# Machine learning force fields and coarse-grained variables in molecular dynamics: application to materials and biological systems

**Paraskevi Gkeka**[1,✉], **Gabriel Stoltz**[2,3,✉], **Amir Barati Farimani**[4], **Zineb Belkacemi**[1], **Michele Ceriotti**[5], **John Chodera**[6], **Aaron R. Dinner**[7], **Andrew Ferguson**[8], **Jean-Bernard Maillet**[9], **Hervé Minoux**[10], **Christine Peter**[11], **Fabio Pietrucci**[12], **Ana Silveira**[6], **Alexandre Tkatchenko**[13], **Zofia Trstanova**[14], **Rafal Wiewiora**[6], and **Tony Lelièvre**[2,3,✉]

[1]Structure Design and Informatics, Sanofi R&D, 91385 Chilly-Mazarin, France
[2]Ecole des Ponts ParisTech, France
[3]Matherials project-team, Inria Paris, France
[4]Carnegie Mellon University, USA
[5]Laboratory of Computational Science and Modelling, Institute of Materials, École Polytechnique Fédérale de Lausanne, Switzerland
[6]Sloan Kettering Institute, USA
[7]Department of Chemistry, The University of Chicago, Chicago, Illinois 60637, USA
[8]Pritzker School of Molecular Engineering, 5640 South Ellis Avenue, University of Chicago, Chicago, Illinois 60637, USA
[9]CEA-DAM, DIF, France
[10]Structure Design and Informatics, Sanofi R&D, 94403 Vitry-sur-Seine, France
[11]University of Konstanz, Germany
[12]Sorbonne Université, UMR CNRS 7590, MNHN, Institut de Minéralogie, de Physique des Matériaux et de Cosmochimie, 75005 Paris, France
[13]Physics and Materials Science Research Unit, University of Luxembourg, Luxembourg
[14]School of Mathematics, The University of Edinburgh, UK

**Machine learning encompasses a set of tools and algorithms which are now becoming popular in almost all scientific and technological fields. This is true for molecular dynamics as well, where machine learning offers promises of extracting valuable information from the enormous amounts of data generated by simulation of complex systems. We provide here a review of our current understanding of goals, benefits, and limitations of machine learning techniques for computational studies on atomistic systems, focusing on the construction of empirical force fields from ab-initio databases and the determination of reaction coordinates for free energy computation and enhanced sampling.**

**Machine learning | Molecular Dynamics | Coarse-graining | Chemical Physics | Force fields | Reaction Coordinates | Collective Variables | Enhanced sampling**

**Correspondence:** *Paraskevi.Gkeka@sanofi.com, gabriel.stoltz@enpc.fr, tony.lelievre@enpc.fr*

## 1. Introduction

The atomistic representation of physical systems offers a precise description of matter. Simplified models based on coarse-grained (CG) representations offer an alternative that can significantly aid in the understanding of the physical properties of the systems under consideration. Such representations can also be used as a surrogate model for enhanced sampling methods (e.g. sampling large conformational changes using reduced models).

Both in the case of biochemical systems as well as in materials, a CG description can be based on distance metrics for structural clustering (1), as well as on reaction coordinates: for instance, the conformational changes of a complex molecule can be modeled by a few key functions of the atomic positions, while a phase transition can be described by a change of the average atomic coordination or box shape. In condensed matter physics, atomic descriptors are employed to summarize the key features of atomic configurations in order to predict forces and energies (2, 3).

In the past, reaction coordinates were defined using empirical methods and chemical intuition, while more systematic approaches were employed for the definition of atomic descriptors (4, 5). During the last decade, the return and rise of Machine Learning (ML) techniques have initiated many efforts focusing on automating the definition of reaction coordinates or descriptors that are able to successfully describe the underlying atomic systems (6–9). The employed methods, both supervised and unsupervised, vary. The most commonly used methods for the identification of reaction coordinates include Principal Component Analysis (PCA) (10), diffusion maps (11, 12), and auto-encoders (13–16). For atomic descriptors, common choices are based on a judicious use of adjacency matrices and their generalizations, or on a large set of feature vectors based on a set of basis functions.

We are witnessing many current attempts for automatically devising intuition-free collective variables, in particular for drug discovery applications (13, 17). Although the initially very high hopes raised by numerical potentials are now mitigated, there have been quite a few systematic studies on the quality of the descriptors obtained by these approaches (18, 19).

A recent CECAM (Center Européen de Calcul Atomique et Moléculaire) discussion meeting[1] brought together a di-

---

[1]See the conference website:
`https://cermics-lab.enpc.fr/cecam_ml_md/`

verse audience of 29 participants from various scientific fields, including chemistry, drug design, condensed matter physics, materials science, and mathematics, to exchange about state-of-the-art techniques for automatically building coarse-grained information on molecular systems. In particular, we believe that the viewpoint and experience of condensed matter physicists in devising atomic descriptors could prove useful insights in devising reaction coordinates in a more systematic way. Mathematics offer, in this framework, a common language for the discussion. One distinctive feature of this CECAM meeting is that the emphasis was on the technical details of the underlying numerical methods.

