


Cite this: *RSC Adv.*, 2022, 12, 25010

Collective variable discovery in the age of machine learning: reality, hype and everything in between†‡

Soumendranath Bhakat *

Understanding the kinetics and thermodynamics profile of biomolecules is necessary to understand their functional roles which has a major impact in mechanism driven drug discovery. Molecular dynamics simulation has been routinely used to understand conformational dynamics and molecular recognition in biomolecules. Statistical analysis of high-dimensional spatiotemporal data generated from molecular dynamics simulation requires identification of a few low-dimensional variables which can describe the essential dynamics of a system without significant loss of information. In physical chemistry, these low-dimensional variables are often called collective variables. Collective variables are used to generate reduced representations of free energy surfaces and calculate transition probabilities between different metastable basins. However the choice of collective variables is not trivial for complex systems. Collective variables range from geometric criteria such as distances and dihedral angles to abstract ones such as weighted linear combinations of multiple geometric variables. The advent of machine learning algorithms led to increasing use of abstract collective variables to represent biomolecular dynamics. In this review, I will highlight several nuances of commonly used collective variables ranging from geometric to abstract ones. Further, I will put forward some cases where machine learning based collective variables were used to describe simple systems which in principle could have been described by geometric ones. Finally, I will put forward my thoughts on artificial general intelligence and how it can be used to discover and predict collective variables from spatiotemporal data generated by molecular dynamics simulations.

Received 13th June 2022
Accepted 20th August 2022

DOI: 10.1039/d2ra03660f

rsc.li/rsc-advances

1 Introduction

Over the past three decades, the major focus of structural biology and biophysics has been to understand the conformational dynamics of biological systems across a broad range of timescales with high spatial resolution (ability to visualise molecular motion). Molecular dynamics (MD) simulation acts as a computational microscope which can capture biologically relevant conformational dynamics across varied timescales with high spatial resolution.¹ Major advances in computational hardware and algorithms led to an ever-growing use of MD simulation in varied biological systems. A typical MD simulation generates high-dimensional spatiotemporal data which captures complex molecular motions. Simulation of increasingly larger molecules at ever increasing time scales leads to the

‘curse of dimensionality’ which can be summarised as follows: “as the size of the biomolecule and the simulation length increases, it also increases the number of explanatory temporal variables (*e.g.* H-bond distances, radius of gyration, RMSD, dihedral angles *etc.*) and the problem of structure discovery using temporal variables gets harder. This is analogous to the problem of variable selection during model fitting in machine learning.”²

Capturing low-dimensional representation (without significant loss of information regarding slow conformational degrees of freedom) from high-dimensional temporal data is an open area of research in biomolecular simulation. An excellent example of this is the large scale flap dynamics in plasmepsin-II and BACE-1.³ In plasmepsin-II, the sine/cos transformation of χ_1 and χ_2 angles of 20 residues present in the flap region generates 58 dimensions. Such a high number of dimensions makes it incredibly difficult to capture slow degrees of freedom which dominate flap opening. However, careful analysis of the probability distributions showed that flipping of χ_1 and χ_2 angles of conserved tyrosine (Tyr) governs the flap dynamics in plasmepsin-II and BACE-1.⁴ Such low-dimensional projections which best captures the conformational changes are known as collective variables (CVs) or order parameters (not to be confused with NMR order parameter which is a measure of conformational entropy). Temporal evolution along CVs provide

Department of Biochemistry and Biophysics, Perelman School of Medicine, University of Pennsylvania, Pennsylvania 19104-6059, USA. E-mail: bhakatsoumendranath@gmail.com; Tel: +1 30549 32620

† Major part of the article was written when the author was affiliated to Division of Biophysical Chemistry, Center for Molecular Protein Science, Department of Chemistry, Lund University, Sweden and Department of Biochemistry and Molecular Biophysics, Washington University School of Medicine, St Louis, Missouri, USA.

‡ Electronic supplementary information (ESI) available. See <https://doi.org/10.1039/d2ra03660f>



thermodynamic and kinetic informations about conformational changes.^{5–7} It is worth mentioning that the use of buzzword “automatic identification” to highlight new methods for CV discovery from MD simulation gives a false impression *i.e.* selection of good CV is a streamlined and well-defined process. But in practice, naive application of automated methods leads to uninterpretable and poor choice of CVs. Further such methods are not necessarily new, but in most cases direct applications or minor modifications of machine learning algorithms developed for classification (binary classifiers and their variants), signal processing (PCA, TICA, autoencoders *etc.*), time-series prediction (neural networks, hidden Markov model *etc.*) and natural language processing (LSTM, transformers) *etc.* Recent years saw a sharp rise in applying machine learning algorithms for CV selection on systems which can be represented by simple geometric CVs *e.g.* dihedral angles, H-bond, RMSDs *etc.* (highlighted in later sections).

Due to overuse of the words ‘AI/machine learning’ and flood of papers reporting new algorithms, the current field of CV discovery in biomolecular simulation is hard to navigate for a beginner or people from other disciplines. In this review we will classify the CV discovery process as following: (1) geometric (dihedral angles, H-bonds, RMSD *etc.*) and (2) abstract (PCA, ICA, neural network *etc.*). One question that is still up in the air: “Has the field of CV discovery in biomolecular simulation reached a plateau?”. I will provide my opinion on the aforementioned question and will put forward some challenges in order to test the robustness and general applicability of CV discovery algorithms in context of biomolecular simulation.

2 Metastability & collective variables

Before we dive into the definition of CV, we must understand the concept of metastability. Let’s consider basic pancreatic trypsin inhibitor (BPTI) as an example and focus on a particular H-bond between Ile18–N–O–Tyr35 (Fig. 1). In crystal conformation the aforementioned H-bond remains in a closed conformation. If

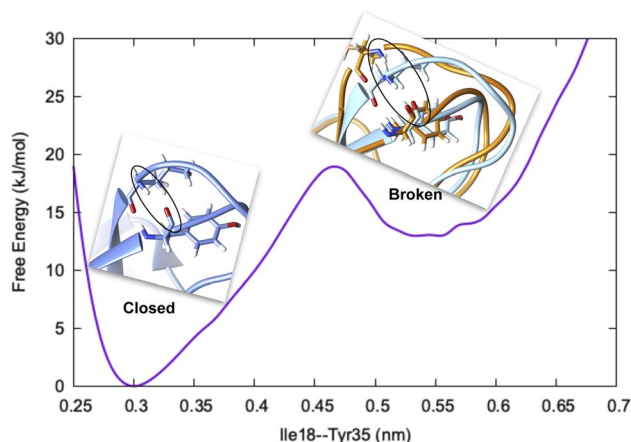


Fig. 1 Free energy surface along Ile18–Tyr35 H-bond distance showing two metastable states, closed (blue) and broken (orange) in BPTI. The backbone atoms involved in H-bond interaction are highlighted in black circle.

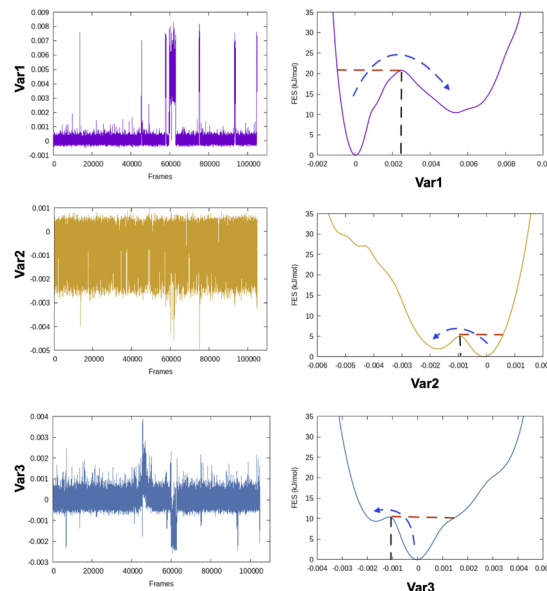


Fig. 2 Fluctuation of three different variables and their corresponding free energy surfaces highlighting the differences in free energy barrier. One can see that the barrier along Var1 is significantly higher compared to other two variables. Heuristically if one has to choose an optimum CV from these three variables, Var1 will be the CV of choice.

one performs a short MD simulation, the stationary probability distribution fluctuates around the closed basin, however, careful analysis of 1 ms MD simulation showed formation of broken conformations. A CV (in this case the H-bond distance) should be able to distinguish the metastable states (such as broken and closed conformations) and estimate apparent free energy profile along it as described by following equation:

$$F(\mathbf{x}) = -k_{\text{B}}T \ln P(\mathbf{x}) \quad (1)$$

where \mathbf{x} is the CV of choice.

The rugged nature of biomolecular energy landscape is by definition characterised by numerous metastable basins. An optimal CV is the one for which two metastable states are separated by high free energy barrier (Fig. 2). In most cases a single CV is not enough to capture the complexity of conformational landscape. Selection of CVs is an essential step to perform free energy/kinetics calculations and drive pathway based sampling methods *e.g.* metadynamics,⁸ steered MD,⁹ umbrella sampling¹⁰ *etc.*

For the sake of simplicity CVs can be divided into two categories: (1) geometric and (2) abstract.

2.1 Geometric CVs

Most commonly used geometric CVs that captures conformational dynamics in biomolecules are:

- (a) Distance: in biomolecular simulation two kinds of distances are most commonly used:
 - (i) distance between two atoms
 - (ii) distance between center of mass (COM) of two group of atoms.



Distance as CVs can be used as an input within adaptive sampling and enhanced sampling methods. For example, distance between COM of ligand and protein can be used as a CV to capture ligand unbinding using well-tempered metadynamics.^{4,11} Similarly H-bond distance can be used as CV (Fig. 3) to capture transition between closed and broken states. Such CV can be especially useful to study transient solvent exposure in protein which leads to hydrogen-deuterium exchange (paper in preparation).

(b) Switching function: a smoother version of distance CV is switching function (SF) which allows a smoother transition between two metastable states along a particular distance (Fig. 3). A typical mathematical form (other functional forms of switching function includes tanh, Gaussian, exponential, cubic *etc*) of SF can be described as follows:

$$s(r) = \frac{1 - \left(\frac{r - r_{ij}}{r_0}\right)^n}{1 - \left(\frac{r - r_{ij}}{r_0}\right)^m} \quad (2)$$

where r is the instantaneous distance between atoms i and j and r_{ij} is the minimum distance between atoms i and j . For $r < r_{ij}$, $s = 1.0$ while $r > r_{ij}$ the function decays smoothly to 0 (zero). r_0 is the value of distance where $s = 0.5$. n and m are the hyper-parameters which decides the steepness of the function.

A more general CV which combines multiple distances with switching function is known as contact map. It calculates the distances between a number of pairs of atoms and convert each distance by a switching function. Such CVs are useful where fluctuation along multiple distances governs conformational dynamics.

