

The repetitive local sampling and the local distribution theory

Pu Tian^{*,†,‡}

[†]*School of Life Sciences, Jilin University, Changchun, China 130012*

[‡]*School of Artificial Intelligence, Jilin University, Changchun, China 130012*

E-mail: tianpu@jlu.edu.cn

Phone: +86 (0)431 85155287

Abstract

Molecular simulation is a mature and versatile tool set widely utilized in many subjects. However, its methodology development has been struggling with a tradeoff between accuracy and speed, significant improvement of both is necessary to reliably substitute many expensive and laborious experiments in molecular biology and nanotechnology. Previously, the ubiquitous severe wasting of computational resources in molecular simulations due to repetitive local sampling was raised, and the local free energy landscape approach was proposed to address it. The core idea is to first learn local distributions, and followed by dynamic assembly of which to infer global joint distribution of a target molecular system. When compared with conventional explicit solvent molecular dynamics simulations, a simple and approximate implementation of this theory in protein structural refinement harvested acceleration of about six orders of magnitude without loss of accuracy. While this initial test revealed tremendous benefits for addressing repetitive local sampling, there are some implicit assumptions need to be articulated. Here, I present a more thorough discussion of repetitive local sampling; potential options for learning local distributions; a more general formulation with potential extension

to simulation of near equilibrium molecular systems; the prospect of developing computation driven molecular science; the connection to mainstream residue pair distance distribution based protein structure studies; and the fundamental difference of the averaging strategy from potential of mean force. This more general development is termed the local distribution theory to release the limitation of strict thermodynamic equilibrium in its potential wide application in general soft condensed molecular systems.

Introduction

Molecular simulation has been utilized in a wide variety of disciplines, including but not limited to chemistry, physics, biology and materials science. Its increasing importance is clearly demonstrated by steady growth of relevant publications as shown in Fig. 1. However, atomistic molecular dynamics (MD) simulations, while being effective in revealing underlying atomic mechanisms for many molecular processes, are extremely computationally intensive.^{1,2} Historically, scientists have developed two lines of algorithms to accelerate molecular simulations, with one being coarse graining³⁻¹² and the other being enhanced sampling.¹³⁻¹⁶ Realizing that there is severe wasting of computational resources due to repetitive local sampling (RLS) in all molecular simulations, the local free energy landscape (LFEL) approach was proposed to eliminate such wasting, and its effectiveness was subsequently demonstrated in an approximate implementation in protein structural refinement.¹⁷ In the initial testing of this new theory, LFEL for amino acid packing in proteins was constructed based on a simple neural network implementation of generalized solvation free energy (GSFE) theory.¹⁸ Further, a computational graph was established through combination of automatic differentiation, coordinate transformation and LFEL cached in trained neural networks. This computational graph was successfully utilized to achieve the only end-to-end and the most efficient protein structural refinement pipeline¹⁷ up to date. The connection among coarse graining, enhanced sampling and LFEL as various forms of applying “dividing and conquering” and “caching” principle in molecular modeling was summarized previously.¹⁹ Like all present protein structure prediction, design and refinement studies,²⁰⁻²⁹ there is an implicit and extremely crude

assumption that all high resolution experimental structures were solved under similar environmental (thermodynamic) conditions. Alternatively, differences in thermodynamic and environmental conditions are deemed not important for all high resolution structural data utilized to train models. Such assumptions neglect rich conformational redistribution of many functional proteins upon environmental stimuli. For example, while most protein structures stay intact when temperature decrease from physiological condition, some proteins denature upon sufficiently low temperature.^{30,31} Many channel proteins may switch between open and closed states upon change of concentration of specific ions³² and mechanical stress^{33,34}. To realize modeling of these functional and dynamic property of biomolecular systems, further development of algorithms that is capable of accounting for diverse environmental stimuli is apparently necessary. Additionally, the LFEL approach as it stands only applies to equilibrium conditions. Here, I explicitly articulate these issues, develop a more general form of the LFEL idea and term it the local distribution theory (LDT). Meanwhile, more concrete discussions of RLS, more options for fitting local distributions, extension of LDT to near-equilibrium scenarios, connection of LDT to present protein structural studies based on deep neural networks, and the difference of LDT from conventional molecular simulation framework based on potential of mean force are presented. It is hoped that this work will intrigue more interest in further development of LDT in general chemical and biomolecular systems, and facilitate advancement of computation driven molecular science.

Repetitive local sampling

In molecular simulations, we have a long history of utilizing RLS in analysis of MD trajectories. For example, when computing a pair distribution function $g(r)$ between oxygen atoms of water molecules, instead of tracking a specific pair of water molecules or water molecules within a given small space and binning distances of oxygen atom pairs, statistics is usually accumulated by counting all pairs of water molecules within a distance of half simulation box to obtain a more smooth curve. Similarly, in calculation of dynamic property such as mean squared displacement, in stead

of focusing on one specific particle and trace its motion with time, displacement of many particles are traced and results are averaged. Such tricks are routinely utilized in various analyses of molecular simulation trajectories. The basis of these manipulations is the belief that all molecules of the same chemical identity and composition are indistinguishable, and ensemble average converges to time average for ergodic systems. All above practices clearly demonstrate that we have been carrying out RLS in essentially all our simulations, except not carefully thinking about its potential utility in saving computational resources in the simulation/sampling stage. The reality that RLS consumes majority of computational resources in regular molecular simulations was raised previously^{17,19} without sufficiently detailed discussions. Some typical examples of RLS in various simulation and/or modeling applications are discussed below.

RLS within a single simulation task

Fig. 2a represents a snapshot for a simulation of aqueous solution comprising a few different types of ions and water molecules, with gas-liquid and liquid-solid interfaces under given thermodynamic conditions. In a production run of this simulation, we may choose to focus our attention on a spherical space A . As the simulation goes on, many different configurations with various number of ions and water molecules within this space will be observed. When the simulation is sufficiently long, a converged LFEL will be obtained. This LFEL is a complex high dimensional distribution that gives correct statistical weight for each thermally accessible structural ensemble (or free energy local minimum) on the one hand, and all possible transition paths connecting these minima with respective statistical significance on the other hand. The exactly same LFEL would have been obtained if another bulk spherical space B with the same volume was taken. As a matter of fact, the exactly same LFEL would have been obtained for all possible bulk spherical spaces with the same volume. However, for each such separate local space, significant computational resource was consumed to obtain the exactly same LFEL, resulting in tremendous wasting of computational resources! This is a typical case of RLS within the same simulation task. Beyond the illustration in Fig. 2a, there are other less obvious forms of RLS. For example, in protein structure

prediction, design and refinement with implicit representation of aqueous solution, each residue in a chain has more or less unique surroundings and no direct RLS over different local spaces seems existing. However, in these tasks, each residue experiences many rounds of adjustment or repacking. Sampled collisions, favorable and unfavorable configurations from each round is partially or completely discarded and performed on the fly in the next round, engendering significant RLS.

RLS across different simulations

Much more computational resource are consumed by RLS across different simulation tasks. Imagine how many times simulations of local packing for water molecules of each popular water force fields have been carried out by thousands of scientists globally! Similarly, molecular packing of amino acids surrounding each of 20 natural amino acids have been carried out numerous times by computational structural bioinformaticians around the world. Such RLS is ubiquitous for simulations of the overwhelming majority of molecular systems in chemistry, biology and materials science.

The generalized solvation free energy perspective

While local spaces near various interface certainly have LFELs different from that of bulk, there are regularities that can be learned as well. Such RLS may be effectively described from a slightly different perspective according to the GSFE theory as shown in Fig. 2b. In GSFE theory, each comprising unit of a molecular system is on the one hand a solute unit solvated by its surrounding units, and on the other hand a comprising solvent unit for each of its surrounding units. As all units with the same chemical identity/structure are indistinguishable, so should be LFEL of their local solvent under given thermodynamic conditions if a simulation trajectory is sufficiently long. When our focus is on LFEL surrounding a given central unit, different scenarios of interfaces are simply different solvent configurations with corresponding statistical weights and no special treatment is required anymore. More specifically, for a water molecule absorbed on wall of a tube filled with water, its solvent units include both water molecules and molecules belong to

the wall surrounding it. To eliminate difficulty of defining interfaces at molecular scales is the very initial motivation for development of the GSFE theory. Additionally, defining local spaces with local coordinates originated from individual molecule is a convenient, efficient and natural choice with two advantages. Firstly, it reduces data requirement and improves accuracy during training/learning of local distributions, and secondly, it facilitates assembly by eliminating the uncertainty of selecting from infinite possible origins for local spaces during inference for global joint distribution (GJD) of a target molecular system.