In the current review, we discuss the following highlights of the meeting:

- **Machine learning force fields and Potential of Mean Force.** ML techniques have been recently employed in the development of force field (FF) parameters based on quantum-mechanical calculations. More generally, ML techniques can be used to define a surrogate model of any quantity that could be obtained from a quantum chemical calculation, as a function of atomic coordinates (e.g. NMR chemical shieldings, IR dipole moments, ...), making it possible to obtain an accurate estimate of experimental observables. Such models are beginning to find merit due to their accuracy and versatility. In Section 2, we review the factors that play an important role in the accuracy and transferability of a force field. Specifically, we report the importance of the input database and the choice of the regression method for the force field construction. The use of prior physico-chemical knowledge in this construction of ML potentials is also discussed.

- **Dimensionality reduction and identification of meaningful collective variables.** Another important issue discussed during the CECAM meeting is the dimensionality reduction and the identification of meaningful CVs using ML techniques (see Section 3). We considered the case when this identification relies on a database which covers the full configuration space of the system under study (obtained for instance by high temperature sampling, steered molecular dynamics, etc), and the case when the data is restricted to a metastable state. Once a reaction coordinate is found, the question of devising a good effective model along this coordinate can also be addressed using machine learning techniques: either approximate free energies (for example by potentials involving only 2, 3 or 4 body interactions), or approximate the terms in the effective dynamics, namely the drift, diffusion coefficient, metric tensor and memory terms, for example using projections *à la* Mori-Zwanzig.

- **Applications of machine learning techniques in biological systems and drug discovery.** In Section 4, we discuss some "real world" applications, where MD simulations coupled with ML techniques enable us to understand the biological complexity at the atomic and molecular levels and provide us with interesting insights about the thermodynamic and mechanistic behaviour of biological processes. In particular, we highlight some examples of ML approaches applied in clustering and construction of Markov state models, we describe how ML methods facilitate enhanced sampling protocols through the use of efficient CVs and we mention some possible applications in the drug discovery process. These examples illustrate the current state and potential of the field of ML in the study of biological systems and drug discovery.

We close the review with some perspectives in Section 5.

## 2. Machine learning force fields and Potential of Mean Force

Interactions between atoms are often modeled using empirical potentials with some prescribed functional forms, as suggested by physical considerations. This provides computationally cheap (with a cost scaling linearly with the number of atoms) but somewhat inaccurate potentials. On the contrary, ab-initio approaches provide more reliable, less uncertain force fields, at the expense however of a large computational cost (typically scaling as the number of electrons to the power 3). The promise of machine learning for force field computations is to predict forces and energies with accuracy arbitrary close to the level of ab-initio approaches (20), but with a much smaller computational cost and scaling as a function of the number of atoms. Ideally, these force fields should be able to describe chemical reactions. This is typically done in practice by setting up a database of configurations with associated forces and energies, summarizing atomic configurations through some descriptors of the local environment, and predicting the forces and energies from these descriptors through a function which has been trained by some (nonlinear) regression procedure to provide good results on the database. The resulting potential is called a "numerical potential".

There are three different factors to discuss the success of ML methods, whose relative importance depend on the aims of the user: accuracy, computational cost, and transferability. The latter concept means that a numerical potential computed for a given material in a given thermodynamic range, can be used outside the fitting domain – for instance because it is used for other materials and systems than the ones it was trained on, and/or in a different thermodynamic range than the one considered for the configurations in the database.

We first discuss in this section elements on the choice of the database, see Section A. We next present various choices for the descriptors and for associated ML regression methods, see Section B. We then discuss in Section C how to incorporate physical insights in order to improve ML techniques, and we give some perspectives in Section D. We end the section by mentioning how ML approaches can also be used to derive CG potentials, see Section E: in this perspective, empirical force fields for all atom models are seen as the reference (they are the counterpart of ab-initio databases in this

context), and effective force fields describing the interaction of coarse-grained variables are sought.

## A. Setting up a database.
One of the key factors that affects the accuracy and transferability of a force field is the database used for its construction. This database defines the envelope of confidence (applicability domain) for the potential as the subsequent regression method is efficient in interpolation. It is often the case that a numerical potential has a poor transferability. Therefore, for condensed matter systems, the database should sample the region of interest, i.e., the thermodynamic conditions where the potential is going to be used. However, this representative part of the configurational space covers only a small fraction of the overall available space. Hence, a systematic exploration is impossible, and physical intuition is often used to constrain the search of new interesting configurations for learning. This makes the construction of the database a rather laborious process. A first application of 'active learning' in this process, also still hand made, is proposed by Artrith and Behler in Ref. 21: two different neural networks are optimized on the same database and, in case their predictions on a new configuration differ too much this configuration should be included in the database. Active learning, based on outlier detection (i.e., definition of a metric to detect parameters corresponding to some extrapolation) is now routinely employed during the database construction (22). In this way, force field accuracy can be improved during the training procedure (23) and the domain of applicability could be extended (24). The bottom line is that 'on the fly' learning (25) enables to perform optimization and prediction at the same time (26). Typically, a trade-off has to be found between the transferability of a potential (its robustness to changes in the database) and its accuracy.