(c) Dihedral angle: tracking temporal evolution of dihedral angles (*e.g.* ϕ , ψ , ω , χ_1 , χ_2) during MD simulation is a well established method to capture conformational dynamics of biomolecules. Most common example includes tracking the conformational sampling in alanine dipeptide by probability distribution along ϕ and ψ angles. Recently, Bhakat and

Söderhjelm⁴ showed that the conformational dynamics of pepsin-like aspartic proteases can be captured by two dihedral angles: χ_1 and χ_2 angles of conserved Tyr (Fig. 4 highlights evolution of χ_1 during MD simulation). However in many cases, conformational dynamics of biomolecules can't be captured by a few dihedral angles. In complex cases, linear combinations of dihedral angles (more on this later) can be used as CVs to capture biomolecular dynamics.

(d) RMSD: RMSD is one of the commonly used CV which measures the similarity between two superimposed atomic coordinates. RMSD is the measure of average distance and in bimolecular simulation it measures deviation of atomic coordinates from the starting conformation using the following equation:

$$\text{RMSD} = \sqrt{\frac{1}{n} \sum_{i=1}^n d_i^2} \quad (3)$$

where d_i is the distance between atom i and a reference structure. RMSD is usually calculated for C α /backbone atoms of the entire protein or for a specific subset. Stock and coworkers argued that RMSD is not an optimal choice to capture local conformational (*e.g.* loop dynamics, domain motions *etc.*) changes at long-timescales.¹² In our unpublished study we have shown that RMSD analysis on a carefully chosen subset can able to capture conformational changes at longer timescales (Fig. 5). RMSD based CVs can be used within enhanced sampling and adaptive sampling frameworks to accelerate sampling of the conformational space.

Besides the aforementioned CVs, several other geometric variables (radius of gyration, eRMSD) are frequently used to analyse MD simulations. Softwares *e.g.* Plumed,¹⁵ gmx plugins integrated with Gromacs, CPPTRAJ,¹⁶ MDAnalysis,¹⁷ MDTraj¹⁸ have build-in capabilities to perform analysis of MD trajectories using geometric CVs.

2.2 Abstract CVs

Abstract CVs are usually linear or non-linear transformations of geometric CVs (*e.g.* dihedral angles, distances, RMSD *etc.*). However the former is often less intuitive compared to the

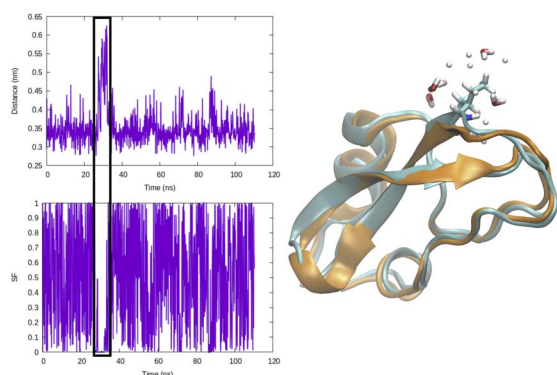


Fig. 3 Temporal evolution of H-bond distance of Ile18–Tyr35 and corresponding switching function ($r_{ij} = 0.32$, $r_0 = 0.06$, $n = 6$, $m = 12$) in BPTI. CV values corresponding broken H-bond conformations is highlighted in black square. Snapshot corresponding broken conformation highlights how breaking of H-bond leads to solvent exposure (highlighting water molecules within 3 radius of Ile18) of Ile18–NH.

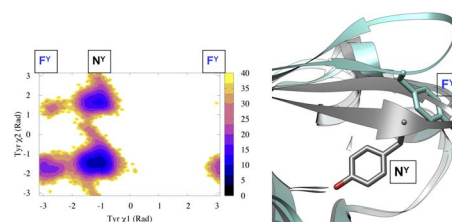


Fig. 4 Projection of χ_1 and χ_2 angle of Tyr77 in apo plasmepsin-II (PDB: 1LF4) during unbiased ~ 600 ns long MD simulation (left panel). Tyr77 predominantly remains in the normal state (N^Y) with rare sampling of the flipped state (F^Y). Conformational snapshots corresponding normal and flipped states with relative position of flap (β hairpin structure) is also highlighted (right panel). In apo plasmepsin-II, rotation of Tyr77 along χ_1 and χ_2 angles dictates the extent of flap opening which governs substrate entry and catalytic activity.



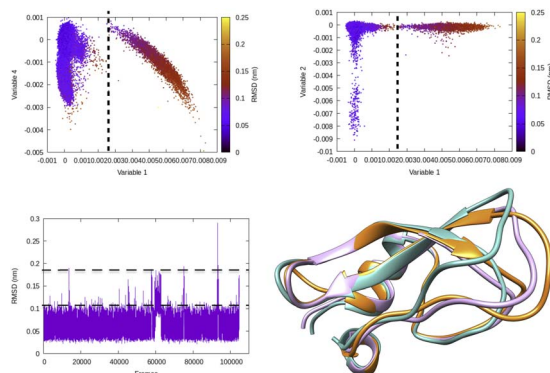


Fig. 5 Projection of different variables (more precisely time independent components as described in Fig. 8) which captures transient loop opening BPT1 as a function of α RMSD of residue 6–56 in BPT1. It can be easily seen that a higher RMSD value separates two metastable states along variable 1. Snapshots corresponding crystal (orange), RMSD ~ 0.15 nm (pink) and RMSD ~ 0.25 nm (sea green) highlights differences in loop conformation. RMSD calculation was performed on M1 state¹³ of 1 ms long conventional MD simulation performed by D. E. Shaw Research.¹⁴

latter. In case of linear transformation, the transformed data is a linear combination of original variables. Whereas non-linear transformations are more complex than that. In this section, I will underscore some of abstract CVs (linear: principal component analysis & independent component analysis; non-linear: kernel trick, diffusion map, t-SNE) that have been regularly used to capture low-dimensional spatiotemporal representation from high dimensional dataset generated by MD simulation. I will further discuss the probabilistic interpretation of variational autoencoders which has been applied to capture compressed representation of temporal variables from molecular simulation. Finally, I will highlight the application of binary classifiers in context of classifying metastable states and drive enhanced sampling simulations to capture state transitions. I will further list the softwares/tools/codes that can be used to generate abstract CVs in context of biomolecular simulation.

2.2.1 Principal component analysis. Principal component analysis (PCA) is an unsupervised dimensionality reduction method which transforms a set of variables $r_1, r_2, r_3, \dots, r_N$ (where $\mathbf{r} = [r_1, r_2, \dots, r_N]^T$) to low dimensional representations y_i which captures as much of the variations as possible. It has been widely used by biomolecular simulation community to capture molecular motions with largest amplitude (high variance).^{12,19} Let $\mathbf{r} \in \mathbb{R}^D$ be a vector of geometric CVs *e.g.* dihedral angles or distances (assuming \mathbf{r} is mean free). First principal component (PC1) of the dataset $r_1, r_2, r_3, \dots, r_N$ is the weighted linear combination of the features that captures the largest variance:

$$y_1 = w_{11}r_1 + w_{21}r_2 + w_{31}r_3 + \dots + w_{N1}r_N \quad (4)$$

where w_1 is the vectors of weights or coefficients (PCA aims to find w which maximises the variance) with elements $w_{11}, w_{21},$

\dots, w_{N1} . The elements are normalised *i.e.* $\sum_{i=1}^N w_{i1}^2 = 1$. The process of generating PCs can be summarised as follows:

- (a) generate a mean free version of the input data *i.e.* a vector of geometric CVs
- (b) compute eigenvectors and eigenvalues from the covariance matrix

$$\mathbf{C} = \frac{1}{N} \bar{\mathbf{r}}^T \bar{\mathbf{r}} \quad (5)$$

where $\bar{\mathbf{r}}$ is the mean free version of \mathbf{r} .

- (c) sort the eigenvalues in descending order and retain k (k is the new subspace $k < D$) eigenvectors that corresponds to k largest eigenvalues. The eigenvector corresponds to the largest eigenvalue capture the greatest variability in the data.

- (d) construct a projection matrix \mathbf{w} with elements w_{1j}, w_{2j}, w_{Nj} (where $\mathbf{w} = [w_{1j}, w_{2j}, \dots, w_{Nj}]^T$).

- (e) generate the PCs, \mathbf{y} by transforming the original geometric vectors \mathbf{r} *via* elements of the projection matrix \mathbf{w} .

Principal components are uncorrelated to each other (as eigenvectors are orthogonal to one another) and sorted by their variance ($\text{PC1} > \text{PC2} > \text{PC3} \dots$). PCs extracted from MD dataset defines direction in feature space along maximal variance (largest amplitude/fluctuations). PCA has been applied routinely on periodic (dihedral angles) and non-periodic (CA atomic positions, RMSDs, distances) degrees of freedom (Fig. 6). The coefficients w corresponding each variables allows

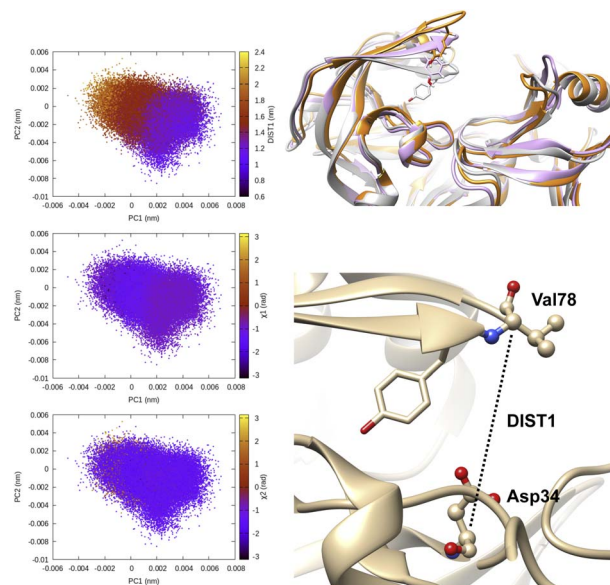


Fig. 6 Principal components project along α - α distance (DIST1) between Asp34-Val78 and χ_1 and χ_2 angles of flap tip Tyr77 in plasmeprin-II. One can see that PCs managed to capture the extent of flap opening which is quantified by DIST1 (DIST1 > 1.8 nm is open conformation). In this case PCA analysis was performed on the α atoms of the protein hence it is unable to capture rotational degrees of freedom along χ_1 and χ_2 angles. Conformational snapshots corresponding crystal (grey, DIST2 = 1.2 nm), semi-open (magenta, DIST1 = 1.6 nm) and open conformation (orange, DIST1 = 2.0 nm) are also highlighted.



incorporating PCs as CVs within enhanced sampling protocols *e.g.* metadynamics and its variants. PCA reduced dimensions can be also used to construct free energy surface of biomolecules. Despite its success, application of PCA in biomolecular simulation has been a subject of controversy. It has been argued that PCA can't capture the slow conformational degrees of freedom within time scales accessible to MD simulations. However, the author believes that PCA on a carefully chosen subset of atomic co-ordinates can capture slow conformational changes in biomolecules.

Softwares: MDAnalysis, Plumed (with capability of using PCs as CVs in metadynamics), MSMBuilder,²⁰ PyEMMA,²¹ pytraj, scikit-learn (can't be directly used on MD trajectories but can be interfaced with MD post-processing tools *e.g.* MDTraj).