Potential benefits of utilizing local distributions to eliminate RLS

Sufficient sampling of complex molecular systems has long been our pursuit in simulation studies. The very fact that we almost always collect statistics from different local spaces and utilize indistinguishable property of molecules for better statistics indicates that we rarely achieve sufficient sampling for a given small space or surrounding of a given single molecule. Therefore, it is likely that more accurate global correlations would have been obtained if sufficient statistics was available for all local regions. Since in construction of GJD by assembly of LFEL, local spaces surrounding each particle is driven by a sufficiently sampled LFEL, consequently facilitating sufficient sampling of local regions. Therefore, the ability to cache and utilize LFEL properly would not only tremendously reduce the need of computational resources, but also potentially improve accuracy due to effectively more sufficient “local sampling”. This is in strong contrast to decades of trade-off in molecular simulations that improved efficiency being always accompanied more or less by reduced accuracy, and increased accuracy being always accompanied by more or less reduction of efficiency! When compared with conventional molecular mechanical force fields^{35–38} or knowledge based potentials,^{39–41} the ability of accounting for many-body correlations is another advantage of LFEL that is likely to contribute to improved accuracy.

The local distribution theory

It is well understood that the folding process and conformational distributions for a given protein depend upon both its sequence and environmental conditions. However, due to lack to data, in both establishment of traditional knowledge based potentials^{39–41} and deep learning studies^{21,22} of protein folding, design and structural refinement, it is widely assumed that all experimental structural data may be deemed as obtained under similar conditions, and details of which may be safely ignored in such tasks. Such simplification was similarly utilized in implementing the LFEL approach in protein structure refinement¹⁷ with focus being on coordinates without attending to thermodynamic and solvent conditions. Should detailed modeling of the variation of interested molecular systems under different environmental and/or thermodynamic conditions be desired, inclusion of these variables in the formulation was essential. Here, previous simplified formulation is extended to deal with such scenarios. Denote environmental and thermodynamic variables (e.g. temperature, pressure, concentrations of relevant molecular species, special restraints) as $\Phi = (\phi_1, \phi_2, \dots, \phi_k)$, molecular coordinates as $X = (x_1, x_2, \dots, x_n)$ and local regions of molecular systems as $R = (r_1, r_2, \dots, r_m)$ ($m \leq n$, $m = n$ is preferred), the GJD may be expressed by local distributions $q(\Phi, r_i)$ and their correlations as:

$$\begin{aligned} P(\Phi, X) &= Q(\Phi, R) \\ &= \frac{Q(\Phi, R)}{\prod_{i=1}^m q(\Phi, r_i)} \prod_{i=1}^m q(\Phi, r_i) \end{aligned} \quad (1)$$

It is important to note that each r_i ($i = 1, 2, \dots, m$) represents a dynamic collection of molecular coordinates for the i th specified region $r_i = \mathcal{M}(x_{i1}, x_{i2}, \dots, x_{il})$, with \mathcal{M} being a translation, rotation and permutation invariant transformation matrix from the global to local coordinate system. Both number l and identity of involved units may change for each local region with propagating trajectories. When ($m = n$) or m is close to n , since each local region contains dozens of or more particles, overlapping among such regions are extensive. Local distributions are essentially LFEL for equilibrium systems. The fraction term $\frac{Q(\Phi, R)}{\prod_{i=1}^m q(\Phi, r_i)}$ includes all complex global correla-

tions among various local regions $r_i (i = 1, 2, \dots, m)$ and is denoted the global correlation factor (GCF) previously.¹⁹ The product term (hereafter “local term”) $\prod_{i=1}^m q(\Phi, r_i)$ is simply to treat all local regions as if they were independent. If the GCF was ignored, then overlapping parts of different r_i may have distinct states. In reality, regardless of how many different local regions a molecule x_i participates, it has a unique physical state at any given instant. So all possible configurations with contradicting molecular states for any molecule participating different local regions have probability density zero. Such correction and additional modification of probability density is achieved by the GCF term. However, direct calculation of GCF is intractable for any realistic complex molecular system. Therefore, equation 1 is not directly useful for understanding and predicting behavior of molecular systems. How to approximately and effectively utilize this equation in practice is an open problem, and likely with multiple potential approximate solutions. The optimal approximation might well have some extent of molecular system specificity.

Probability density (free energy in equilibrium) of a specific configuration may be decomposed into three approximately independent contributions. The first is the short range contribution (F_{SR}) that measures the extent of structural stability/compatibility within each local region and is quantified by the local term in equation 1. The second contribution is from mediated interactions (F_{MED} , Fig. 3ab) that measures the extent of compatibility among all overlapping local regions, and the third contribution measures direct long range (F_{LR} , Fig. 3b) compatibility within the whole molecular system. Both the second and the third contributions are contained in the GCF term. With the assumption that mediated interactions are independent from long-range interactions, the GCF may be approximately split into F_{MED} and F_{LR} as shown below.

$$\frac{Q(\Phi, R)}{\prod_{i=1}^m q(\Phi, r_i)} \approx \exp(-\sum F_{MED}(\Phi, R)) \exp(-\sum F_{LR}(\Phi, R)) \quad (2)$$

The summations are over all mediated and long-range interactions in the given configuration R . In practical computation, separation of F_{SR} and F_{MED} is challenging on the one hand and inefficient on the other hand. In the previous implementation¹⁷ F_{SR} and F_{MED} were merged. Specifically, As

shown in Fig 3b, at any given instant, a molecule (particle) in the system experiences free energy driving force additively from local distributions centered on each of its directly interacting neighbors within a preset cutoff. This is in strong contrast to regular MD simulations in which a particle experience direct forces from its directly interacting neighbors. While F_{LR} was not accounted for previously, it may be added in for each particle in each or every few propagation step(s). So in equation 1, local interactions are separated from the GCF, which may be approximately decomposed into mediated and long range interactions. However, local and mediated interactions were computed together in the previous implementation. This choice is somewhat counter intuitive but is feasible and efficient. Since an analytically clean mathematical factorization of the GCF is not available, it is likely that the above approximation is just one of many possible ways to realize practical computation. Distinct molecular systems may have different correlation characteristics and the optimal approximation is likely to be system specific. Nonetheless, the overall idea is quite clear, that is to first train local distributions, which are subsequently to be assembled to compose the GJD according to suitable approximation of the equation 1. The core idea of the LDT is to use local distributions to eliminate RLS.

Combining equations 1 and 2, we have the following equation:

$$Q(\Phi, R) \approx \prod_{i=1}^m q(\Phi, r_i) \exp(-\sum F_{MED}(\Phi, R)) \exp(-\sum F_{LR}(\Phi, R)) \quad (3)$$

with $q(\Phi, r_i)$ being neural networks parameterized by W , $\exp(-\sum F_{MED}(\Phi, R))$ enforce mediated interactions via sampling constraints, and $\exp(-\sum F_{LR}(\Phi, R))$ being calculated through a selected direct long-range interaction calculation method independent of W . In the training stage, local distributions are learned by optimizing parameter set W with a given data set $D = \{X_1, X_2, \dots, X_d\}$

$$W^* = \operatorname{argmax}_W \prod_{j=1}^d Q_j(\Phi, R(X_j)) \quad (4)$$

with W^* representing the optimal parameter set and subscript j indicating the j th record with coordinate X_j in the data set. In the free energy optimization stage, W^* is fixed, with Q^* indicating

local distributions in which being parameterized by W^* , we have:

$$X^* = \underset{X}{\operatorname{argmax}} Q^*(\Phi, R(X)) \quad (5)$$

as each unit x participates in many local distributions, this optimization process essentially is competition of each unit for its best position given its solvent environment, each comprising unit of which is simultaneously seek respective best position, thus drive the propagation of a target molecular system.