The representation of the database should also be meaningful: finding a proper space for this representation allows to define an envelope of confidence for the potential. When the potential is used, each new configuration can rapidly be plotted in this space to check if it belongs to the database envelope (applicability domain), i.e., if the potential is used in interpolation or in extrapolation. It then becomes a useful criterion for outlier detection.

What is globally accepted is that the methods should systematically be validated on test data, different from the training data. In any case, one should be very careful about the quality of the model for extrapolation.

## B. Descriptors and regression methods.
We present in this section the technical approaches to fit a potential on a database. We distinguish the representation of the atomic configurations through descriptors, and the subsequent regression allowing to fit the parameters of the chosen model. Typically, a very simple descriptor, based on physical/chemical intuition or moment estimates for atomic densities, should be combined with a complex regression such as a neural network; on the other hand, more educated descriptors, for instance based on convolutional neural networks

and a scattering transform (27), can be fed into quite simple (bi)linear regression models.

***B.1. Representing atomic configurations.*** It is almost never appropriate to use the Cartesian coordinates of atoms in a structure as the input of a machine-learning scheme (28), because Cartesian coordinates do not conform with the invariance of the target properties, e.g. permutation of the indices of identical atoms, rigid translations, rotations and reflections. For this reason, several different schemes have been devised to map atomic configurations onto vectors of features that fulfil these symmetry requirements. Usually, it is desirable for this mapping to be differentiable and smooth, particularly in applications where one needs to compute forces as the derivative of a machine-learning potential or CG force field.

One can roughly partition methods to represent atomic configurations into two classes. *Descriptors* are often highly simplified representations of a structure, usually of much smaller dimensionality than the number of degrees of freedom and incorporating some degree of chemical intuition, or a heuristic understanding of the behavior of the system being studied. Cheminformatics schemes to characterise the connectivity of a molecule, such as SMILES (29) strings, are useful when dealing with databases of organic compounds. Steinhardt parameters (30) are often used to characterize the coordination of liquids and solids. Backbone dihedral angles, or more complex indicators of secondary structure (31) can be utilized to discard information on the side chains of polypeptides. The dimensionality reduction that is intrinsic to this family of methods typically induce loss of information, which may be desirable (when it discards irrelevant details) or problematic: in the latter case, it is often more effective to use a more complete description and then proceed with an automatic dimensionality reduction algorithm, some of which will be discussed in Section 3.

*Representations*, on the other hand, attempt to provide a complete description of a configuration. This family of features is typically used when building regression models for energy and properties. Most of the time (particularly for condensed-phase applications, but often also for isolated molecules) representations are not built for an entire structure, but are instead used to describe atom-centered environments. This is advantageous, because - by representing a structure as a collection of compact groups of atoms, and assuming that the overall property can be computed as a sum of local contributions - it becomes possible to train models that can be easily transferred between systems of different sizes, and from simple to more complex configurations. Many of these systematic representations - including e.g., SOAP (bi)spectrum (32), Behler-Parrinello symmetry functions (33), moment tensor potentials (18), FCHL kernels (34) - can be seen as projections on different basis of n-body correlation functions (35), and offer a systematic and completely general way to describe atomic configurations, that can be applied equally well to condensed phases, gas-phase molecules and polypeptides (36).

***B.2. Choosing the regression method.*** Once the atomic descriptor has been chosen, the choice of the regression method to determine the force field is crucial and greatly depends on the system under study (37). A distinction should be made between learning based on neural networks, and other regression methods based on kernels or (bi)linear methods. Training neural networks is a complex non-convex optimization problem in very high dimension (generally thousands of parameters are needed to parameterize the networks under consideration). Already the computation of the gradient of the objective function is non trivial and relies on clever numerical tricks, such as backpropagation. Kernel-based methods or (bi)linear regression techniques lead, on the other hand, to much better behaved optimization problems, which can even be solved analytically through some matrix inversion on the Euler equation defining the minimizer.

The choice of the regression method also determines whether error estimators are available. For example a variance can be associated with a prediction when a kernel method is used, whereas error quantification is harder using neural networks. Moreover, the robustness of the potential depends on the regression method and its associated regularization (used to alleviate overfitting issues). A simple (bi)linear method may be less accurate but more robust. It may also be sufficient if the descriptors already provide enough information on the system, as is the case for the descriptors obtained via convolutional neural networks in Ref. 27.