2.2.2 Independent component analysis. Independent component analysis is a dimensionality reduction method which transforms a set of vectors (*e.g.* distances, RMSDs, dihedral angles *etc.*) into maximally independent linear combinations (independent components). Imagine $\mathbf{y}(t) = [y_1(t), \dots, y_N(t)]^T$ is a linear mixture of high-dimensional mean-free data $\mathbf{r}(t) = [r_1(t), \dots, r_N(t)]^T$ such that $\mathbf{y}(t) = \mathbf{A}\mathbf{r}(t)$ where \mathbf{A} is the square mixing matrix and the components of $\mathbf{y}(t)$ are mutually independent. Two components can be said to be independent if their joint distribution is equal to the product of their marginals (a corresponding measure of independence is Kullback–Leibler divergence):

$$p(y_1(t), y_2(t)) = p(y_1(t))p(y_2(t)) \quad (6)$$

where $p(y_1(t), y_2(t))$ are the joint probability distributions and $p(y_1(t))p(y_2(t))$ is the marginal along two components $y_1(t)$ and $y_2(t)$. However the aforementioned definition of independence has a drawback: let's consider a signal without time auto-correlation and a second signal which is equal to the first signal but shifted in time (often known as signal with time-delay). If one applies eqn (6) the two signals will appear to be mutually independent. The time-delayed signal can be called statistically independent if all the time-delayed correlations are zero (second-order ICA).²² In biomolecular simulation second-order ICA is known as time-lagged independent component analysis (TICA).^{23,24} The first step of second-order ICA/TICA is to introduce a time-lagged (or time-delayed) correlation matrices of the input variables $\mathbf{r}(t)$:

$$\mathbf{C}^r(\tau) = \langle \mathbf{r}(t)\mathbf{r}(t + \tau)^T \rangle \quad (7)$$

where τ is the lag-time or time delay between two signals. The entries of $\mathbf{C}^r(\tau)$ (should be diagonal for all τ) are denoted as $C_{ij}^r(\tau)$. The common practice is to express $\mathbf{C}^r(\tau)$ as a symmetrized version of correlation matrices:

$$\mathbf{C}^r(\tau) = \frac{1}{2} [\langle \mathbf{r}(t)\mathbf{r}(t + \tau)^T \rangle + \langle \mathbf{r}(t + \tau)\mathbf{r}(t)^T \rangle] \quad (8)$$

Symmetrization is a mathematical trick which allows applying the algorithm to the reversible dynamics (as $\langle \mathbf{r}(t + \tau)\mathbf{r}(t)^T \rangle = \langle \mathbf{r}(t)\mathbf{r}(t - \tau)^T \rangle$). It also makes sure that the eigenvalues

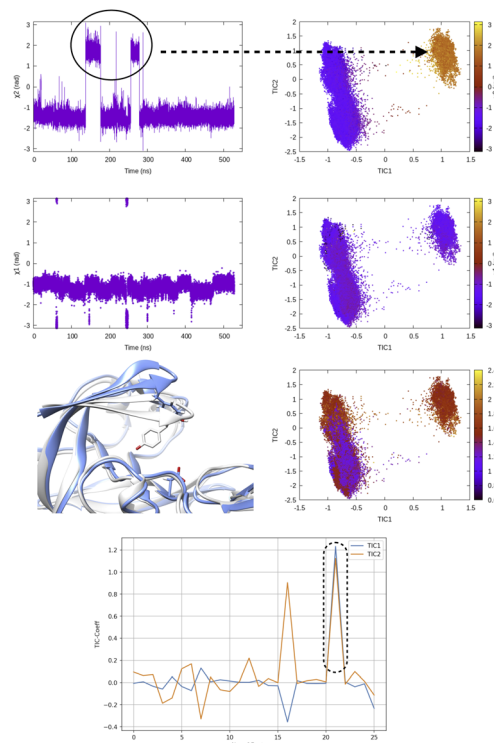


Fig. 7 Evolution of χ_1 and χ_2 angles of Tyr77 of plasmepsin-II during MD simulation. TICA analysis (lag time 10 ns) was performed on χ_1 and χ_2 angles (sin/cos modification) of residue 74–84. Projection of TIC1 and TIC2 along χ_1 and χ_2 angles of Tyr77 shows that TIC1 managed to capture the slowest degree of motion which is the χ_2 rotation of Tyr77. Further, projection along TICs failed to separate fluctuation along DIST1 as it is defined by distance between backbone $C\alpha$ atoms. Fluctuation of TIC coefficients (eigenvalues) also highlights feature no 21 (sin χ_2 of Tyr77) as the dominant feature among 25 features (Table 1 in ESI†). Snapshots corresponding crystal (grey) and flipped ($\chi_2 \sim 2$ rad) are also highlighted.

are real and two eigenvectors that comes from distinct eigenvalues are orthogonal. Finally the TICA problem can be formulated a generalised eigenvalue problem (can be solved by second-order blind source separation algorithms *e.g.* AMUSE):

$$\mathbf{C}^r(\tau)\mathbf{A} = \mathbf{C}^r(0)\mathbf{A}\mathbf{\Lambda} \quad (9)$$

where $\mathbf{A} = (a_1, \dots, a_N)$ is the orthogonal matrix of generalised eigenvectors and $\mathbf{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_N)$ is the diagonal matrix of eigenvalues (the eigenvalues $\lambda_1, \dots, \lambda_N$ are associated with eigenvectors of a_1, \dots, a_N). \mathbf{A} contains the independent components (ICs) the original data $\mathbf{r}(t)$ can be projected on the TICA space as: $\mathbf{y}(t) = \mathbf{A}\mathbf{r}(t)$. The scalar components of λ_i captures the magnitude of the auto-covariance where smaller λ_i captures largest auto-covariance:

$$\lambda_1 < \lambda_2 < \dots < \lambda_N; \lambda_1 \text{ captures largest auto-covariance and so on} \quad (10)$$

The eigenvalues, λ_i are associated with the relaxation time-scales of a biomolecular process by:



$$t_i = -\frac{\tau}{\ln|\lambda_i|} \quad (11)$$

TICA has been routinely used on temporal data from MD simulation to identify slow conformational degrees of freedom such as flipping of side-chain dihedral angle (Fig. 7), transient exposure of protein interior which leads of solvent exposure (Fig. 8) and drive enhanced sampling simulations.²⁵ Recently Schultze and Grubmüller compared projection of TICs between high-dimensional data from protein dynamics and random walk.²⁶ However the authors disregarded any discussion surrounding the choice of input features and its effect on TICA. Fig. 7 highlights the importance of input features in describing conformational dynamics along TIC space. Lack of separation aka population gap along TIC space in ref. 26 (see Fig. 7, left panel) is a sign that either TICA analysis was performed on a poorly chosen feature space or the sampling was not sufficient enough to capture slow conformational dynamics in proteins.

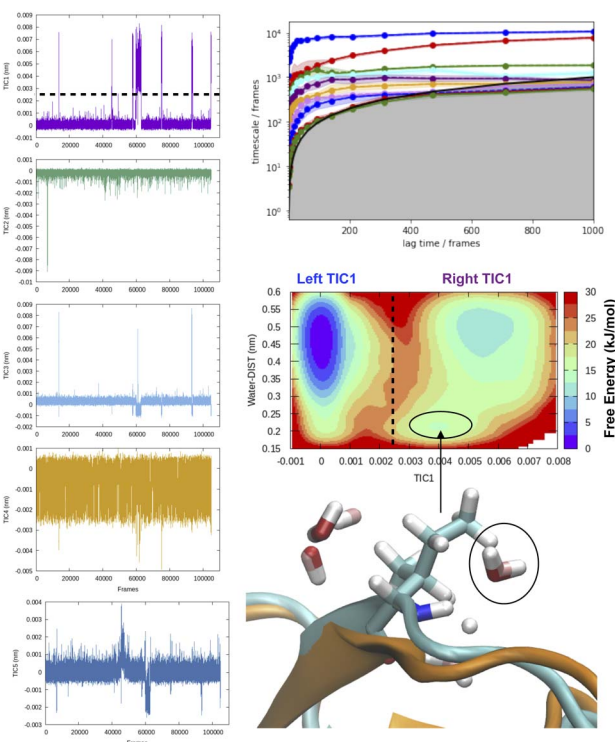


Fig. 8 Time-series projection of first 5 TICs and corresponding implied timescales showing how first few TICs captured slow conformational degrees of freedom (transient loop opening as depicted in Fig. 5) in BPTI. TICA analysis was performed on C α atomic positions of residues 6–56 using a lag time of 500 frames (all calculations were performed on M1 state as described by ref. 13). Transient loop opening in BPTI leads to breaking of H-bond interaction between Ile18–NH and Tyr35–O which makes the amide of Ile18 solvent exposed and enables hydrogen–deuterium exchange.²⁹ Projection of water distance from Ile18–NH shows how TIC1 captures a metastable solvent exposed basin where a water molecule comes within 0.25 nm of amide hydrogen. Snapshot corresponding solvent exposed conformation highlights the relative position of backbone amide and closest (with 0.3 nm) waters.

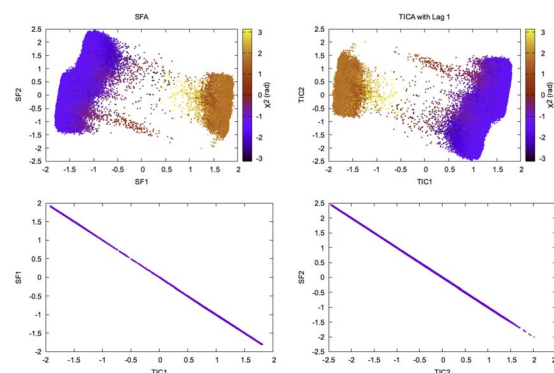


Fig. 9 Projection of slow features and TICs with lag time 1 along χ_2 angle of Tyr77 shows similarity of TICA and SFA in separating slow conformational degree of freedom. The squared Pearson correlation coefficient, R^2 between SFA and TICA with lag time 1 is 1.00. This example is a demonstration of mathematical concepts described in ref. 27 in context of biomolecular simulation. It further raises possibilities of combining the two algorithms and use it to develop Markov state model.

Wiskott and co-workers²⁷ have shown how the objective function of TICA with lag time 1 is formally equivalent to slow feature analysis (SFA) which captures slowly varying features from high-dimensional data (Fig. 9). SFA and TICA are based on two different principles: slowness and statistical independence. However, the similarity between SFA and TICA (second order ICA with lag-time 1) opens up possibilities to combine these two algorithms²⁸ for capturing slowly varying statistically independent components from high-dimensional temporal data generated by MD simulation.

Softwares: MSMBuilder,²⁰ PyEMMA²¹ and deeptime³⁰ comes with built in TICA functionality. A scikit-learn style implementation of SFA can be found here: <https://github.com/wiskott-lab/sklearn-sfa>.