In a proper implementation of LDT, a target molecular system may be propagated similarly as in the case of MD simulations except for the two differences. The first difference is that empirical potentials driving MD is replaced by approximate GJD assembled from LDTs. The second is that a learning rate α_a , which is implicitly related to temperature, needs to be given. It is important to note that LDTs are utilized to replace RLS, not global sampling. To accelerate global sampling of a given molecular system, the propagation may be carried out in different temperatures other than the one corresponding to the training data. Methodologies such as simulated annealing⁴² may be realized just as in regular MD or MC simulations simply by assign a proper scheme of temperature cycles specified by corresponding gaussian noise term with variance α_b . In practice, α_a and α_b need not be identical in the following Langevin equation:

$$X_{t+1} = X_t - \alpha_a \frac{\partial(\sum F_{SR} + \sum F_{MED} + \sum F_{LR})}{\partial X} + \epsilon, \epsilon \sim \mathcal{N}(0, \alpha_b) \quad (6)$$

Challenges and options for fitting local distributions

Training/learning of local terms is by no means trivial. In reality, strictly normalized local distributions is beyond reach and we may approximate them by complex high dimensional unnormalized potential functions. The direct consequence of lacking normalization is that resulting free energy unit is arbitrary and is different for different molecular systems. When direct long range interactions are to be added, or comparison of results among different molecular systems are essential,

this uncertainty has to be resolved. If long-range interactions with fixed unit may be calculated accurately, then it can serve as a unit-defining quantity among different molecular systems.

Construction of local distributions is essentially a density estimation problem in high dimensional space. Firstly, each local region need to be represented mathematically in a translation, rotation and permutation invariant way for its probability density to be effectively fit. Such processing of molecular coordinates is accomplished by descriptor functions, which have accompanied development of neural network force fields,^{43,44} and are quite well understood. One possible way of defining a local region is to utilize the position of an given particle as the origin for the local coordinates, so $r_i = (x_{i-c}, y_{i-s})$, with x_{i-c} being the origin of the local coordinates defined by a given unit and y_{i-s} being the coordinates of all surrounding molecules within a preset cutoff. It is important to note that the number of molecules may fluctuate and so is the dimensionality of y_{i-s} , and padding is a feasible way to address it. The distribution of a local region within a molecular system under given environmental conditions Φ may be decomposed into a local prior $q(\Phi, y_{i-s})$ and a local likelihood $q(\Phi, x_{i-c}|\Phi, y_{i-s})$ as shown below:

$$\begin{aligned} q(\Phi, r_i) &= q(\Phi, x_{i-c}, y_{i-s}) \\ &= q(\Phi, x_{i-c}|\Phi, y_{i-s})q(\Phi, y_{i-s}) \end{aligned} \tag{7}$$

The likelihood term measures extent of match between the particle at the origin (x_{i-c}) and its surroundings. The prior term represent structural stability of the surrounding under given environmental conditions. In the protein structure refinement implementation,¹⁷ identities of the central amino acids were utilized as labels to train a simple neural network representing likelihood terms, and prior terms were approximated with simple weights. This strategy is likely to be not very useful for general molecular systems. For example, in a typical molecular system of dilute aqueous solution, the fraction of water molecules is the overwhelming majority. Training with identity will face extremely unbalanced data and important differences among minority molecular/ionic species are likely to be lost. To improve fitting of local distributions, accurate description of both

likelihood and prior terms are essential.

Like any density estimation application, fitting of local distributions may be carried out directly without decomposing into likelihood and prior terms. As a matter of fact, density estimation problem is of fundamental importance in both statistics and machine learning. Not surprisingly, many neural network architectures have been developed to tackle density estimation in high dimensional space where conventional methods (e.g. kernel density estimators⁴⁵) are not effective. The most widely utilized two types are autoregressive models⁴⁶ and normalizing flows.^{47,48} The former decompose a target joint density into product of conditional densities, which are modeled by parametric densities (e.g. mixture of gaussians) with trainable parameters. The later utilizing invertible neural network architectures to realize a direct quantitative map from a known density (e.g. uniform or gaussian) to the target density space. Establishment of proper correlations among different parametric densities is a highly challenging task for autoregressive models. The invertibility requirement in normalizing flow methodology imposes heavy restrictions on neural network architecture and hence its representation power. One outstanding application example of normalizing flow in modeling molecular system is the Boltzmann generator.⁴⁹ However, application of Boltzmann generator in complex molecular system remain to be tested. The fundamental difference between Boltzmann generator and LDT is that the former aims to directly model GJD for target molecular systems while the later decompose the problem into fitting and assembly of local distributions. Therefore RLS across different tasks is not addressed by Boltzmann generator, which as a results loses transferability of computed results among different molecular systems. A recent more general approach, Roundtrip,⁵⁰ was proposed to overcome weakness of these two density estimation methodology. However, it takes an expensive sampling step to finalize the density estimation. Each available class of methods has its pros and cons, and no theory is available for selection of proper density estimation methodology presently. It might well be that better methods will arise in future. For fitting local distributions in specific complex molecular system, many tests are likely necessary to construct a proper neural network model. Different molecular systems may have distinct structural distributions and case by case exploration is probably necessary to achieve

high accuracy.

Energy based models (EBM)^{51,52} are good candidates for fitting local distributions, either as a whole or when decomposed into prior and likelihood terms. In EBM, an energy is trained to be associated with a given configuration, thus eliminating the need of a normalization constant, which is a core challenge in fitting local distributions. Present tests of EBMs are mainly in conventional machine learning application scenarios such as computer vision or natural language processing.^{53–56} Density distributions for such systems are quite different from complex molecular systems of condensed matter. Since LDT is a new development, significant effort is necessary to search for both proper loss functions, neural network architectures, optimization algorithms and their combinations for EBM to facilitate fitting local distributions in our interested molecular systems.

While neural networks have been black boxes with exceptional fitting capability up to date, and have been utilized with a wide variety of architectures. Efforts are undergoing for building white box neural networks.⁵⁷ To realize more physically interpretable and mathematically elegant fitting of local distributions transparently is certainly an attractive potential direction to explore.

Connection to conventional AI driven protein structure studies

Contact map has played a critical role in development of protein structure prediction.²⁹ Earlier contact was a simple binary assignment (contact or not) defined by a cutoff distance based mostly on C_β atoms,²⁹ later on it evolved into residue pair distance distributions (RPDD).^{20,24,25,27} Significant effort has been invested in investigating impact of various input information and neural network architectures on RPDD prediction with great progress in understanding. As the only known fully end-to-end and the most efficient protein structure refinement and dynamic simulation pipeline, GSFE-refinement¹⁷ has a distinct overall pipeline from RPDD based algorithms of protein structure prediction/refinement. With the common goal of describing protein structures, these seemingly very different procedures have to be somehow connected. Fundamentally, all methodologies targeting protein structures reflect their underlying free energy landscape from cer-

tain perspective. In GSFE-refinement, the GJD assembled from local distributions (or LFEL) lacks direct long-range correlations beyond spatial range of mediated interactions (Fig. 3) as the method stands now. Certainly, addition of long-range correlations is feasible as already discussed above, and is in fact one important task in our future development plan. Sequence information is limited to the target protein itself in contrast to RPDD based methods, where multiple sequence alignment information is included as a critical part of input. In AlphaFold,²⁰ AlphaFold2⁵⁸ and many other RPDD based studies,^{21,22,24–29,59,60} the core information obtained is explicit protein (family) specific RPDD, which are in fact marginalization of the GJD after integrating away all other variables except the distance between the concerning residues. While marginalization in general is an extremely difficulty task in high dimensional space, it is trivial for an approximate GJD represented by a trajectory of configurations with heavy statistical weights confined within the corresponding manifold. Complex neural networks in RPDD based methods essentially realize a fitting from input information (protein sequence and multiple sequence alignment) to these marginal distributions without explicit construction of the GJD, approximation of which is the very goal of LDT based methods/models. As shown in Fig. 4, mapping from GJD to RPDD is readily achievable through marginalization. It is important to note that it takes some number of propagation steps (depending upon ruggedness of the underlying FEL) to obtain approximate GJD of sufficient accuracy assuming the underlying local distributions are sufficiently accurate. Marginalization is a deterministic procedure with significant loss of information, specifically correlations among different RPDD. Conversely, with RPDD, one may in principle construct GJD with sufficient sampling and optimization with necessary restraints. However, since correlations among different RPDD are absent, resulting GJD is highly dependent upon parameters and algorithms utilized in the corresponding reconstruction process. Present mainstream AI-based protein structural prediction/refinement neural networks implicitly cache some projections of local distributions and rules for assembling them into RPDD, each comes with its own loss of information that is hard to retrieve. LDT theory aims to first directly and explicitly learn local distributions, which are subsequently dynamically assembled to construct the most comprehensive GJD. LDT thus has the full potential to perform dynamic

modeling of relevant molecular processes as long as local distributions were fit for corresponding conditions. However, extending GSFE-refinement for accurately modeling dynamic protein folding is certainly not trivial as data on intermediate states are scarce presently. Nevertheless, LDT is a general theory applicable to any soft condense matter as long as fitting of corresponding local distributions is accomplished.