In principle, both neural network (NN) and non-linear kernel regression models are sufficiently sophisticated to obtain a trustworthy representation of scalar potential-energy surfaces (PES) or vector force fields of arbitrary complexity. However, in practice, choices have to be made for the similarity measure between atomic configurations (in both kernel regression methods and NN) or for the architecture of the neural network. The optimal choices are not the same for different systems, i.e., descriptors/parameters that work well for solids are not easily transferable to biological molecules and vice versa. Hence, many ML developments are currently specific to either organic molecules or materials. That being said, there is currently a growing interest in understanding the advantages and limitations of the different existing approaches (18, 27, 32, 33, 38–41) and developing truly general frameworks for learning complex PES or force fields that work seamlessly for both organic and inorganic matter.

***B.3. Current methods and their performances.*** We list some key methods in Table 1. The first successful ML approaches were developed to describe PES of defectless materials and their surfaces (32, 33, 38) with the goal to enable efficient and accurate Molecular dynamics (MD) of large supercells of elementary or binary materials. The Behler-Parrinello NN approach (33) or the kernel-based GAP approach of Csanyi (32) are both able to achieve accuracies of 1-2 meV/atom for some solids (C, Si, Cu, TiO2, among others). There are several key differences between these two methods, the main ones being the NN vs kernel approach and the different similarity measures between atomic configurations. Both approaches typically require on the order of tens to hundreds of thousands

reference calculations at the DFT level for constructing the training dataset, in order to achieve 1-2 meV/atom accuracy. Recently, PES-fitting methods based on deep networks have also been developed (41, 42). These approaches often do not require any *a priori* definition of the similarity measure; they are instead able to learn the similarity measure from the training data.

Constructing ML models for organic molecules is a field that faces somewhat different challenges compared to ML models for solids and materials. While DFT calculations are often deemed to provide sufficiently accurate reference data for solids, this is not the case for organic molecules. The "gold standard" is coupled cluster CCSD(T) computations. Quantum-chemical CCSD(T) calculations are however computationally expensive and it is only possible to carry hundreds of such calculations even for simple molecules such as aspirin. Early successful non-linear PES models were based on permutationally-invariant polynomials (PIP) (39). More recent developments include the so-called gradient-domain machine learning (GDML) approach (7, 40) for constructing molecular force fields. The GDML approach learns an explicit force field and obtains the PES via integration, instead of the more conventional approach to learning a PES and then taking its gradient to drive MD. This has two advantages: (i) the usage of an explicit Hessian kernel that provides the maximum flexibility, minimizes noise and prevents artifacts between forces and energies in the learning process; (ii) a significant gain in data efficiency, since globally accurate force fields for small molecules (accuracy of 0.2 kcal/mol and 1 kcal/mol/Å) can now be constructed using only a few hundred molecular conformations for training. This data efficiency currently enables the construction of essentially exact force fields for molecules with up to 30-40 atoms (7).

**C. Synergy between physics, chemistry, mathematics and ML approaches.** ML approaches used to construct accurate PES and force fields have already been successful and have enabled simulations of molecules and materials that were previously considered impossible. Ultimately, it would be worthwhile to achieve an optimal balance between physics-based models and ML approaches to enable not only faster and more accurate simulations, but also obtain insights into interactions of complex quantum-mechanical molecules and materials. For example, the GAP, Behler-Parrinello, GDML, and PIP approaches discussed above already incorporate translational, rotational, and permutational symmetries of molecules and materials in their internal representation of atomic interactions. Such symmetries were also made precise in the mathematical literature (18). In addition, by learning simultaneously energy and forces such that the latter are (minus) the gradient of the former, all of these methods enforce exactly energy conservation.

However, many more physical symmetries can and should be incorporated in ML approaches. For example, exact constraints are known for asymptotic forms of atomic interaction potentials. Also, some analytic and empirical results are known for series expansions of interatomic potentials. Finally, there are mathematical results which provide rigorous

| Method | Short description | Ref. |
|---|---|---|
| Kernel-based Gaussian approximation potentials (GAP) | Combines a structural descriptor and a kernel establishing the link between structure and energy | 32 |
| Behler-Parrinello NN | Feed-forward NNs for each atom. The potential energy is constructed as the sum of local atomic energies | 33, 38 |
| Deep NN (DTNN) | No a priori similarity definition needed, similarity is learned | 41, 42 |
| Permutationally-invariant polynomials (PIP) | Uses polynomials of Morse variables in fitting PES | 39, 43 |
| Gradient-domain ML (GDML) | Learns an explicit FF and obtains the PES via integration | 7, 40 |

**Table 1.** Summary of some key learning methods for force field (FF) development.

statements on the behavior of the potential energy functions in terms of the locality of the interactions (19). The incorporation of such prior knowledge could improve the efficiency and accuracy of ML potentials and ultimately also lead to novel analysis tools that offer new insights into the complex nature of atomic interactions (44).