2.2.3 Kernel trick. MD simulation generates complex and non-linear representation of biomolecular dynamics. Kernel trick is a mathematical transformation which maps the original data into a higher dimensional feature space which is then used to find linear projections using PCA (kernel PCA) or tICA (kernel TICA³¹). For data points r_i, r_j in the input space N , kernel function $k(r_i, r_j)$ generates modified inner products which maps $N \rightarrow Z$

$$k(r_i, r_j) = \langle \phi(r_i), \phi(r_j) \rangle_Z \quad (12)$$

where ϕ is the mapping function and $\langle \cdot, \cdot \rangle_Z$ must be proper inner product. One doesn't need to explicitly compute ϕ as the kernel matrix $k(r_i, r_j)$ (modified inner product) can be easily computed by a variety of functions:

(a) Polynomial kernel: $k(r_i, r_j) = (\gamma \cdot r_i^T r_j + c)^d$, $\gamma > 0$

(b) Sigmoid kernel: $k(r_i, r_j) = \tanh(\gamma \cdot r_i^T r_j + r)$.

(c) Gaussian kernel: $k(r_i, r_j) = \exp(-\gamma \cdot \|r_i - r_j\|^2)$, $\gamma > 0$. where r, d and γ are kernel hyper-parameters. A drawback of this approach is that when we map the data into higher dimensions, we may overfit the model. Hence the choice of right kernel functions are of utmost importance. Schwartz and colleagues



used kernel PCA (using polynomial kernel function) to identify CVs in lactate dehydrogenase. Further, kernel TICA has been used to identify CVs which captures folded to unfolded transition in small folded proteins.³²

Softwares: scikit-learn and MODE-TASK³³ have in-built kernel PCA functionality. MSMBuilder has a built in kernel TICA algorithm which can be applied on high-dimensional features from MD dataset for discovering low-dimensional representations aka CVs.

2.2.4 Diffusion map. Diffusion map is a non-linear dimensionality reduction technique. It combines the concept of random walk Markov chain and diffusion process by projecting the input data in a low-dimensional space where the distance (e.g. Euclidean) between data points resembles the diffusion distance in the original high-dimensional space. When applied to MD dataset, diffusion map generated vector space (CVs) are constructed in such a way that conformations which are kinetically closed are placed together whereas kinetically distant (separated by high free-energy barrier) conformations are placed far apart.

In a random walk Markov model the connectivity between data points r_i and r_j is defined as the probability of jumping from state r_i to r_j in one step of random walk. The connectivity can be also expressed as a normalised kernel function $k(r_i, r_j)$ (similar to the functional form of Gaussian kernel. One can choose other measures of distances as kernel functions e.g. Mahalanobis distance, Jensen–Shannon divergence *etc.*) which measures similarity between data points:

$$k(r_i, r_j) = \exp\left(-\gamma\|r_i - r_j\|^2\right), \gamma > 0$$

$$= \exp\left(-\frac{\|r_i - r_j\|^2}{\alpha}\right) \quad (13)$$

where the value of α decides the size of neighbourhood. It can also be seen as a hyper-parameter which chooses the conformations that are kinetically closed. In practice conformations closer than the value of α are relevant for $k_{ij} = k(r_i, r_j)$ as the contribution from distant conformations decays exponentially. The diffusion kernel satisfies the following two properties:

(a) k is symmetric: $k(r_i, r_j) = k(r_j, r_i)$. This allows spectral analysis of the distance matrix k_{ij}

(b) $k(r_i, r_j) \geq 0$.

The transition probability $p(r_i, r_j)$ can be expressed as:

$$p(r_i, r_j) = \frac{1}{D} k(r_i, r_j) \quad (14)$$

where D is the normalisation factor. The normalised transition matrix \mathbf{P} with $P_{ij} = p(r_i, r_j)$ (where $\sum_j P_{ij} = 1$) allows spectral (in other words eigen) decomposition:

$$\mathbf{P}\mathbf{A} = \mathbf{A}\mathbf{\Lambda} \quad (15)$$

where a_1, \dots, a_N ($\mathbf{A} = (a_1, \dots, a_N)$) are real valued eigenvectors corresponding to eigenvalues $\lambda_1, \dots, \lambda_N$ ($\mathbf{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_N)$) where $1 = \lambda_1 > \lambda_2 \geq \lambda_N$. Similarly to TICA, the eigenvalues generated from normalised transition matrix \mathbf{P} shows spectral gap which can be expressed as the difference between two

largest eigenvalues $\lambda_1 - \lambda_2$ or more generally $1 - \max\{|\lambda_2|, |\lambda_N|\}$. Diffusion map at time t can be approximated as mapping between the original space and the latent space of first k eigenvectors:

$$\mathbf{A}(\mathbf{r}) = (\lambda_1^t a_1(r), \lambda_2^t a_2(r), \dots, \lambda_k^t a_k(r)) \quad (16)$$

The diffusion distance (analogous to the functional form of kinetic distance as proposed by Noe and Clementi ref. 34) at time t can be expressed as a function of eigenvectors. Kevrekidis and coworkers³⁵ has shown that the diffusion distance can be approximated as Euclidean distance on the diffusion map space:

$$D_t^2(r_i, r_j) = \sum_l \lambda_l^{2t} (a_l(r_i) - a_l(r_j))^2 = \|\mathbf{A}_t(r_i) - \mathbf{A}_t(r_j)\|^2 \quad (17)$$

Eqn (16) provides justification of using Euclidean distance in the diffusion map space for clustering. Diffusion distance can also be interpreted as a measure of how kinetically close are the two conformations. The distance is small if there are multiple high probability transition pathways between two conformations. Diffusion map has been applied in biomolecular simulation to capture slow transitions and guide enhanced sampling simulations.^{36–38} Recently the functional form of eqn (16) in combination with maximum entropy based CV selection method SGOOP was used to capture kinetically relevant low dimensional projection in small peptides.³⁹

Softwares: MDAnalysis has an integrated diffusion map module which can be applied on selected features from MD simulation.

2.2.5 t-distributed stochastic neighbour embedding. t-Distributed stochastic neighbour embedding (t-SNE) is a non-linear dimensionality reduction method which performs embedding of high-dimensional data to low-dimensional space such that neighbouring data-points are assigned highest probability while distant points are assigned lower probability.^{40,41} For a high-dimensional space r_1, \dots, r_N the distance between r_j and r_i can be expressed as:

$$p_{j|i} = \frac{\exp\left(-\|r_i - r_j\|^2 / 2\sigma_i^2\right)}{\sum_{k \neq i} \exp\left(-\|r_i - r_k\|^2 / 2\sigma_i^2\right)} \quad (18)$$

where $p_{j|i}$ is the conditional probability which captures the similarity between data points r_j with r_i such that $p_{i|i} = 0$ and $\sum_j p_{j|i} = 1$. The functional form of eqn (18) is similar to Gaussian kernel in eqn (13). The conditional probability can be transformed into a joint probability distribution of symmetrized matrix P_{ij} :

$$P_{ij} = \frac{p_{j|i} + p_{i|j}}{2N} \quad (19)$$

where $P_{ij} = P_{ji}$, $P_{ii} = 0$ and $\sum_{i,j} P_{ij} = 1$.

The Gaussian bandwidth σ_i has been set in such a way that the perplexity of the conditional distribution equals to the perplexity provided by the user (⁴⁰ highlighted that there is a monotonically increasing relationship between perplexity and bandwidth). The perplexity is defined as:

$$\text{Perp}(p_i) = 2^{H(p_i)} \quad (20)$$



where p_i is the conditional probability distribution over all the data-points given r_i and $H(p_i)$ is the Shannon entropy $H(p_i) = -\sum_j p_{ji} \log_2 p_{ji}$. Perplexity can be thought as a measure of nearest neighbours and the choice of perplexity heavily determines the outcome of t-SNE (Fig. 10).

t-SNE aims to learn a low dimensional manifold $r'_1 \dots r'_N$ in such a way that the new conditional probability p'_{ji} reflects similarity with p_{ji} . p'_{ji} can be expressed as:

$$p'_{ji} = \frac{(1 + \|r'_i - r'_j\|^2)^{-1}}{\sum_{l \neq k} (1 + \|r'_k - r'_l\|^2)^{-1}}, \quad p'_{ii} = 0 \quad (21)$$

The distance based metrics $(1 + \|r'_i - r'_j\|^2)^{-1}$ is a heavily tailed distribution. In t-SNE, Student t-distribution has been used to measure the similarity between low-dimensional data-points so that the points that are far apart have p'_{ji} which are invariant of perturbation. The algorithm aims to minimise the Kullback–Leibler divergence (other measures of similarity e.g. Jensen–Shannon divergence, Bhattacharya distance can be explored to improve the embedding result) by comparing joint probability distributions P' (low-dimensional) and P (high-dimensional):

$$\text{KL}(P||P') = \sum_{i \neq j} p_{ij} \log \frac{p_{ij}}{q_{ij}} \quad (22)$$

The gradient descent algorithm is usually effective for small datasets but performs poorly in case of larger datasets. Recently a time-lagged version of t-SNE⁴² was proposed using inspiration from time-lagged ICA (TICA). However, the stochastic nature (mainly due to perplexity) of t-SNE (as shown in Fig. 10 and discussed in ref. 43) prevents its use as a meaningful CV for reconstructing trustworthy free energy surface and calculating kinetics from MD simulation. Clustering on the top of t-SNE embedding can also produce artificial clusters which might

trick the user of thinking that it discovers new metastable states but in reality they belong to the same free energy basin.

Softwares: time-lagged version of t-SNE can be accessed here: <https://github.com/spiwokv/tltsne>. Popular machine learning package scikit-learn also has a t-SNE module: <https://scikit-learn.org/stable/modules/generated/sklearn.manifold.TSNE.html>. MODE-TASK <https://github.com/RUBI-ZA/MODE-TASK> a python toolkit for analysing MD simulation has an integrated t-SNE functionality. Tensorboard embedding projector <https://projector.tensorflow.org> has a graphical user interface to perform t-SNE.

In recent years, several other dimensionality reduction methods such as spectral gap optimization of order parameters (SGOOP),⁴⁴ isomap⁴⁵ (<https://scikit-learn.org/stable/modules/generated/sklearn.manifold.Isomap.html>), dynamic mode decomposition (DMD)⁴⁶ (<https://mathlab.github.io/PyDMD/>, see ref. 47 for similarity between DMD and TICA and their variants), multi-dimensional scaling (MDS),⁴⁵ UMAP^{48,49} (<https://umap-learn.readthedocs.io/en/latest/>, Fig. 11), iVIS⁵⁰ (<https://bering-ivis.readthedocs.io/en/latest/>) which are similar to some of the aforementioned algorithms, were applied on temporal data from MD simulation. Further, sparse regression based dimensionality reduction method sparse identification of nonlinear dynamical systems (SINDy)⁵¹ can possibly be applied on temporal data from MD simulation to discover linear combinations of features which best captures conformational dynamics in macromolecules. Recent review by Glielmo *et al.*⁵² summarises mathematical concepts, strengths and limitations as well as applicability of few such methods in analysing high-dimensional data from MD simulations.