Potential extension to near equilibrium scenarios

At molecular scale, temperature, pressure and concentration of comprising molecules have significant fluctuations. In conventional MD simulations, temperature and pressure are usually controlled by various thermostats and barostats⁶¹ with equilibrium assumption. If we have a heterogeneous cell being heated at one side, specifying temperature and pressure at different locations within it is a challenge. It might well be that both temperature and pressure are heterogeneous in a live cell (sometimes or always) and we just have no proper way of measuring. To specify temperature and pressure with thermostats and barostats is difficult in such scenarios since we have no information on heterogeneous temperature in the first place. The probabilistic description of both molecular coordinates and thermodynamic/environmental variables can be of great utility. To apply LDT for these problems, we need to assume that target molecular systems are near-equilibrium. More specifically, all local distributions in a target molecular system are well approximated by local distributions trained from equilibrium data despite the global molecular system is off equilibrium (e.g. having temperature/pressure gradient). In such scenarios, we need thermodynamic variables to be associated with each local distribution. If the number of local regions was defined as the same as number of molecules/particles, we would have a set of relevant variables associated with each particle $\Phi_i = (\phi_{i1}, \phi_{i2}, \dots, \phi_{ik})$ and denote the environmental conditions as $\Phi = (\Phi_1, \Phi_2, \dots, \Phi_n)$. The equation 1 may be expanded as shown below:

$$Q(\Phi, R) = \frac{Q(\Phi, R)}{\prod_{i=1}^m q(\Phi_i, r_i)} \prod_{i=1}^m q(\Phi_i, r_i) \quad (8)$$

With near-equilibrium assumption, we may safely learn local distributions from data collected in equilibrium states and relevant environmental conditions. However, propagation of global molecular systems by dynamic assembly of such local distributions is significantly more challenging. Continuity restraints of relevant Φ variables is probably necessary, this may be realized through smoothing within certain spatial range. For equilibrium system, propagation of a molecular system under thermal fluctuation may be carried out with Langevin equation (equation 6) with a white noise term associated with a given temperature. However, in near equilibrium scenario, two choices need to be made for propagating the molecular system. The first is utilize either a maximum likelihood or bayesian approach to determine control variable at each molecule, with later being significantly more expensive. The second choice is to select a proper smoothing procedure to prevent large variance in control variables during the inference process. With these issues taken care of, assuming that local distributions $q(\Phi_i, r_i)$ have been learned with high accuracy, similar assembly and propagation procedures may be utilized as in the equilibrium case except with Φ included and stochastic forces added according to corresponding temperature at each molecule. Large variance of parameters such as temperature and pressure may derail such simple treatment. Significant exploration and development is necessary in these regards. Nonetheless, this opens a potential highly efficient and probabilistic pathway for treatment of near equilibrium massive complex molecular systems (e.g. a cell).

Rapid automatic search for implicit manifold

Due to both local and long range interactions/correlations in condensed molecular systems, the real dimensionality of which is significantly smaller than the number of degrees of freedom (DOF). Correlations inevitably reduce dimensionality. For example, a two dimensional system with variables (x, y) satisfying $x^2 + y^2 = 1$ is fundamentally a one dimensional manifold on a circle with a radius of 1. Similarly, considering 1000 rigid water model molecules (each with 6 DOFs) in a fixed rigid box. The number of DOF for the molecular system is 5997 (three DOFs are removed by

fixed center of mass), but its real dimensionality is reduced to a unknown but significantly smaller number due to complex correlations that are dependent upon environmental variables (e.g. temperature, pressure, container material). Van der Waals interactions, hydrogen bonding networks, dipolar and multipolar interactions all contribute to correlations and dimensionality reduction in water. Conventional way of understanding underlying manifolds for molecular systems is to perform dimensionality reduction analysis on sufficiently sampled trajectories. However, popular principal component analysis does not treat nonlinear correlations properly, and many nonlinear algorithms have their own limitations.⁴⁴ More importantly, these dimensionality reduction methodologies are usually utilized as a post processing step for understanding molecular systems after expensive sampling dominated by RLS has been performed. So the goal is to understand manifolds as one of terminal goals, rather than utilizing manifolds to reduce computational cost. Learned local distributions fundamentally maps to local manifolds with significant statistical weights and the remaining configurational space with negligible weights. Dynamic assembly of local distributions is, therefore, an implicit manifold assembly and search process on the one hand, and utilizes manifolds to reduce consumption of computational resources on the other hand. This is because negative gradients of local distributions always point to configurational regions of heavier statistical weights. Upon assembly of local distributions in propagation driven by derivatives of approximate instantaneous GJD density with respect to coordinates, a molecular system either stay on its manifold (free energy valleys) with fluctuations dependent upon temperature or rapidly return to the manifold when being away from it. To state alternatively, construction of GJD by assembly of local distributions according to equation 1 is equivalent to construction of global manifold by stitching together local manifolds embedded in local distributions without any manual intervention.

It is interesting to note that when viewed from the manifold perspective, LDT is effectively a completely automatic, significantly more accurate and efficient implicit counterpart of Metadynamics when local distributions were fit accurately and assembled properly. In Metadynamics, one first guess or compute for guiding collective variables (CVs), which is essentially an explicit and significantly simplified representation of the manifold for a target molecular system in a given

coordinate system. This is a highly challenging task, usually some iterative process is necessary but accuracy of resulting CVs has no guarantee, and no systematic theory is available for explicit searching of CVs. Subsequently explicit biases are accumulated to compute probability density of visited segments along CVs. In a properly implemented LDT, a target molecular system in propagation is automatically and implicitly maintained on its manifold, so the challenge of searching for CVs is met implicitly. Additionally, no bias is necessary and an unnormalized probability density is directly computed for each visited configuration.

Toward computation driven molecular sciences

Recent neural network force fields has demonstrated significant improvement in accuracy,^{44,62–64} albeit with accompanying reduction of efficiency when compared with conventional atomistic MD simulations. With further development of density estimation/fitting, local distributions may be built from all atom simulations based on neural network force fields of near quantum accuracy, or directly from highly accurate density functional theory *ab initio* simulations, and subsequently utilized to compose global distributions via dynamic assembly of local distributions as described by the LDT. Such combination may realize long-desired near-quantum accuracy and superior efficiency beyond conventional coarse grained models. With corresponding dramatic and simultaneous improvement of both accuracy and efficiency brought by LDT, molecular biology and nanotechnology research may experience a transition from experiment driven to computation driven as spatial and time scales will be accessible by present and computational facility expected in a few years.

Most proteins are dynamic molecules with their function supported by rich conformational transformations in response to environmental stimuli. However, if different solvent and thermodynamic conditions and corresponding conformational distributions for proteins are considered, available experimental data for any specific condition are certainly not sufficient for learning. Deficiency of structural data is even more severe for denatured states of proteins, nucleic acids and

other biomolecular systems (e.g. membranes). Presently, modeling of diverse thermodynamic and solvent conditions and denatured states relies heavily on all atom conventional MD simulations, which on the one hand is not sufficiently accurate, and on the other hand are limited to micro-second time scales in routine investigations of typical proteins for small research groups, and simulation of large complexes and more extensive biomolecular systems is much more challenging. Development of LDT for efficient and accurate construction of local distributions, when combined with one-time near quantum level MD simulations for general biomolecular systems has the potential of bridging this gap, and realize routine simulations of large molecular complexes on realistic time scales (milli-seconds and longer). Many present experiment dominated molecular biology research (e.g. protein-protein interactions and protein-drug interactions) may experience transition to computation driven with dramatically improved efficiency. This is especially true for proteins and other biomolecules that are marginally stable and hard to express and store under regular experimental conditions.

