistently agitates the HB network of a crystalline ice. It gradually destroys the resilient HB network and leads the system into the total melting. Secondly, I investigated how water mixes with methanol in the molecular level at different water concentrations (X). The C K-edge XAS spectral intensity corresponding to the local structure around methyl group changes at X=0.3 and 0.7. MD simulation shows the percolation of the HB network among only water molecules and the mixture is broken above X=0.3 and 0.7, respectively. Thus, I clarified that the local structure of the mixture changes non-linearly governed by the global HB network behavior. In addition, I developed a conformational factorization method to improve the sampling efficiency for estimating the free energy of ligand binding for macro biomolecules, which is likely to facilitate screening calculations associated with drug design.

In this thesis, I mainly focused on HBs network feature and showed that the HBs network is very resilient to change its phase, even the initial conformation is in the metastable state. Accordingly, at moderate temperature, the solid-liquid transition (melting) is hardly achieved without the entangled mechanism to prevent the recrystallization and to promote the phase transition. The resilient behavior of HB network attributes to the character of local structure that each water molecule attempts to form tetra-coordinated HBs with its neighbors for potential energy stabilization. This HBs network character is still available in the aqueous mixture with the smallest amphiphile molecule (methanol) and the local structure changes around methyl group against the aqueous concentration corresponds to the global HBs network behavior composed of water and hydroxyl group. Thus, I have performed the MD simulations and carefully analyzed the detail structure and dynamics of HBs network in the process of ice melting and in the aqueous mixture, then revealed that the local HBs structure strongly correlates with the global HBs network. In addition, I developed the new calculation approach for the ligand-binding free energy, which must be useful to evaluate water contribution to the ligand binding process.

In the future work, the molecular mechanism, I found in the homogeneous melting of ice, must be explored in the melting processes of other atomic/molecular solids. The role of the interstitial defects and the entanglement of the network can be essential for not only ice but also other solids. The knowledge on the detail molecular mechanism of the phase transitions will be

useful to understand how the complex molecular systems such as proteins and many biological molecular system can exhibit functionalities; for example, a biomolecular system often undergoes sequential unique reactions with a very small energy (ratchet-type mechanism), and a drastic conformational change of ligand-protein system is initiated by a small ligand molecule or a small local signal (amplification).

Acknowledgements

Completion of this doctoral dissertation was possible with the support of several people. I would like to express my sincere gratitude to all of them. First of all, I am extremely grateful to my research guide, Professor Iwao Ohmine, Director General of Institute for Molecular Science and Vice President of National Institutes of Natural Sciences, for his valuable guidance, scholarly inputs and consistent enthusiasm I received throughout the research work. I would never have come back to academic without his help in the first place. I am also grateful to Professor Nobuhiro Kosugi, Deputy Director General of Institute for Molecular Science, for his academic supports, patience and the facilities provided to carry out the research work at the Institute. Collaboration with experiments must extend my research field. I am sincerely grateful to Dr. Masakazu Mastumoto, Associate Professor at Okayama University, for his continual support and encouragements. His epochal research techniques and ideas always surprised me. Especially I learned how to describe or express “Structure”. My sincere thanks also goes to Professor David Wales at Cambridge University for giving me the opportunity to visit Cambridge and telling me how to treat bio-molecules. I thank the group members in Kosugi Lab. and Wales Lab. for their advices and discussions. Finally, I am very much indebted to my family, my wife Chiemi, sons Rintaro and Kyuji, parents Hidenori and Ikumi, grandparents Teruhisa, Nawoe, Kiyoshi and Kazue, brothers Kazuhiro and Toshikuni, parents in law Shigeru and Yasue, and all relatives, for their spiritual supports throughout my life. I managed to complete this doctoral thesis here by a lot of great help.