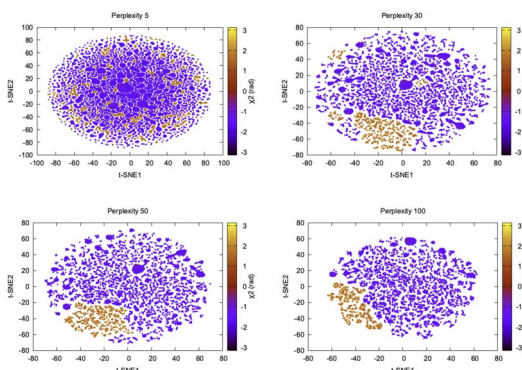


Fig. 10 Projection of first two t-SNE components along χ_2 angle of Tyr77 in plasmepsin-II shows how the embedding changes with perplexity. Further t-SNE performed on the high dimensional torsional features (Table 1 in ESI†) didn't manage to separate flipping along χ_2 which is a slow degree of freedom. It further shows that clustering on top of t-SNE reduced dimensions will generate artificial clusters which in reality belong to same metastable state.

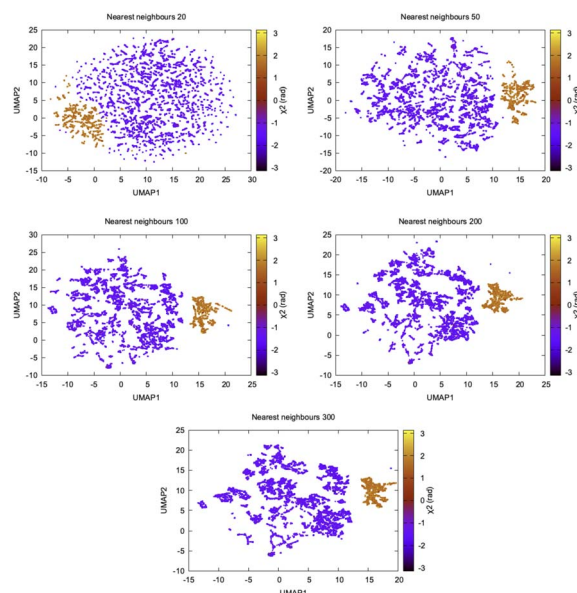


Fig. 11 UMAP components projected along χ_2 of Tyr77 highlights how the shape of UMAP vary as a function of nearest neighbour hyperparameter. UMAP converges with nearest neighbour 200. UMAP does a better job compared to t-SNE in separating conformational space along χ_2 . I believe clustering on top of UMAP reduced dimensions might be useful in understanding conformational heterogeneity within a broad metastable basin.



2.2.6 Learning co-ordinates for dynamics using variational autoencoder. Variational autoencoder (VAE) is a dimensionality reduction method which takes high dimensional data (e.g. multi-variate geometric or abstract CVs) as inputs and learns latent (compressed) representations that captures minimal essential information necessary to describe the dynamics of the system. VAE is a type of autoencoder whose principals are deeply rooted in variational statistics.^{53,54} In this article we will use the language of probability theory to describe VAE and try to establish a connection between free energy and the loss function of VAE.

VAE takes high dimensional vector $\mathbf{r} = [r_1, r_2, \dots, r_N]^T$ as inputs. Each components of \mathbf{r} are probability distributions along geometric or abstract CVs. The encoder part of the VAE encodes the high dimensional data into the latent variables \mathbf{z} . The joint probability distribution of input and latent variables, $p(\mathbf{r}, \mathbf{z})$ can be expressed as:

$$p(\mathbf{r}, \mathbf{z}) = p(\mathbf{r}|\mathbf{z})p(\mathbf{z}) \quad (23)$$

The generative process can be expressed as:

(a) sampling latent variables, \mathbf{z}_i from the prior distribution $p(\mathbf{z})$ and

(b) sampling of data point r_i from the likelihood $p(\mathbf{r}|\mathbf{z})$ which is conditional on the latent variables \mathbf{z} .

The goal of VAE is to infer good approximation of the latent variables (\mathbf{z}) given input data (\mathbf{r}) which is equivalent to calculate the posterior distribution $p(\mathbf{z}|\mathbf{r})$ as following:

$$p(\mathbf{z}|\mathbf{r}) = \frac{p(\mathbf{r}|\mathbf{z})p(\mathbf{z})}{p(\mathbf{r})} = \frac{p(\mathbf{r}, \mathbf{z})}{p(\mathbf{r})} \quad (24)$$

where

$$p(\mathbf{r}) = \int p(\mathbf{r}|\mathbf{z})p(\mathbf{z})d\mathbf{z} \quad (25)$$

Estimation of $p(\mathbf{r})$ is a computationally expensive process as it requires integrating over all the possible values of \mathbf{z} . VAE approximates the posterior distribution $p(\mathbf{z}|\mathbf{r})$ by a new distribution $q(\mathbf{z}|\mathbf{r})$ (tractable distribution). If $q_\lambda(\mathbf{z}|\mathbf{r})$ is similar to $p(\mathbf{z}|\mathbf{r})$ then we can use it to approximate $p(\mathbf{r})$. λ is the hyper-parameter which indicates type of distribution. If $q(\mathbf{z}|\mathbf{r})$ a Gaussian distribution then $\lambda_{ri} = (\mu_{ri}, \sigma_{ri}^2)$ where μ is the mean and σ^2 is the variance of latent variables of each input component.

Now we want to measure how well $q_\lambda(\mathbf{z}|\mathbf{r})$ approximates $p(\mathbf{z}|\mathbf{r})$. VAE uses Kullback–Leibler (KL) divergence to measure information loss between two probability densities as described below:

$$\begin{aligned} \text{KL}(q_\lambda(\mathbf{z}|\mathbf{r})||p(\mathbf{z}|\mathbf{r})) &= \mathbb{E}_q \left[\log \frac{q_\lambda(\mathbf{z}|\mathbf{r})}{p(\mathbf{z}|\mathbf{r})} \right] \\ &= \mathbb{E}_q [\log q_\lambda(\mathbf{z}|\mathbf{r}) - \log p(\mathbf{z}|\mathbf{r})] \\ &= \mathbb{E}_q [\log q_\lambda(\mathbf{z}|\mathbf{r})] - \mathbb{E}_q [\log p(\mathbf{r}, \mathbf{z})] + \log p(\mathbf{r}); \text{ taking log of eqn 24} \\ &= -(\mathbb{E}_q [\log p(\mathbf{r}, \mathbf{z})] - \mathbb{E}_q [\log q_\lambda(\mathbf{z}|\mathbf{r})]) + \log p(\mathbf{r}) \\ &= -\text{ELBO}(\lambda) + \log p(\mathbf{r}) \end{aligned} \quad (26)$$

We can reformulate eqn (26) as follows:

$$\log p(\mathbf{r}) = \text{ELBO}(\lambda) + \text{KL}(q_\lambda(\mathbf{z}|\mathbf{r})||p(\mathbf{z}|\mathbf{r})) \quad (27)$$

From eqn (27) we can see that minimising KL divergence (KL divergence is always greater or equal to zero) is equivalent to maximising Evidence Lower Bound (ELBO). It allows us to bypass the hard task of minimising KL divergence between the approximate $q(\mathbf{z}|\mathbf{r})$ and true posterior $p(\mathbf{z}|\mathbf{r})$, instead we maximise ELBO. ELBO term can be further expressed as follows:

$$\begin{aligned} \text{ELBO}(\lambda) &= \mathbb{E}_q [\log p(\mathbf{r}, \mathbf{z})] - \mathbb{E}_q [\log q_\lambda(\mathbf{z}|\mathbf{r})] \\ &= \mathbb{E}_q [\log p(\mathbf{r}|\mathbf{z})] + \mathbb{E}_q [\log p(\mathbf{z})] - \mathbb{E}_q [\log q_\lambda(\mathbf{z}|\mathbf{r})]; \\ &\quad \text{log transformation of eqn 23} \\ &= \mathbb{E}_q [\log p(\mathbf{r}|\mathbf{z})] - \text{KL}(q_\lambda(\mathbf{z}|\mathbf{r})||p(\mathbf{z})) \end{aligned} \quad (28)$$

Eqn (28) is known the VAE loss function (L):

$$L = \mathbb{E}_q [\log p(\mathbf{r}|\mathbf{z})] - \text{KL}(q_\lambda(\mathbf{z}|\mathbf{r})||p(\mathbf{z})) \leq \log p(\mathbf{r}) \quad (29)$$

In biomolecular simulation deep neural network (DNN) encodes input data and computes λ which approximates $q_w(\mathbf{z}|\mathbf{r}, \lambda)$. The decoder takes $p(\mathbf{z})$ as input and maps into original distribution $p_w(\mathbf{r}|\mathbf{z})$. w and w' acts as a neural network weights. w transforms the input data within the neural network hidden layers and the resulting latent variable is a combination of linear transformations that are modified by non-linear activation functions. Weights are learnable parameters during VAE training. The value of weight dictates the importance of a variable. A higher weight value corresponding an input component indicates that it will have a significant influence on the output.

Mathematically speaking, the reason behind VAE's popularity in biomolecular simulation community is due to the interconnectivity between loss function and variational free energy (F):

$$\begin{aligned} \text{KL}(q_\lambda(\mathbf{z}|\mathbf{r})||p(\mathbf{z}|\mathbf{r})) &= \int q_\lambda(\mathbf{z}|\mathbf{r}) \log \frac{q_\lambda(\mathbf{z}|\mathbf{r})}{p(\mathbf{z}|\mathbf{r})} d\mathbf{z} \\ &= \int q_\lambda(\mathbf{z}|\mathbf{r}) \log \frac{q_\lambda(\mathbf{z}|\mathbf{r})p(\mathbf{r})}{p(\mathbf{r}, \mathbf{z})} d\mathbf{z} \\ &= \int q_\lambda(\mathbf{z}|\mathbf{r}) \log \frac{q_\lambda(\mathbf{z}|\mathbf{r})}{p(\mathbf{r}, \mathbf{z})} d\mathbf{z} + \int q_\lambda(\mathbf{z}|\mathbf{r}) \log p(\mathbf{r}) d\mathbf{z} \\ &= \int q_\lambda(\mathbf{z}|\mathbf{r}) \log \frac{q_\lambda(\mathbf{z}|\mathbf{r})}{p(\mathbf{r}, \mathbf{z})} d\mathbf{z} + \log p(\mathbf{r}); \\ \text{as } \int q_\lambda(\mathbf{z}|\mathbf{r}) d\mathbf{z} &= 1 \\ &= - \int q_\lambda(\mathbf{z}|\mathbf{r}) \log \frac{p(\mathbf{r}, \mathbf{z})}{q_\lambda(\mathbf{z}|\mathbf{r})} d\mathbf{z} + \log p(\mathbf{r}) \\ &= -F + \ln p(\mathbf{r}) \end{aligned} \quad (30)$$

where

$$\int q_\lambda(\mathbf{z}|\mathbf{r}) \log p(\mathbf{r}) d\mathbf{z} = 1 \quad (31)$$

By comparing eqn (26) with (31) we can conclude that the variational free energy (F) is equal to the ELBO(λ). Thus



minimising KL divergence is equivalent to maximising the variational free energy (F).