Establishment of a chain of tools from high level first principle calculations to simulation of large complex molecular systems has been long standing wish for molecular simulation community. Conventionally, coarse-graining has been the only available option and has made great contributions. Development and implementation of LDT in various general molecular systems provides a potential alternative pathway in this regard. However, to realize this goal, significant effort is necessary for development of algorithms in fitting local distributions for a wide variety of molecular systems. Condensed matter in general, and biological systems in particular, are organized in hierarchical structures with distinct correlation patterns over different length and time scales. Such characteristics were well summarized by Anderson⁶⁵ decades ago and significant efforts have been invested in multi-scale algorithm development in many subjects.⁶⁶⁻⁶⁹ As discussed above, local distributions are essentially manifolds of local regions under various composition and environmental conditions. The specific meaning of “local” is dependent upon definition of comprising unit on the one hand, and upon length scales on the other hand. Implementation of LDT on multiple scales, and how should it interact with coarse-graining or evolve independently, is a fully open

field awaiting intensive exploration.

Two distinct ways of averaging

Conventional FF parameterization is fundamentally a construction of potential of mean force (PMF)^{70,71} by integration/averaging as shown below:

$$U(x) = \int U(x, y) dy \quad (9)$$

For fitting of atomistic FF from *ab initio* calculations, y correspond to electronic DOFs, for fitting of coarse-grained FF from atomistic simulations, y correspond to all atomic DOFs other than coarse-grained sites. PMF accurately reproduce behavior of variable x when a time scale separation exists between x and y . Therefore, conventional molecular simulation framework is based on the idea of PMF.

Local distributions are clearly results of statistical averaging based on data obtained from expensive local sampling, either through experimental or computational approaches. Essentially, relative frequency of visiting many different configurations are recorded. However, there is no explicit reduction of variables in this process as in the case of PMF integration in FF parameterization (i.e. resolution is maintained). These statements seem to be contradictory as the process of averaging inevitably results in annihilation of some details. One would certainly like to know what is annihilated during the averaging process of fitting local distributions. In the mapping from complete molecular DOFs in local regions of a molecular configuration to local probability, based on correlations among different DOFs, some of which are implicitly and adaptively eliminated by neural networks. Such implicit process may be alternatively explained in terms of configurational space discretization (CSD). In computers there is no strictly “continuous” variables anymore as everything is stored by discrete “boxes” in CPU registers, memory chips and hard drives. So all modeling in computer is performed on lattices defined by float point number discretization! In fitting of local distributions via neural networks, while input of molecular configurations has the res-

olution of lattices defined by selected float point digits, there is probably further implicit merging (coarse graining) of different lattice boxes not necessarily uniformly both on different dimensions and on different positions of the same dimension. Such implicit and adaptive annihilation of resolution on various places of the configurational space by the fitting machinery (neural networks) is schematically illustrated in Fig. 5. More specifically, for the given data set $D = \{X_1, X_2, \dots, X_d\}$, X_i may be partitioned into n_i local regions, denote the transformed coordinate vector for the j th local region from the i th record as r'_{ij} . Local distribution is effectively obtained by an summation established by the trained neural network:

$$q(\Phi, r) \propto \sum_{i=1}^d \sum_{j=1}^{n_i} f_{\Phi, W}(r'_{ij} - r) 1(r'_{ij} - r) \quad (10)$$

$$1(r'_{ij} - r) = \begin{cases} 1, & \text{if } (r'_{ij} - r) \leq g_{cut}(r) \\ 0, & \text{otherwise} \end{cases} \quad (11)$$

with W being neural network parameters. $f_{\Phi, W}(r'_{ij} - r)$ is an implicit weight function for evaluating contribution of each data point occurring represented by the indicator function $1(r'_{ij} - r)$. $f_{\Phi, W}(r'_{ij} - r)$ decays with $|r'_{ij} - r|$ but not necessarily isotropically, with specifics determined by W and the neural network architecture. The qualitative trend is that $f_{\Phi, W}(r'_{ij} - r)$ approaches a delta function for configurational space region with sufficiently high data density, and becomes flatter for configurational space region with lower data density. $g_{cut}(r)$ is a r dependent cutoff determined by the trained neural network. The indicator function maybe assimilated into a new weight function f' :

$$q(\Phi, r) \propto \sum_{i=1}^d \sum_{j=1}^{n_i} f'_{\Phi, W}(r'_{ij} - r) \quad (12)$$

Therefore, LDT adopts a distinctive path of averaging based on implicit adaptive configurational space discretization (CSD) instead of explicit integrating out selected DOFs adopted by PMF.

Ultimately specifics of such heterogeneous CSD are likely to be determined by details of loss function, network architecture, optimization process and their interactions. However, presently,

how such implicit process relates to corresponding neural networks is not transparent. There is no published research on neural networks regarding this topic to the best of my knowledge. Understanding such implicit CSD is likely an essential step to be accomplished in constructing transparent white box neural networks. Manual configuration space discretization has been performed to facilitate free energy analysis.⁷² However, proper CSD strategy is usually different for distinct molecular systems and is not necessarily achievable even after significant human efforts. Therefore, to develop transparent, easy-to-manipulate and automatically adaptive schemes for CSD in fitting behavior of neural networks is an important open field to explore.

LDT and neural network force fields

Neural networks, with their strong fitting capability, are at the core of representing both local distributions in LDT and neural network force fields. In these two methodologies, many body correlations are accounted for naturally and the accuracy ceiling due to fixed functional forms (describing molecular interactions) is lifted. However, local distributions are fundamentally different from neural network force fields. The former is trained to reproduce local probability distributions while the latter is trained to reproduce exact energy for each given atomic configuration. By incorporating both energetic and entropic contributions within local regions, local distributions when assembled according to LDT essentially eliminates RLS and significantly increases efficiency. Neural network force fields replace conventional molecular fields in established framework of “force fields + sampling” without considering RLS, and improve accuracy at the cost of decreased efficiency. However, as discussed above, fitting local distributions is essentially a density estimation problem, which is in general significantly more difficult than supervised training of neural network force fields. Presently, development of neural network force fields is the mainstream of research bridging artificial intelligence and molecular simulations.

Conclusions and prospects

RLS in molecular simulations consumes large amount of computational resources on the one hand and slows down exploration of relevant research fields dramatically on the other hand. The LFEL approach was developed to address RLS previously. However, the formulation and its exemplary implementation in protein structural refinement, while demonstrated tremendous potentials, is limited to a single implicit set of given environmental conditions. Here I propose the local distribution theory to generalize LFEL for addressing variable environmental conditions and near-equilibrium application scenarios. As a matter of fact, essentially all biological systems are off equilibrium to various extent. Despite the simple theoretical proposal presented here, extending implementation of LDT to near-equilibrium poses great challenges and significant exploratory efforts are necessary. Theoretical connection and fundamental differences of LDT with metadynamics, with neural network force fields, with RPDD based AI-driven protein structural research, and with PMF based framework of conventional molecular simulation in general are discussed. It is hoped that discussions and speculations herein stimulate more interest and attract more scientists in further development and application of the local distribution theory.

Abbreviations

CSD: Configurational Space Discretization. CV: Collective Variable. DOF: Degree of Freedom. FF: Force Fields. GJD: Global Joint Distribution. GSFE: Generalized Solvation Free Energy. MD: Molecular Dynamics. MC: Monte Carlo. LDT: Local Distribution Theory. LFEL: Local Free Energy Landscape. PMF: Potential of Mean Force. RLS: Repetitive Local Sampling. RPDD: Residue Pair Distance Distribution.

Acknowledgement

I thank Professor Jingkai Gu for encouragement and Professors Zhonghan Hu and Yaoqi Zhou for comments when this theory was conceived.