It is also worth noting that electronic interactions in complex molecules and materials can be rather long-ranged. For example, electrostatic interactions and plasmon-like electronic fluctuations in molecules and nanostructures can lead to interatomic potentials extending to at least 20-30 nanometers (45, 46). Most current ML models explicitly or implicitly cut off interactions at an interatomic distance of 5-6 Å. Hence, by construction, these ML approaches are not able to capture interactions extending over larger length scales. For this reason, it is ultimately necessary to couple ML approaches that excel at capturing complex short-range chemical bonding with explicit physics-based approaches to non-covalent interactions. It is important to note that such physics-based models can also employ ML approaches to learn short-range interaction parameters based on datasets of electrostatic moments and polarizabilities. The recently developed IPML approach lies the foundation for unifying ML force fields and physics-based interatomic potentials (47). An alternative approach based on the definition of structure representations that incorporate long-range correlations with the correct asymptotic behavior (48) can simplify the simultaneous description of the multiple length scales contributing to molecular interactions.

**D. Perspectives for ML approaches to the determination of force fields.** We gather in this section some mathematical and numerical perspectives, as well as open problems, on ML methods for force fields:

- A first perspective is the use of ML to learn the difference between already acceptable empirical force fields and DFT models, as some form of preconditioning. Such an approach greatly depends on the regression method. For example, for kernel methods, it has been shown that a potential can be built on top of pre-existing two-body and three-body classical potentials, improving the overall accuracy (49, 50). On the

contrary, fitting differences between a good classical potential and an ab-initio potential with a linear regression yields very poor results, since the difference is small (almost noisy) and rugged (not smooth). It is observed that a simpler starting guess, such as the Ziegler–Biersack–Littmark potential (51), yields better results, since this increases the numerical stability and improves the accuracy.

- A question related to the robustness of these learning techniques is whether it would make sense to optimize potentials on a Pareto curve, where various properties of interest are weighted in different manners in the cost function. Indeed, the optimization is usually performed on a multi-objective cost function (including energy, force, stress, and sometimes bond distances, ...). The so-obtained potential is a result of the user arbitrary choice of the weighting parameters – infinitely many 'optimal' potentials can be obtained depending on the choice of the weights. The naturally rising question here is: is it possible to have a unified way of defining cost functions?

- An important practical concern is the sensitivity of the learnt parameters relatively upon the data (for instance depending on the fraction of elements used for training vs. testing).

- Another more theoretical question is: What is the numerical stability induced by machine learning potentials on the time integration of Hamiltonian dynamics and its variations? Indeed, some preliminary results suggest that machine learning potentials may be smoother than current empirical potentials.

- For reasons which remain to elucidate, predicting intensive (as opposed to extensive) properties seems to be very challenging.

**E. Bottom-up coarse-graining force fields: From PES to FES.** A classical particle-based coarse grained (CG) simulation model, where several atoms are grouped together, can be viewed as a reduction of the dimensionality of the classical phase space (see Figure 1). It requires the determina-
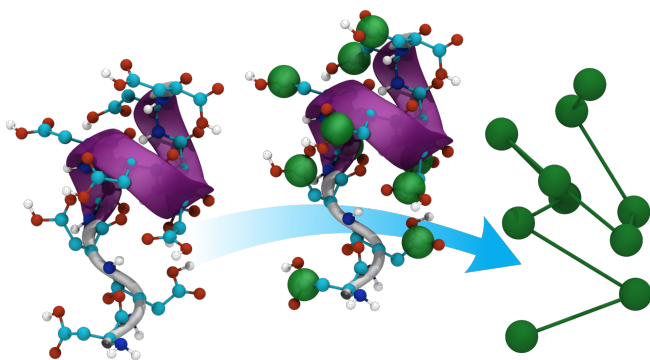
**Fig. 1.** Particle-based coarse-graining: high dimensional free energy surfaces (FES) can be extract from atomistic data and used as a basis for CG models (52, 53).

tion of an effective Hamiltonian that allows the model to explore the phase space in the same way as an atomistic simulation would. Thus, in the so-called bottom up coarse-graining strategies, the interactions in the CG model are devised such that an accurate representation of a (known) atomistic sampling of the configurational phase space (mapped to the CG representation) is achieved. These methods use the underlying multidimensional potential of mean force (PMF) derived from the atomistic simulation data as parameterization target, i.e., they try to reproduce a (typically high-dimensional) free-energy surface (FES) as opposed to a PES. Naturally, this is of particular relevance to the simulation of soft matter problems such as liquid state systems, soft materials and biological systems, where entropic effects, disorder and heterogeneity dominate the overall properties of the system.

Free energies and potentials of mean force are not a direct output of a MD simulation. They can be calculated by Boltzmann inversion of a (high-dimensional) probability density distribution obtained from sampling configurations in phase space or from mean forces acting on the interaction sites in the CG representation. In the past, several bottom-up coarse-graining methods have been derived which - while all aiming for an effective Hamiltonian that approximates a multidimensional PMF/FES - differ in terms of both the actual parameterization target (multidimensional PMFs/probability density distributions, structure functions as low-dimensional representations of these PMFs; mean forces in the direction of selected CVs or relative entropies) and the type of CG interactions which are typically represented by low-dimensional potentials, i.e., pair interactions, or three-body interactions) (54–58). Since these coarse-graining methods derive interactions from atomistic reference simulations, they are intrinsically data driven. Consequently, ML-based approaches yield new types of reference atomistic data and new types of CG interactions and parameterization methods. On the one hand, ML methods can be used to determine dimensionality-reduced representations of the phase space and to derive or validate CG models by matching the sampling of a (relatively complex) FES as opposed to low-dimensional target functions/properties. On the other hand, ML methods can also be employed to identify suitable CVs that describe the states and the dynamics of a system, which can then either be di-

rectly used in the CG potentials or be employed to identify optimal CG representations and learn CG interactions. This is discussed at length in Section 3.