Fig. 12 shows a typical architecture of VAE in context of biomolecular simulation. The encoder part of the VAE acts as a non-linear dimensionality reduction of input \mathbf{r} . Whereas the decoder DNN acts an input reconstruction. VAE has been primarily used a dimensionality reduction method to compress multi-variate probability distributions. The weights of the encoder layer which maps the input data onto the latent layer has been further used to drive enhanced sampling simulations such as metadynamics.⁵⁵ As the latent layer of VAE represents a non-linear combination of input data hence it reduces the need of multiple CVs as inputs in metadynamics (traditionally metadynamics is limited to a maximum of two CVs). However in terms of exploration of the conformational space, metadynamics bias applied to multiple CVs using parallel-bias metadynamics framework will be equivalent to using VAE's latent variable as CV in a 1D metadynamics. Several different flavours of VAEs have been developed which mainly differs in their architecture (e.g. β -VAE which uses an additional hyperparameter β to learn disentangled latent variables by controlling the KL divergence⁵⁶) and types of inputs (e.g. VAE-SNE which takes output of tSNE as inputs, DMD-VAE which can take multiple dynamic modes as inputs, SINDy-VAE which takes outputs from SINDy algorithm as inputs). However, combining autoencoder based machine learning architecture with enhanced sampling algorithms or MD codes is a difficult task in

practice. Recently, Chen and co-workers⁵⁷ proposed a user-friendly tool called MLCV which uses an autoencoder architecture to transform arbitrarily defined CVs to autoencoder learned CVs and integrate it with MD engines for enhanced sampling calculations. In future, MLCV can be extended by adding different variants of VAE to capture conformational heterogeneity and drive enhanced sampling calculations.

Softwares & implementations: VAEs for post-processing high-dimensional time-series data generated from MD simulation can be implemented using popular machine learning libraries e.g. PyTorch, Tensorflow, Keras *etc.* Recently Pande and co-workers⁵⁸ (variational dynamic encoder: <https://github.com/msmbuilder/vde>) as well as Tiwary and colleagues⁵⁹ (RAVE: <https://github.com/tiwarylab/RAVE>) applied VAE to compress high-dimensional temporal data and able to capture non-linear dynamics in context of protein folding, protein-ligand binding/unbinding *etc.* A collection of prebuilt VAE architectures implemented with PyTorch (<https://github.com/AntixK/PyTorch-VAE>) opens up the possibility to evaluate VAE variants in context of MD simulation.

2.2.7 Classifiers as CVs. Classifiers are supervised learning algorithms which uses feature vectors and feature values to map input data into specific categories.⁶⁰ In MD simulation, the classifiers are used as CVs to recognise different metastable states (can also be used to identify features that distinguishes wild and mutant variants) based on a common set of feature values that distinguishes them (Fig. 13). Recently different classes of classifiers (e.g. support vector machine, logistic regression, artificial neural network, linear discriminant analysis *etc.*) have been applied to categorise different metastable states and drive enhanced sampling simulations.^{61,62} In this article, I will discuss support vector machine (SVM) and linear discriminant analysis (LDA) and their applications in enhanced sampling.

2.2.7.1 Support vector machine. Let A and B are two metastable states which can be represented by unique points in space where each point belongs to high-dimensional feature vector $\mathbf{r} = [r_1, \dots, r_N]$. The goal of SVM algorithm is to find a separating hyperplane that maximises the minimum distance between closest points (support vectors) of the two metastable states (otherwise known as classes):

$$d_H = \frac{\mathbf{w}^T \mathbf{r} + b}{\|\mathbf{w}\|_2} \quad (32)$$

where d_H is the distance of a point to hyperplane, \mathbf{w} is the vector of coefficients (it's direction gives us the predicted class: positive or negative), b is the intercept and $\frac{1}{\|\mathbf{w}\|_2}$ is the normalisation term. The SVM algorithm then maximise the minimum distance $\mathbf{w}^* = \arg_{\mathbf{w}} \max[\min_{\mathbf{r}} d_H]$. In context of biomolecular simulation, eqn (32) acts a CV which can separate metastable states and drive enhanced sampling simulations (Fig. 14). Sultan and Pande showed how d_H can be used as CV within metadynamics framework to sample conformational space of alanine dipeptide and chignolin folding.⁶¹ SVM algorithm can be further extended for non-linear classification by using kernel trick as explained before.

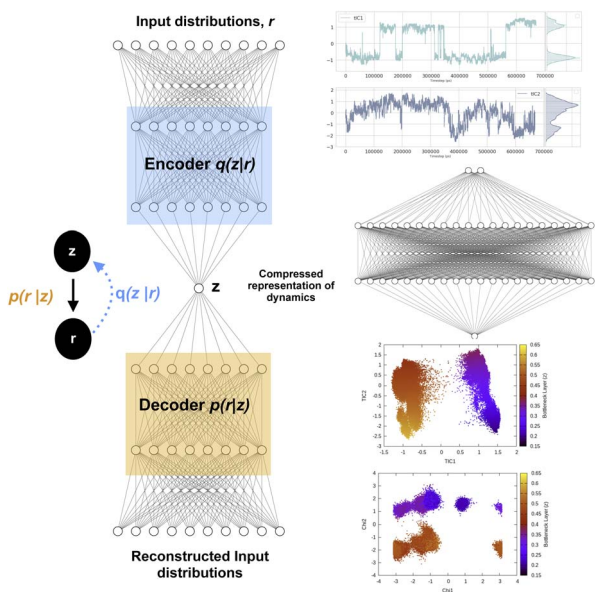


Fig. 12 Left panel showing a typical architecture of deep neural network based variational autoencoder. The input distributions r can be probability distributions of geometric variables, TICs, DMD *etc.* Right panel shows how latent layer of variational autoencoder learned the dynamics of the system using TICs as inputs. TICs were generated using χ_1 and χ_2 angles of residues present in the flap region of plasmeptsin-II (Table 1 in ESI†) with lag-time 1000. Following hyperparameters were used during the training process of VAE, no of hidden layers = 2, no of neurons in each hidden layer = 20, epochs = 100, batch size = 500, learning rate = $1e-2$.



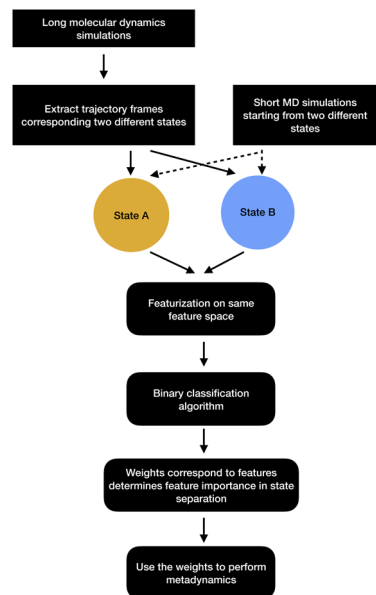


Fig. 13 Schematic representation of showing how binary classifiers can be used to separate metastable states and drive enhanced sampling calculations.

2.2.7.2 Linear discriminant analysis. Linear discriminant analysis (LDA) is a supervised dimensionality reduction method which finds linear combinations of features that best separates two or more metastable states.⁶⁴ Imagine two metastable states with common features \vec{r} have means $\vec{\mu}_A$ and $\vec{\mu}_B$ and covariances Σ_A and Σ_B . The linear combinations of features $\vec{w} \cdot \vec{r}$ (\vec{w} is known as LDA coefficients) will have means $\vec{w} \cdot \vec{\mu}_{A,B}$ and variances $\vec{w}^T \Sigma_{A,B} \vec{w}$. LDA objective function separates two distributions by taking ratio between variance between the classes to variance within classes:

$$s(\text{LDA}) = \frac{(\vec{w} \cdot (\vec{\mu}_B - \vec{\mu}_A))^2}{\vec{w}^T (\Sigma_A + \Sigma_B) \vec{w}} \quad (33)$$

Mathematically it can be shown that the maximum separation occurs when $\vec{w} \propto \frac{(\vec{\mu}_B - \vec{\mu}_A)}{(\Sigma_A + \Sigma_B)}$. LDA algorithm can be generalised to more than one classes (Multi-class LDA). Eqn (33) acts as a CV in metadynamics to drive transition between state A and B. Recently Parrinello and co-workers used a harmonic version of LDA (HLDA) as a CV to study folding of a small protein and protein-ligand binding/unbinding.^{62,65} The framework of LDA can be patched with neural network (Deep-LDA) as well as kernel trick in order to introduce non-linearity. Recently Deep-LDA framework has been applied in SAMPL5 host-guest systems to accurately calculate binding free energy and to understand the role of water molecule in ligand binding.⁶⁶

In order to apply classifier algorithms to separate metastable states one first need to sample different metastable states. However, sampling of metastable states separated by high entropic barrier often requires ultra-long MD simulations or enhanced sampling simulations *e.g.* parallel-tempering.

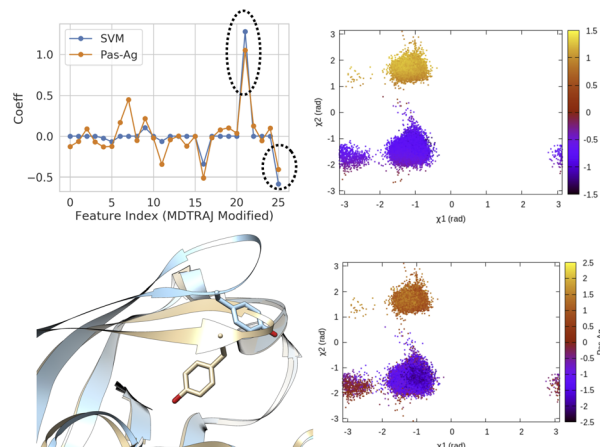


Fig. 14 Coefficients (weights) correspond to linear SVM and passive aggressive classifier⁶³ trained on χ_1 and χ_2 angles of residue 74–84 (Table 1 in ESI†) in plasmepsin-II showing the importance of $\sin \chi_2$ (feature 21) and $\cos \chi_2$ (feature 25) in separating two metastable basins, normal ($\chi_2 \sim -1$ rad) and flipped ($\chi_2 \sim 2$ rad). 50 ns trajectories from each basins were used for training purpose using the scheme described in Fig. 13. Further χ_1 and χ_2 angles of Tyr77 projected along classifier decision boundaries demonstrated separation of metastable states. Snapshots corresponding normal (grey) and flipped (blue) states are also highlighted.

Further, one needs to carefully choose a set of features that can separate a set of metastable states. For complex systems, selection of such features are not trivial and often needs significant amount of trial and error.

Softwares: scikit-learn's supervised learning module incorporates both LDA and SVM algorithms for training purpose. It further integrates various other supervised learning algorithms that can be applied to classify metastable states and act as CVs in enhanced sampling simulations.

2.2.8 Algorithms to predict temporal evolution of CVs: proposing a challenging dataset. Recently neural network based time-series prediction models such as long short-term memory (LSTM) and transformers were proposed to predict rare events and extract kinetics and thermodynamics.^{67,68} These methods often work quite well on simpler systems where the training data has multiple recrossing. However in order to correctly predict kinetics from time-series one first needs to accurately predict the frequency and lifetime of rare events. In case of Fig. 15, temporal evolution along H-bond and dihedral angle CVs showed multiple recrossing whereas evolution along TIC1 only captures transient recrossing between metastable states. Intuitively one can say predicting temporal evolution along TIC1 (rare fluctuation) is far more challenging task when compared with H-bond or dihedral angle. No such study is currently available which compares the time-series prediction capability of neural networks in context of different CVs. In future, such study is necessary to understand the limit of neural network based time-series prediction algorithms⁶⁹ in terms of capturing temporal evolution of rare events.