References

- (1) Dror, R. O.; Dirks, R. M.; Grossman, J.; Xu, H.; Shaw, D. E. Biomolecular Simulation: A Computational Microscope for Molecular Biology. *Annual Review of Biophysics* **2012**, *41*, 429–452, PMID: 22577825.
- (2) Bedrov, D.; Piquemal, J.-P.; Borodin, O.; MacKerell, A. D., Jr.; Roux, B.; Schroeder, C. Molecular Dynamics Simulations of Ionic Liquids and Electrolytes Using Polarizable Force Fields. *CHEMICAL REVIEWS* **2019**, *119*, 7940–7995.
- (3) Rudzinski, J. F.; Noid, W. G. Coarse-graining entropy, forces, and structures. *Journal of Chemical Physics* **2011**, *135*.
- (4) Marrink, S. J.; Tieleman, D. P. Perspective on the martini model. *Chemical Society Reviews* **2013**, *42*, 6801–6822.
- (5) Noid, W. G. Perspective: Coarse-grained models for biomolecular systems. *Journal of Chemical Physics* **2013**, *139*.
- (6) Saunders, M. G.; Voth, G. A. Coarse-graining methods for computational biology. *Annual Review of Biophysics* **2013**, *42*, 73–93.
- (7) Ruff, K. M.; Harmon, T. S.; Pappu, R. V. CAMELOT: A machine learning approach for Coarse-grained simulations of aggregation of block-copolymeric protein sequences. *Journal of Chemical Physics* **2015**, *143*, 1–19.

- (8) Kmiecik, S.; Gront, D.; Kolinski, M.; Wieteska, L.; Dawid, A. E.; Kolinski, A. Coarse-Grained Protein Models and Their Applications. *Chemical Reviews* **2016**, *116*, 7898–7936.
- (9) Hafner, A. E.; Krausser, J.; Saric, A. Minimal coarse-grained models for molecular self-organisation in biology. *CURRENT OPINION IN STRUCTURAL BIOLOGY* **2019**, *58*, 43–52.
- (10) Sambasivan, R.; Das, S.; Sahu, S. K. A Bayesian perspective of statistical machine learning for big data. *COMPUTATIONAL STATISTICS* **2020**, *35*, 893–930.
- (11) Joshi, S. Y.; Deshmukh, S. A. A review of advancements in coarse-grained molecular dynamics simulations. *Molecular Simulation* **2020**, *0*, 1–18.
- (12) Gkeka, P. et al. Machine Learning Force Fields and Coarse-Grained Variables in Molecular Dynamics: Application to Materials and Biological Systems. *Journal of Chemical Theory and Computation* **2020**, *16*, 4757–4775, PMID: 32559068.
- (13) Bernardi, R. C.; Melo, M. C. R.; Schulten, K. Enhanced sampling techniques in molecular dynamics simulations of biological systems. *BIOCHIMICA ET BIOPHYSICA ACTA-GENERAL SUBJECTS* **2015**, *1850*, 872–877.
- (14) Mlynsky, V.; Bussi, G. Exploring RNA structure and dynamics through enhanced sampling simulations. *CURRENT OPINION IN STRUCTURAL BIOLOGY* **2018**, *49*, 63–71.
- (15) Yang, Y. I.; Shao, Q.; Zhang, J.; Yang, L.; Gao, Y. Q. Enhanced sampling in molecular dynamics. *Journal of Chemical Physics* **2019**, *151*.
- (16) Wang, A.-h.; Zhang, Z.-c.; Li, G.-h. Advances in Enhanced Sampling Molecular Dynamics Simulations for Biomolecules. *CHINESE JOURNAL OF CHEMICAL PHYSICS* **2019**, *32*, 277–286.
- (17) Cao, X.; Tian, P. Molecular free energy optimization on a computational graph. *RSC Adv.* **2021**, *11*, 12929–12937.

- (18) Long, S.; Tian, P. A simple neural network implementation of generalized solvation free energy for assessment of protein structural models. *RSC Advances* **2019**, *9*, 36227–36233.
- (19) Cao, X.; Tian, P. “Dividing and Conquering” and “Caching” in Molecular Modeling. *International Journal of Molecular Sciences* **2021**, *22*.
- (20) Senior, A. W. et al. Improved protein structure prediction using potentials from deep learning. *Nature* **2020**, *577*, 706–710.
- (21) Suh, D.; Lee, J. W.; Choi, S.; Lee, Y. Recent Applications of Deep Learning Methods on Evolution- and Contact-Based Protein Structure Prediction. *International Journal of Molecular Sciences* **2021**, *22*.
- (22) Pakhrin, S. C.; Shrestha, B.; Adhikari, B.; KC, D. B. Deep Learning-Based Advances in Protein Structure Prediction. *International Journal of Molecular Sciences* **2021**, *22*.
- (23) Skolnick, J.; Gao, M. The role of local versus nonlocal physicochemical restraints in determining protein native structure. *Current Opinion in Structural Biology* **2021**, *68*, 1–8, Protein-Carbohydrate Complexes and Glycosylation Sequences and Topology.
- (24) Ju, F.; Zhu, J.; Shao, B.; Kong, L.; Liu, T.-Y.; Zheng, W.-M.; Bu, D. CopulaNet: Learning residue co-evolution directly from multiple sequence alignment for protein structure prediction. *Nature Communications* **2021**, *12*, 2535.
- (25) Jing, X.; Xu, J. Improved protein model quality assessment by integrating sequential and pairwise features using deep learning. *Bioinformatics* **2020**, *36*, 5361–5367.
- (26) Zhao, K.-l.; Liu, J.; Zhou, X.-g.; Su, J.-z.; Zhang, G.-j. Structural bioinformatics MMpred : a distance-assisted multimodal conformation sampling for de novo protein structure prediction. *Bioinformatics* **2021**, 1–8.
- (27) Xia, Y.-h.; Peng, C.-x.; Zhou, X.-g.; Zhang, G.-j. Structural bioinformatics A Sequential

- Niche Multimodal Conformational Sampling Algorithm for Protein Structure Prediction. *Bioinformatics* **2021**, 1–9.
- (28) Wu, F.; Xu, J. Deep template-based protein structure prediction. *PLoS computational biology* **2021**, *17*, 1–18.
- (29) Zhang, H.; Bei, Z.; Xi, W.; Hao, M.; Ju, Z. Evaluation of residue-residue contact prediction methods : From retrospective to prospective. *PLOS Comput. Biol.* **2021**, *17*, 1–33.
- (30) Privalov, P. L. Cold Denaturation of Protein. *Critical Reviews in Biochemistry and Molecular Biology* **1990**, *25*, 281–306.
- (31) Sanfelice, D.; Temussi, P. A. Cold denaturation as a tool to measure protein stability. *Biophysical Chemistry* **2016**, *208*, 4–8, SIBPA 2014 - XXII SIBPA Congress.
- (32) Jones, S. W. Overview of Voltage-Dependent Calcium Channels. *Journal of Bioenergetics and Biomembranes* **1998**, *30*, 299–312.
- (33) Wu, J.; Lewis, A. H.; Grandl, J. Touch, Tension, and Transduction – The Function and Regulation of Piezo Ion Channels. *Trends in Biochemical Sciences* **2017**, *42*, 57–71.
- (34) Jiang, Y.; Yang, X.; Jiang, J.; Xiao, B. Structural Designs and Mechanogating Mechanisms of the Mechanosensitive Piezo Channels. *Trends in Biochemical Sciences* **2021**, *46*, 472–488.
- (35) Mackerell Jr., A. D. Empirical force fields for biological macromolecules: Overview and issues. *Journal of Computational Chemistry* **2004**, *25*, 1584–1604.
- (36) Kumar, A.; Yoluk, O.; MacKerell Jr., A. D. FFParm: Standalone package for CHARMM additive and Drude polarizable force field parametrization of small molecules. *Journal of Computational Chemistry* **2020**, *41*, 958–970.
- (37) Oweida, T. J.; Kim, H. S.; Donald, J. M.; Singh, A.; Yingling, Y. G. Assessment of AMBER Force Fields for Simulations of ssDNA. *Journal of Chemical Theory and Computation* **2021**, *17*, 1208–1217, PMID: 33434436.