Following the methodology of inferring all-atom potential energy functions from corresponding quantum mechanical data, John and Csanyi have extended the Gaussian Approximation Potential (GAP-CG) approach to coarse-graining of simple liquid systems (59). In this case, the many-body PMF is described via local multibody terms, based on local descriptors and multidimensional functions which are determined by Gaussian process regression from atomistic training data (instantaneous collective forces or mean forces). In a similar vein, Zhang et al. developed a scheme, called the Deep Coarse-Grained Potential (DeePCG), which uses a NN to construct a many-body CG potential for liquid water (60). The network is trained with atomistic data in a manner similar to the force matching in the multi-scale coarse-graining method (61), and in such a way that it preserves the natural symmetries of the system. While the described two methods are related to the force-matching type of bottom-up coarse-graining and use ML to significantly extend the complexity of the CG interactions, Lemke and Peter follow a different strategy (52). A NN is used to extract high-dimensional FES from atomistic MD simulation trajectories. The NN is trained to predict conformational free energies by creating a classification problem between real MD conformations and fake conformations of a known distribution. With such a classification based procedure it is possible to train the NN to return probability densities without requiring any binning or normalization – which circumvents the problem of binning in high dimensional space (62). By using the NN probability densities directly in a Monte Carlo type of sampling of conformations, a (relatively) high-dimensional FES is thus used as effective CG Hamiltonian. This NN network model was successfully tested for several homo-oligopeptides (53). By employing a convolutional NN architecture, the NN model could be simultaneously trained on data of different chain lengths and could even make meaningful predictions for polymers with chain lengths different from the ones in the training data. Thus, such an approach is promising for the simulation of polymer systems where naturally training data are restricted to chain lengths that are shorter than the intended polymers.

Coarse-graining of potential energy functions into free energy type interactions has a well founded statistical interpretation. A difficult question is however whether some dynamical properties are also preserved in this coarse-graining process, and to which extent.

## 3. Dimensionality reduction and identification of collective variables

The objective of this section is to discuss various techniques to identify collective variables. After some general considerations in Section A, we first present the main two ideas to build collective variables in Section B, namely looking for high-variance or slow degrees of freedom. We then discuss how this can be used to enhance the sampling of the canonical ensemble on the example of diffusion maps in Section C,

before discussing dynamical aspects in Sections D and E.

**A. General considerations.** Molecular systems are characterized by the fact that their long-time dynamical behavior is typically governed by a small number of emergent collective variables (CVs) (63–65). These collective modes arise from cooperative couplings between the constituent atoms induced by interatomic forces (e.g., covalent bonds, electrostatics, van der Waals interactions) and possibly external fields (e.g., electric fields, hydrodynamic flows), and which render the effective dimensionality of the system far lower than that of the full-dimensional phase space in which the system Hamiltonian and equations of motion are formulated (64, 65). In a dynamical systems sense, the long-time evolution of the system is restrained to a low-dimensional attractor or intrinsic manifold and its dynamics over these time scales may be described within the Mori-Zwanzig projection operator formalism as evolving within a subspace of slow collective variables to which the remaining degrees of freedom are effectively slaved (64).

Traditional unbiased MD is not able to efficiently explore the whole kinetic landscape with time scales spanning over orders of magnitude, from picoseconds to milliseconds. In this scenario, one relies on extensive simulations together with some clever strategy to escape metastable states. Such a strategy can only be devised if one is able to identify what defines a "long-lived" state, which is equivalent to discovering meaningful collective variables (CVs) or reaction coordinates (66). The methods described below aim at finding these CVs or states. As will become clear later, depending on the objective, the focus may be different: gain insight/intuition on the system, bias to exit metastable states, compute a free energy profile, set up a coarse-grained dynamics simulation, cluster/classify configurations, etc.

**B. Data-driven discovery of high-variance and slow collective variables.** The inherently multi-body and emergent nature of the CVs means that they are exceedingly challenging to intuit for all but the most trivial systems, and data-driven techniques present a powerful means to systematically estimate them from molecular simulation data. The origins of this data-driven approach can be traced back to pioneering work in the early 1990's by Toshiko Ichiye and Martin Karplus (67), Angel Garcia (68) and Andrea Amadei, Antonius Linssen and Herman Berendsen (69) who applied PCA to molecular simulations of protein folding. Since that time there has been an explosion of interest in the use of data science and machine learning techniques to estimate CVs from molecular simulation data and the subsequent use of these CVs to inform new understanding, perform molecular design, and guide enhanced sampling.