Softwares: Tiwary and co-workers⁶⁷ used LSTM based neural network to predict temporal evolution in model systems. Their



code can be accessed here: <https://github.com/tiwarylab/LSTM-predict-MD>. Tensorflow (<https://www.tensorflow.org/text/tutorials/transformer> and <https://www.tensorflow.org/guide/keras/rnn>) and Pytorch (<https://github.com/jdb78/pytorch-forecasting>) based time-series forecasting modules can be used for this purpose. Zeng and co-workers⁶⁸ used LSTM/transformers to predict temporal evolution of different CVs in alanine dipeptide (<https://github.com/Wendysigh/LSTM-Transformer-for-MD>).

3 Use of machine learning/AI in CV selection: the hype

Recently numerous papers reported different combinations of abstract order parameters to capture conformational dynamics and molecular recognition of macromolecules. In one such case neural network based LDA (Deep-LDA) approach has been applied on SAMPL5 host–guest systems to estimate binding free energy.⁶⁶ Previously funnel metadynamics using simple

distance CV combined with artificial restraint along binding-site solvation managed to accurately predict the binding free energy for the same host–guest systems.⁷⁰ We can ask a question: was deep-LDA framework necessary for a relatively simple problem such as host–guest binding/unbinding which was previously solved using simple geometric CVs? One can argue that deep-LDA framework is more natural as it didn't impose any artificial restraint along solvent degrees of freedom. In SAMPL5 host–guest systems, the solvation (when the binding site is occupied by long-lived water molecule) of the binding site is a slow process (Fig. S2 in ref. 70). In principal TICA/SGOOP (which by design capable of capturing slow degrees of freedom) could have solved such a problem. Further, deep-LDA method is a slightly fancier application of previously described work which uses binary classifiers as CVs to drive conformational sampling.⁶¹ This is a classic example where the developers of a CV discovery method used buzzwords to solve a trivial problem without pointing out how exactly such deep learning based CV is superior compared to other geometric or abstract order parameters.

In another case authors used TICs/SGOOPs as inputs within a VAE and used the VAE's latent layer as CV for enhanced sampling.^{55,71} Such an approach is advantageous if an enhanced sampling algorithm is limited to driving along one CV. Further, the eigenvalues corresponding each TIC/SGOOP can be used to filter geometric CVs which then can be used as an input in VAE. This approach can be seen as an extension of previous approach where one inputs TICs/SGOOPs directly into VAE. However a critic can ask: is there any advantage in terms of sampling if one uses VAE's latent layer vs. multiple CVs (using parallel-bias metadynamics) within metadynamics (Fig. 16) and more importantly can you call such method artificial intelligence (AI)?

Most of the algorithms described in the previous section either require extensive sampling or prior knowledge about metastable states in order to capture conformational changes in biomolecules. These algorithms can be best described as data driven machine learning algorithms which solves a deterministic closed set finite problems using large amount of training data.

In case of complex biological process *e.g.* protein–ligand unbinding, protein–protein binding, protein folding involving multiple metastable states and lack of prior knowledge about transition states, an iterative machine learning based CV discovery method combined with enhanced sampling can accelerate sampling of the underlying free energy surface.^{73,74} However such protocol also requires initial guess of CVs. An near-term artificial general intelligence (AGI) on other hand will learn about the dynamical system on the fly without the external supervision (*e.g.* feature selection, choice of algorithms, tuning of hyper-parameters *etc*) and extract observables that can be confirmed by biophysical experiments. Development of such self-aware AGI will require significant innovation in terms of novel algorithms and software design. Until that time, the biomolecular simulation community should refrain from using the word AI in context of identifying low-dimensional representation to approximate kinetics and thermodynamics of biomolecular systems.

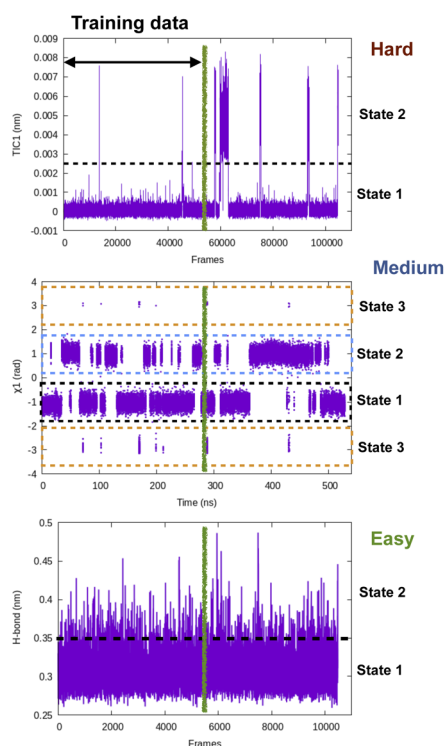


Fig. 15 Pictorial representation showing time-series fluctuation of different CVs and their hypothetical prediction difficulty using LSTM, transformers style neural network. Along H-bond distance one can see multiple recrossing between two states which makes it an easy case for time-series prediction algorithms. In case of dihedral flipping along χ_1 , training data observes several recrossing between state 1 and 2 however, there are transient sampling of rare transitions to state 3 which makes it a medium level difficult task. In case of TIC1, the training data contains rare transitions between state 1 and 2 hence I believe it will be a hard task for time-series prediction algorithms. One way to compare time-series projection results along a CV is compare mean first passage transition time (MFPT) as well as lifetime of different states with temporal data from long MD simulation.



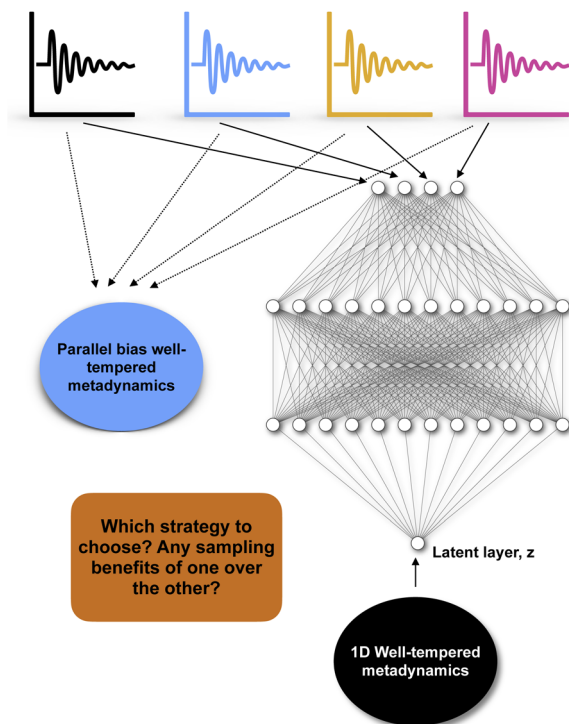


Fig. 16 Describes dilemma of choosing an enhanced sampling strategy over other: in case of four different CVs (depicted by black, blue, yellow and magenta) one can either choose to perform a parallel-bias⁷² variant of well-tempered metadynamics simulation or feed these four CVs within variational autoencoder and uses the latent layer of variational autoencoder as CV in 1D well-tempered metadynamics.

4 Conclusions

In recent years data-driven machine learning models impacted the field of biomolecular simulation and have been applied in context of protein folding, protein-ligand/protein binding and capturing rare conformational changes during molecular simulation. However a real challenge in comparing the power of different machine learning models is due to absence of a common dataset. We believe that the basic pancreatic trypsin inhibitor (BPTI) can act as a reasonable dataset to test different algorithms. The reasons behind that are accessibility of (a) 1 ms long MD simulation and (b) experimental observables such as NMR order parameter, kinetics of ring flipping (rare events) and protection factors (measure of amide hydrogen-deuterium exchange). Further transient loop opening of BPTI which leads to amide hydrogen-deuterium exchange²⁹ and aromatic ring flipping^{75,76} are the rare events that can be captured by both simulation and experiments. Transient loop opening leads to amide hydrogen-deuterium exchange which is represented by a metric called protection factor. A good machine learning model for CV discovery combined with an enhanced sampling framework should sample the rare events and predict experimental observables (rate of ring flipping and protection factors). Recent metadynamics based investigation with geometric CVs (dihedral angles) highlighted the challenge of sampling slow

ring flipping in BPTI. In order to live up to the hype, machine learning based abstract CVs should do better in sampling such rare events and predict kinetics.

A major leap forward in combining machine learning with MD will be prediction of multi-dimensional spatiotemporal evolution of biomolecular dynamics. Recently, transformer network has been applied to model crowd motion dynamics which managed to achieve state-of-the-art performance on commonly used pedestrian prediction datasets.⁷⁷ In future extension of such work can be applied in context of MD simulation to predict temporal evolution of CVs and corresponding spatial representation. However prediction of spatiotemporal evolution of biomolecular structure without correctly capturing the water dynamics doesn't model physical reality as water plays integral role in molecular recognition and conformational dynamics of biomolecules. We hope breakthroughs in quantum computer together with novel machine learning algorithms will make spatiotemporal prediction of solvated biomolecular system a reality, until then we have to carefully select geometric/abstract CVs to capture dynamics of biomolecules from atomistic simulations.

Author contributions

SB conceived the idea, performed all the calculations and wrote the manuscript.

Conflicts of interest

There are no conflicts to declare.

Acknowledgements

The computations were performed on computer resources provided by the Swedish National Infrastructure for Computing (SNIC) at LUNARC (Lund University) and HPC2N (Umeå University). The author thanks Born Novak for carefully reading the manuscript.

Notes and references

- 1 S. A. Hollingsworth and R. O. Dror, *Neuron*, 2018, **99**, 1129–1143.
- 2 J. Fan and J. Lv, *Stat. Sin.*, 2010, **20**, 101–148.
- 3 S. Bhakat, *RSC Adv.*, 2021, **11**, 11026–11047.
- 4 S. Bhakat and P. Söderhjelm, *J. Chem. Inf. Model.*, 2022, **62**, 914–926.
- 5 G. Bussi and A. Laio, *Nat. Rev. Phys.*, 2020, **2**, 200–212.
- 6 Y. Wang, J. M. Lamim Ribeiro and P. Tiwary, *Curr. Opin. Struct. Biol.*, 2020, **61**, 139–145.
- 7 M. Chen, *Eur. Phys. J. B*, 2021, **94**, 211.
- 8 O. Valsson, P. Tiwary and M. Parrinello, *Annu. Rev. Phys. Chem.*, 2016, **67**, 159–184.
- 9 J. S. Patel, A. Berteotti, S. Ronsisvalle, W. Rocchia and A. Cavalli, *J. Chem. Inf. Model.*, 2014, **54**, 470–480.
- 10 W. You, Z. Tang and C.-e. A. Chang, *J. Chem. Theory Comput.*, 2019, **15**, 2433–2443.