- (38) Huai, Z.; Shen, Z.; Sun, Z. Binding Thermodynamics and Interaction Patterns of Inhibitor-Major Urinary Protein-I Binding from Extensive Free-Energy Calculations: Benchmarking AMBER Force Fields. *Journal of Chemical Information and Modeling* **2021**, *61*, 284–297, PMID: 33307679.
- (39) Alford, R. F. et al. The Rosetta All-Atom Energy Function for Macromolecular Modeling and Design. *Journal of Chemical Theory and Computation* **2017**, *13*, 3031–3048, PMID: 28430426.
- (40) Sasse, A.; de Vries, S. J.; Schindler, C. E. M.; de Beauchêne, I. C.; Zacharias, M. Rapid Design of Knowledge-Based Scoring Potentials for Enrichment of Near-Native Geometries in Protein-Protein Docking. *PLOS ONE* **2017**, *12*, 1–19.
- (41) Narykov, O.; Bogatov, D.; Korkin, D. DISPOT: a simple knowledge-based protein domain interaction statistical potential. *Bioinformatics* **2019**, *35*, 5374–5378.
- (42) Dowsland, K. A.; Thompson, J. M. In *Handbook of Natural Computing*; Rozenberg, G., Bäck, T., Kok, J. N., Eds.; Springer Berlin Heidelberg: Berlin, Heidelberg, 2012; pp 1623–1655.
- (43) Behler, J. Perspective: Machine learning potentials for atomistic simulations. *Journal of Chemical Physics* **2016**, *145*.
- (44) Gkeka, P. et al. Machine Learning Force Fields and Coarse-Grained Variables in Molecular Dynamics: Application to Materials and Biological Systems. *Journal of Chemical Theory and Computation* **2020**, *16*, 4757–4775, PMID: 32559068.
- (45) Parzen, E. On Estimation of a Probability Density Function and Mode. *The Annals of Mathematical Statistics* **1962**, *33*, 1065 – 1076.
- (46) Uria, B.; Côté, M.-A.; Gregor, K.; Murray, I.; Larochelle, H. Neural Autoregressive Distribution Estimation. *Journal of Machine Learning Research* **2016**, *17*, 1–37.

- (47) Kobyzev, I.; Brubaker, M. A. Normalizing Flows: Introduction and Ideas. *ArXiv* **2019**, 1–35.
- (48) Papamakarios, G.; Nalisnick, E. Normalizing Flows for Probabilistic Modeling and Inference. *ArXiv* **2019**, 1–60.
- (49) Noé, F.; Olsson, S.; Kohler, J.; Wu, H. Boltzmann generators: Sampling equilibrium states of many-body systems with deep learning. *Science* **2019**, *365*, eaaw1147.
- (50) Liu, Q.; Xu, J.; Jiang, R.; Hung, W. Density estimation using deep generative neural networks. *Proceedings of the National Academy of Sciences* **2021**, *118*.
- (51) Cun, Y.; Huang, F. Loss functions for discriminative training of energy-based models. AIS-TATS 2005 - Proceedings of the 10th International Workshop on Artificial Intelligence and Statistics. 2005; pp 206–213, 10th International Workshop on Artificial Intelligence and Statistics, AISTATS 2005 ; Conference date: 06-01-2005 Through 08-01-2005.
- (52) BakIr, G.; Hofmann, T.; Schölkopf, B.; Smola, A. J.; Taskar, B.; Vishwanathan, S. *Predicting Structured Data*; The MIT Press, 2007.
- (53) Zhao, J.; Mathieu, M.; Lecun, Y.; Artificial, F. Energy-based generative adversarial networks. *ArXiv* **2017**, 1–17.
- (54) Liu, M.; Yan, K.; Oztekin, B.; Ji, S. GraphEBM: Molecular Graph Generation with Energy-Based Models. *ArXiv* **2020**,
- (55) Mordatch, I. Compositional Visual Generation with Energy Based Models. *ArXiv* **2020**,
- (56) Grathwohl, W.; Duvenaud, D. Your classifier is secretly an energy based model and you should treat it like one. *ArXiv* **2020**, 1–23.
- (57) Chan, K. H. R.; Yu, Y.; You, C.; Qi, H.; Wright, J.; Ma, Y. ReduNet: A White-box Deep Network from the Principle of Maximizing Rate Reduction. 2021.
- (58) Jumper, J. et al. Highly accurate protein structure prediction with AlphaFold. *Nature* **2021**,

- (59) Jing, X.; Xu, J. Fast and effective protein model refinement using deep graph neural networks. *Nature Computational Science* **2021**, *1*, 462–469.
- (60) Baek, M. et al. Accurate prediction of protein structures and interactions using a three-track neural network. *Science* **2021**, *373*, abj8754.
- (61) Frenkel, D.; Smit, B. In *Understanding Molecular Simulation (Second Edition)*, second edition ed.; Frenkel, D., Smit, B., Eds.; Academic Press: San Diego, 2002; pp 139–163.
- (62) Zhang, L.; Han, J.; Wang, H.; Car, R.; E, W. Deep Potential Molecular Dynamics: A Scalable Model with the Accuracy of Quantum Mechanics. *Phys. Rev. Lett.* **2018**, *120*, 143001.
- (63) Bedolla, E.; Padierna, L. C.; Castañeda-Priego, R. Machine learning for condensed matter physics. *Journal of Physics: Condensed Matter* **2020**, *33*, 053001.
- (64) Lu, D.; Wang, H.; Chen, M.; Lin, L.; Car, R.; E, W.; Jia, W.; Zhang, L. 86 PFLOPS Deep Potential Molecular Dynamics simulation of 100 million atoms with ab initio accuracy. *Computer Physics Communications* **2021**, *259*, 107624.
- (65) Anderson, P. More is different. *Science* **1972**, *177*, 393–396.
- (66) Franco, A. A.; Rucci, A.; Brandell, D.; Frayret, C.; Gaberscek, M.; Jankowski, P.; Johansson, P. Boosting Rechargeable Batteries R&D by Multiscale Modeling: Myth or Reality? *Chemical Reviews* **2019**, *119*, 4569–4627, PMID: 30859816.
- (67) Henning, P.; Peterseim, D. Oversampling for the Multiscale Finite Element Method. *Multiscale Modeling & Simulation* **2013**, *11*, 1149–1175.
- (68) Abdulle, A.; Weinan, E.; Engquist, B.; Vanden-Eijnden, E. The heterogeneous multiscale method. *Acta Numerica* **2012**, *21*, 1?87.
- (69) Ramstead, M. J. D.; Kirchhoff, M. D.; Constant, A.; Friston, K. J. Multiscale integration: beyond internalism and externalism. *Synthese* **2021**, *198*, 41–70.

- (70) Kirkwood, J. G. Statistical Mechanics of Fluid Mixtures. *Journal of Chemical Physics* **1935**, 3, 300–313.
- (71) Roux, B. The calculation of the potential of mean force using computer simulations. *Computer Physics Communications* **1995**, 91, 275–282.
- (72) Wang, K.; Long, S.; Tian, P. Configurational space discretization and free energy calculation in complex molecular systems. *Scientific Reports* **2016**, 6, 22217.

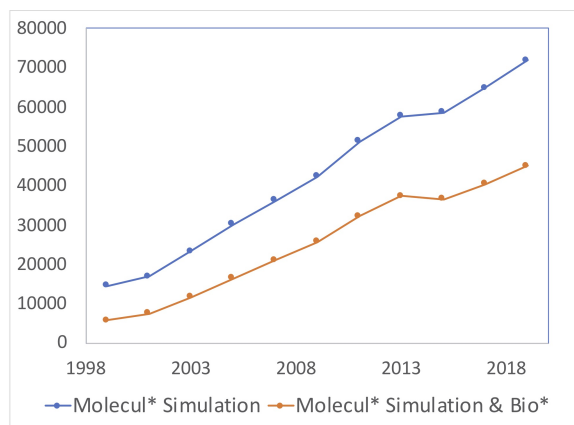


Figure 1: The number of publications retrieved from web of science on Jun. 1st 2021 with subject word "molecul* simulation" and "molecul* simulation & bio" respectively. The corresponding time frame is every two years starting from 1999. The first data point is the number of papers published in year 1999 and 2000.