Data-driven CV discovery typically employs unsupervised learning techniques that seek low-dimensional parameterizations of the geometry of the data in the high-dimensional phase space of atomic coordinates (70). This procedure can usually be cast as an optimization problem that maximizes some objective function, or equivalently minimizes some loss function, over the data. The techniques can be categorized into linear and nonlinear methods. Linear techniques are restricted to discovering CVs that are linear combinations of the input features, whereas nonlinear techniques can discover more general nonlinear functional relations. The more powerful and general nonlinear techniques are typically better suited to the estimation of the complex emergent CVs in molecular systems, but linear techniques should not be discounted since they are typically more robust, interpretable, and less data hungry, and can also admit nonlinearities through feature engineering or the kernel trick (71). The importance of the choice of features in which the molecular system is represented to the CV discovery tool should not be underestimated. Feature sets that contain and foreground the important molecular behaviors and respect fundamental symmetries (e.g., translation, rotation, permutation) can be critical to the success of CV discovery (particularly in the case of linear techniques), whereas poor choices that mask or discard essential information or contain spurious symmetries can easily produce poor performance. What constitutes a good choice of feature set is strongly system dependent and is typically reliant on some combination of intuition, experience, and exploratory trial-and-improvement. We refer for example to Ref. 72 for a discussion on the importance of the choice of the representation of the data.

Although the details and specifics differ, most CV discovery techniques can be placed in one of two categories: those that seek high-variance CVs and those that seek slow CVs (see Figure 2).

High variance CVs maximally preserve the configurational variance in the high-dimensional data upon projection into the low-dimensional space spanned by these CVs. Slow (i.e., maximally autocorrelated) CVs define a low-dimensional space that maximally preserves the long-time kinetics of the system. Frequently the slow and high-variance collective modes are related, but this is not always the case. Importantly, the estimation of slow CVs requires data arranged in time series (e.g., MD trajectories) whereas the estimation of high-variance CVs can be applied to data sampled without temporal ordering (e.g., Monte Carlo trajectories). Notice however that methods exist to recover dynamical information according to some artificial dynamics (e.g. reversible purely diffusive dynamics) upon non-time ordered data to render it amenable to temporal analysis techniques (73).

Let us also mention that recent advances in deep reinforcement learning (DRL) in robotics opens up new avenues for deploying DRL to atomic and molecular systems. In all DRL algorithms, a reward function, state and action space should be defined. In atomic systems, state space can be atomic coordinate, action space can be the movement of atoms, and reward can be defined as energy. DRL can be suitable replacement for finding transition paths and can potentially be used to strengthen the string or nudged-elastic-band method (74, 75).

Before giving more details about the high-variance and slow CVs, let us mention that a widespread definition of an optimal *scalar-valued* reaction coordinate in the rare event-field is
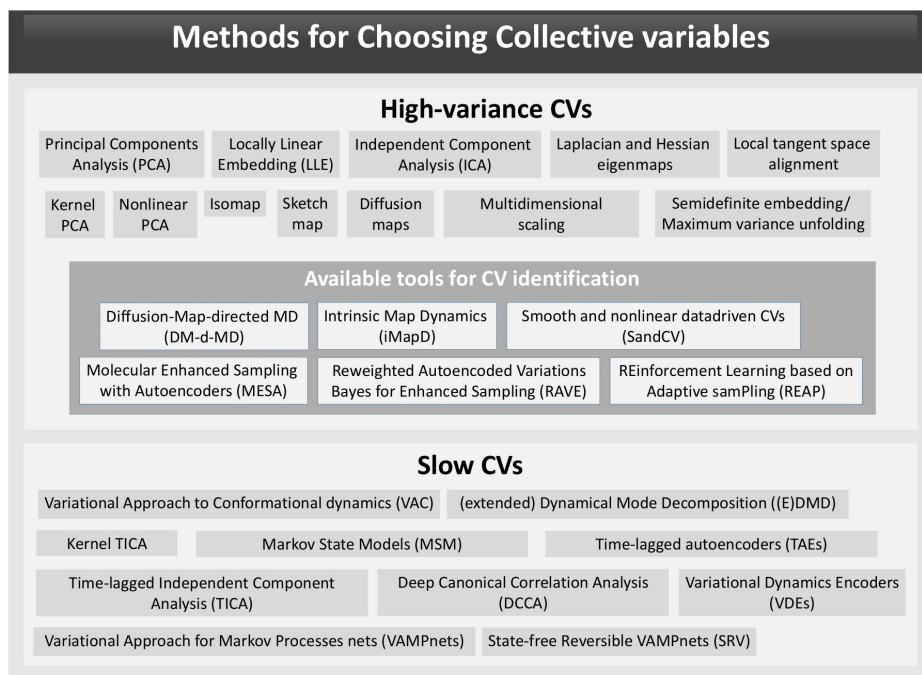
**Fig. 2.** Representative methods for CV identification. All related citations are in the main text.

the committor function, i.e., in a system with two metastable states, the probability that a given atomic configuration will evolve towards the products before reaching the reactants. Such probability can in principle be estimated by generating a huge number of MD simulations from each configuration of interest: even if such a procedure cannot be applied in practice to the whole configuration space, the committor represents an ideal reaction coordinate in some sense (we refer the reader to (76) or (77, p.126) for example) and provides tests and optimization strategies for candidate CVs (5, 17, 76, 78–80).