- 11 L. S. Dodda, J. Tirado-Rives and W. L. Jorgensen, *J. Phys. Chem. B*, 2019, **123**, 1741–1748.
- 12 F. Sittel and G. Stock, *J. Chem. Phys.*, 2018, **149**, 150901.
- 13 F. Persson and B. Halle, *J. Am. Chem. Soc.*, 2013, **135**, 8735–8748.
- 14 D. E. Shaw, P. Maragakis, K. Lindorff-Larsen, S. Piana, R. O. Dror, M. P. Eastwood, J. A. Bank, J. M. Jumper, J. K. Salmon, Y. Shan and W. Wriggers, *Science*, 2010, **330**, 341–346.
- 15 G. A. Tribello, M. Bonomi, D. Branduardi, C. Camilloni and G. Bussi, *Comput. Phys. Commun.*, 2014, **185**, 604–613.
- 16 D. R. Roe and T. E. Cheatham, *J. Chem. Theory Comput.*, 2013, **9**, 3084–3095.
- 17 N. Michaud-Agrawal, E. J. Denning, T. B. Woolf and O. Beckstein, *J. Comput. Chem.*, 2011, **32**, 2319–2327.
- 18 R. T. McGibbon, K. A. Beauchamp, M. P. Harrigan, C. Klein, J. M. Swails, C. X. Hernández, C. R. Schwantes, L.-P. Wang, T. J. Lane and V. S. Pande, *Biophys. J.*, 2015, **109**, 1528–1532.
- 19 F. Sittel, T. Filk and G. Stock, *J. Chem. Phys.*, 2017, **147**, 244101.
- 20 M. P. Harrigan, M. M. Sultan, C. X. Hernández, B. E. Husic, P. Eastman, C. R. Schwantes, K. A. Beauchamp, R. T. McGibbon and V. S. Pande, *Biophys. J.*, 2017, **112**, 10–15.
- 21 M. K. Scherer, B. Trendelkamp-Schroer, F. Paul, G. Pérez-Hernández, M. Hoffmann, N. Plattner, C. Wehmeyer, J.-H. Prinz and F. Noé, *J. Chem. Theory Comput.*, 2015, **11**, 5525–5542.
- 22 L. Molgedey and H. G. Schuster, *Phys. Rev. Lett.*, 1994, **72**, 3634–3637.
- 23 G. Pérez-Hernández, F. Paul, T. Giorgino, G. De Fabritiis and F. Noé, *J. Chem. Phys.*, 2013, **139**, 015102.
- 24 C. R. Schwantes and V. S. Pande, *J. Chem. Theory Comput.*, 2013, **9**, 2000–2009.
- 25 M. M. Sultan and V. S. Pande, *J. Chem. Theory Comput.*, 2017, **13**, 2440–2447.
- 26 S. Schultze and H. Grubmüller, *J. Chem. Theory Comput.*, 2021, **17**, 5766–5776.
- 27 T. Blaschke, P. Berkes and L. Wiskott, *Neural Comput.*, 2006, **18**, 2495–2508.
- 28 T. Blaschke and L. Wiskott, *Independent Component Analysis and Blind Signal Separation*, Berlin, Heidelberg, 2004, pp. 742–749.
- 29 F. Persson and B. Halle, *Proc. Natl. Acad. Sci. U. S. A.*, 2015, **112**, 10383–10388.
- 30 M. Hoffmann, M. Scherer, T. Hempel, A. Mardt, B. de Silva, B. E. Husic, S. Klus, H. Wu, N. Kutz, S. L. Brunton and F. Noé, *Mach. learn.: sci. technol.*, 2021, **3**, 015009.
- 31 C. R. Schwantes and V. S. Pande, *J. Chem. Theory Comput.*, 2015, **11**, 600–608.
- 32 M. P. Harrigan and V. S. Pande, *bioRxiv*, 2017, DOI: [10.1101/123752](https://doi.org/10.1101/123752).
- 33 C. Ross, B. Nizami, M. Glenister, O. Sheik Amamuddy, A. R. Atilgan, C. Atilgan and Ö. T. Bishop, *Bioinformatics*, 2018, **34**, 3759–3763.
- 34 F. Noé and C. Clementi, *J. Chem. Theory Comput.*, 2015, **11**, 5002–5011.
- 35 B. Nadler, S. Lafon, I. Kevrekidis and R. Coifman, *Advances in Neural Information Processing Systems*, 2006.
- 36 W. Zheng, M. A. Rohrdanz and C. Clementi, *J. Phys. Chem. B*, 2013, **117**, 12769–12776.
- 37 A. L. Ferguson, A. Z. Panagiotopoulos, P. G. Debenedetti and I. G. Kevrekidis, *Proc. Natl. Acad. Sci. U. S. A.*, 2010, **107**, 13597–13602.
- 38 A. L. Ferguson, A. Z. Panagiotopoulos, I. G. Kevrekidis and P. G. Debenedetti, *Chem. Phys. Lett.*, 2011, **509**, 1–11.
- 39 S.-T. Tsai, Z. Smith and P. Tiwary, *J. Chem. Theory Comput.*, 2021, **17**, 6757–6765.
- 40 L. van der Maaten and G. Hinton, *J. Mach. Learn. Res.*, 2008, **9**, 2579–2605.
- 41 M. Wattenberg, F. Viégas and I. Johnson, *Distill*, 2016, DOI: [10.23915/distill.00002](https://doi.org/10.23915/distill.00002).
- 42 V. Spiwok and P. Kříž, *Front. Mol. Biosci.*, 2020, **7**, 132.
- 43 T. Chari, J. Banerjee and L. Pachter, *bioRxiv*, 2021, DOI: [10.1101/2021.08.25.457696](https://doi.org/10.1101/2021.08.25.457696).
- 44 P. Tiwary and B. J. Berne, *Proc. Natl. Acad. Sci. U. S. A.*, 2016, **113**, 2839–2844.
- 45 B. Ghogjogh, A. Ghodsi, F. Karray and M. Crowley, *Multidimensional Scaling, Sammon Mapping, and Isomap: Tutorial and Survey*, 2020.
- 46 J. H. Tu, C. W. Rowley, D. M. Luchtenburg, S. L. Brunton and J. N. Kutz, *J. Comput. Dyn.*, 2014, **1**, 391–421.
- 47 S. Klus, F. Nüske, P. Koltai, H. Wu, I. Kevrekidis, C. Schütte and F. Noé, *J. Nonlinear Sci.*, 2018, **28**, 985–1010.
- 48 L. McInnes, J. Healy and J. Melville, *UMAP: Uniform Manifold Approximation and Projection for Dimension Reduction*, 2020.
- 49 F. Trozzi, X. Wang and P. Tao, *J. Phys. Chem. B*, 2021, **125**, 5022–5034.
- 50 H. Tian and P. Tao, *J. Chem. Inf. Model.*, 2020, **60**, 4569–4581.
- 51 S. L. Brunton, J. L. Proctor and J. N. Kutz, *Proc. Natl. Acad. Sci. U. S. A.*, 2016, **113**, 3932–3937.
- 52 A. Glielmo, B. E. Husic, A. Rodriguez, C. Clementi, F. Noé and A. Laio, *Chem. Rev.*, 2021, **121**, 9722–9758.
- 53 D. P. Kingma and M. Welling, *Found. Trends Mach. Learn.*, 2019, **12**, 307–392.
- 54 D. P. Kingma and M. Welling, *Auto-Encoding Variational Bayes*, 2014.
- 55 M. M. Sultan, H. K. Wayment-Steele and V. S. Pande, *J. Chem. Theory Comput.*, 2018, **14**, 1887–1894.
- 56 C. P. Burgess, I. Higgins, A. Pal, L. Matthey, N. Watters, G. Desjardins and A. Lerchner, *Understanding disentangling in β -VAE*, 2018.
- 57 H. Chen, H. Liu, H. Feng, H. Fu, W. Cai, X. Shao and C. Chipot, *J. Chem. Inf. Model.*, 2022, **62**, 1–8.
- 58 C. X. Hernández, H. K. Wayment-Steele, M. M. Sultan, B. E. Husic and V. S. Pande, *Phys. Rev. E*, 2018, **97**, 062412.
- 59 J. M. L. Ribeiro, P. Bravo, Y. Wang and P. Tiwary, *J. Chem. Phys.*, 2018, **149**, 072301.
- 60 N. Mohanty, A. L.-S. John, R. Manmatha and T. Rath, *Handbook of Statistics*, Elsevier, 2013, vol 31, pp. 249–267.
- 61 M. M. Sultan and V. S. Pande, *J. Chem. Phys.*, 2018, **149**, 094106.
- 62 D. Mendels, G. Piccini, Z. F. Brotzakis, Y. I. Yang and M. Parrinello, *J. Chem. Phys.*, 2018, **149**, 194113.



- 63 K. Crammer, O. Dekel, J. Keshet, S. Shalev-Shwartz and Y. Singer, *J. Mach. Learn. Res.*, 2006, **7**, 551–585.
- 64 R. A. Fisher, *Ann. Hum. Genet.*, 1936, **7**, 179–188.
- 65 R. Capelli, A. Bochicchio, G. Piccini, R. Casasnovas, P. Carloni and M. Parrinello, *J. Chem. Theory Comput.*, 2019, **15**, 3354–3361.
- 66 V. Rizzi, L. Bonati, N. Ansari and M. Parrinello, *Nat. Commun.*, 2021, **12**, 93.
- 67 S.-T. Tsai, E.-J. Kuo and P. Tiwary, *Nat. Commun.*, 2020, **11**, 5115.
- 68 W. Zeng, S. Cao, X. Huang and Y. Yao, *A Note on Learning Rare Events in Molecular Dynamics using LSTM and Transformer*, 2021.
- 69 B. Lim and S. Zohren, *Philos. Trans. R. Soc., A*, 2021, **379**, 20200209.
- 70 S. Bhakat and P. Söderhjelm, *J. Comput.-Aided Mol. Des.*, 2017, **31**, 119–132.
- 71 S. Pant, Z. Smith, Y. Wang, E. Tajkhorshid and P. Tiwary, *J. Chem. Phys.*, 2020, **153**, 234118.
- 72 J. Pfendtner and M. Bonomi, *J. Chem. Theory Comput.*, 2015, **11**, 5062–5067.
- 73 S. Pant, Z. Smith, Y. Wang, E. Tajkhorshid and P. Tiwary, *J. Chem. Phys.*, 2020, **153**, 234118.
- 74 J. M. Lamim Ribeiro and P. Tiwary, *J. Chem. Theory Comput.*, 2019, **15**, 708–719.
- 75 U. Weininger, K. Modig and M. Akke, *Biochemistry*, 2014, **53**, 4519–4525.
- 76 M. Kulkarni and P. Söderhjelm, *bioRxiv*, 2021, DOI: [10.1101/2021.01.07.425261](https://doi.org/10.1101/2021.01.07.425261).
- 77 C. Yu, X. Ma, J. Ren, H. Zhao and S. Yi, *Spatio-Temporal Graph Transformer Networks for Pedestrian Trajectory Prediction*, 2020.