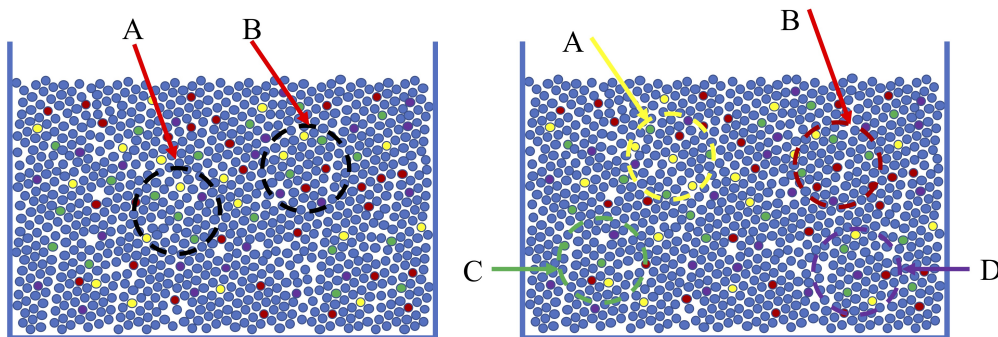


Figure 2: Schematic illustration of RLS. Left: the spatial perspective. A) and B) are two different spherical bulk spaces. We expect the same local distributions after sufficiently long simulations of the whole molecular system. In such cases, spherical and partial spherical spaces near or on interfaces have different local distributions from that of the bulk, special treatment of such spherical spaces engenders significant difficulty. Right: indistinguishable particle and GSFE perspective. All particles of the same species are indistinguishable, so should be local distributions of local regions defined by spherical spaces with such a particle as the origin. This removes the need for special treatment of all interfacial issues as different interfaces may be simply defined as more cases of particle packing surrounding a given particle with well defined statistical weight under given thermodynamic and environmental conditions. A), B), C) and D) are examples of surrounding local regions of different particle species.

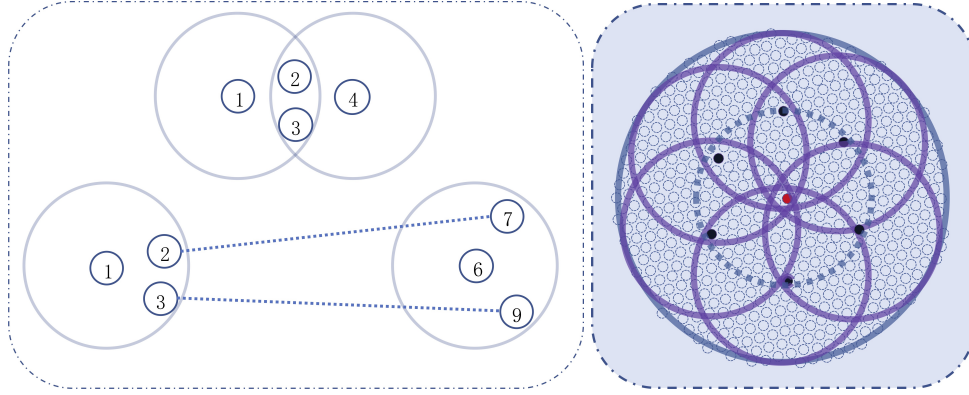


Figure 3: Schematic representation of the short range, mediated and long range interactions as implemented in ref. Left: particles (1,2,3), (2,3,4) and (6,7,9) are directly interacting with short range interactions. (1,4) are interacting through mediation by (2,3), (2,7) and (3,9) have direct long range interactions. Right: here the focus is the central red particle, which define a region with boundary being shown as a dotted partially transparent blue line. Each of all other particles within this region defines a local distribution, six of the most further of such regions are represented as purple circles. The central red particle experience forces from all of local distributions surrounding each of its neighbors. In this way, short range and mediated interactions are effectively accounted for simultaneously. In summary, for the central red particle, it experiences short range interactions from particles within the dotted partial transparent blue circle, mediated interactions from particles between the dotted blue circle and large solid blue circle, and long range interactions from the region outside the large blue circle.

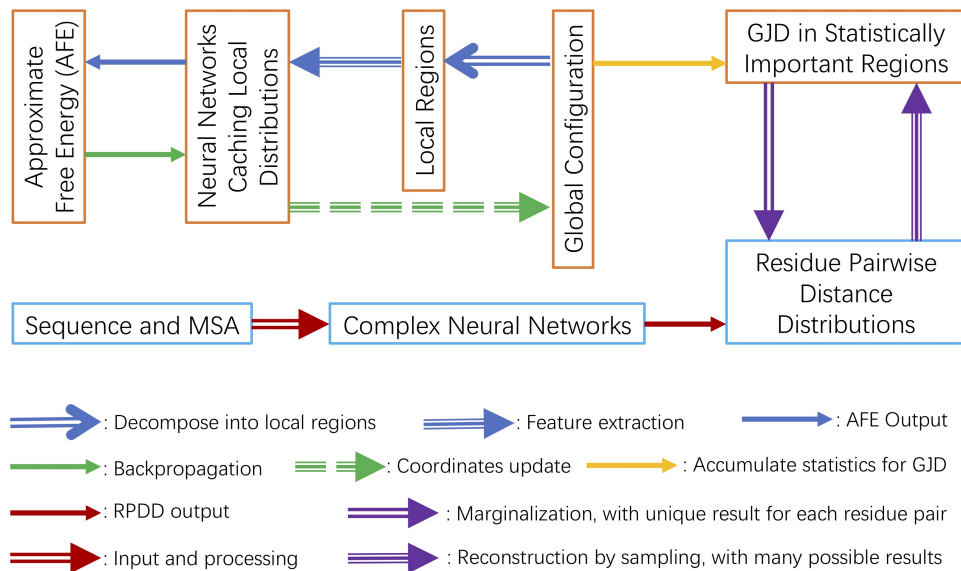


Figure 4: Schematic comparison between LDT based end-to-end protein structure modeling (top orange boxes) and mainstream RPDD based protein structure prediction and refinement schemes (bottom blue boxes). It is important to note that LDT based modeling aims to generate the GJD, which is the most comprehensive information for any complex molecular systems and is generally applicable. The marginalization from the GJD to pairwise residue distance distributions is an irreversible process with deterministic results and significant information loss on correlations among different pairwise distances. The converse process is a highly expensive process with sampling and optimization involved, due to complexity of correlations among different distances, resulting global distribution is highly dependent both on initialization and the optimization procedures.

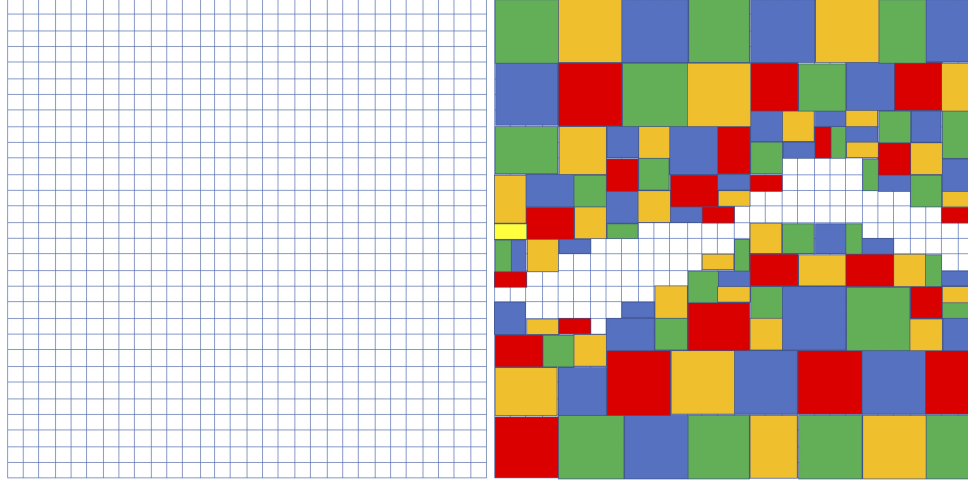


Figure 5: Schematic illustration of CSD. Left: natural discretization of two dimensional configurational space by float point digits. Right: a imagined heterogeneous CSD resulted from fitting of neural network on local distributions, and the highest density is supposedly in the white region where CSD is as fine as lattices determined by float point digits. Qualitatively, finer discretization corresponds to region with high data density and coarser discretization corresponds to region with lower data density.