***B.1. High-variance CV estimation.*** The best known high-variance CV estimation technique is PCA (10), also known as the Karhunen-Loève transform (81–84), or proper orthogonal decomposition (85, 86). This approach discovers an orthogonal transformation of the input data to define a hyperplane approximation that preserves most of the variance in the data. Popular nonlinear techniques for high-variance CV estimation include kernel and nonlinear PCA (87–90), independent component analysis (ICA) (91), multidimensional scaling (92), sketch map (93) locally linear embedding (LLE) (94, 95), Isomap (96–98), local tangent space alignment (99), semidefinite embedding / maximum variance unfolding (100), Laplacian and Hessian eigenmaps (101, 102), and diffusion maps (11, 103). These approaches differ in their mathematical details, but can be broadly conceived of as nonlinear analogs of principal component analysis that pass curvilinear manifolds through the data to define nonlinear projections into a low-dimensional subspace spanned by the learned CVs. Specialized techniques for molecular simulations that integrate iterative high-variance CV discovery and accelerated sampling of configurational space have been developed in recent years (13–15, 104–114).

The techniques described above can be coupled with enhanced sampling methods, which use the uncovered CV's to help the system leave metastable states. In this case, one actually relies on CV estimates based on partial sampling (73). Let us describe a few methods in that direction.

Diffusion-map-directed MD (DM-d-MD) uses diffusion maps to identify CVs spanning the range of explored system configurations and then initializes new simulations at the frontiers of this domain to drive sampling of new system configurations (113, 114). Intrinsic map dynamics (iMapD) employs diffusion maps to construct a nonlinear embedding of the high dimensional simulation trajectory and then uses boundary detection algorithms with a local principal components analysis to extrapolate into new regions of phase space at which to seed new simulations (105). The Smooth And Nonlinear Data-driven Collective Variables (SandCV) approach identifies nonlinear CVs using Isomap, expands them within basis functions centered on a small number of landmark points, and then passes this parameterization to the adaptive biasing force accelerated sampling technique to drive sampling along these coordinates (109). Molecular enhanced sampling with autoencoders (MESA) employs autoencoding neural networks to discover nonlinear CVs for enhanced sampling without the need for approximate basis function expansions (13, 14). Reweighted Autoencoded Variational Bayes for Enhanced Sampling (RAVE) employs variational autoencoders to discover nonlinear CVs that are compared at the level of their probability distributions with an ensemble of physical candidate variables to identify physical coordinates for accelerated sampling (15). REinforcement learning based Adaptive samPling (REAP) employs reinforcement learning to identify the dynamically-varying relative importance in driving exploration of configurational

space of each CV within a candidate set and then adaptively seeds new simulations from configurations with high reward functions (104).

***B.2. Slow CV estimation.*** The identification of slow CVs is valuable and informative from many perspectives. From a mechanistic perspective, these CVs reveal the collective modes that dictate the metastable states of the system and the transitions between them. From a design perspective, they can offer a blueprint for the structural, thermodynamic, and dynamic properties of the system. From an enhanced sampling perspective, they provide good variables in which one can apply biases to accelerate barrier crossing and improve exploration of configurational phase space.

A number of approaches have been proposed to analyze MD time series to estimate slow CVs. The theoretical basis for these techniques is founded in the variational principle of conformational dynamics (VAC) (115), or in the (extended) dynamical mode decomposition ((E)DMD) (116, 117) that, respectively, frame the recovery of the slow CVs as a variational optimization or regression problem (16, 118). Shortly, VAC estimates the slowest modes as linear combinations of *a priori* defined basis functions of the input coordinates. In Time-lagged independent component analysis (TICA) these basis functions are the coordinates themselves (115, 119–125). In Markov state models, the slow CVs are approximated in a basis of indicator functions defined over the data (118, 126) (see also the recent special issue Ref. 127 for the latest developments on Markov state models). Perron cluster analysis can be used to reduce the large number of states uncovered by clustering methods along the trajectory, to a few metastable states, see Ref. 128–130. Combining TICA with the kernel trick yields kernel TICA (kTICA) that is capable of approximating the slow CVs with nonlinear functions of the input features (115, 131). Deep canonical correlation analysis (DCCA) (132), the variational approach for Markov processes nets (VAMPnets) (133), and state-free reversible VAMPnets (SRV) (134) all employ Siamese neural networks to learn nonlinear featurizations of the input coordinates as basis functions with which to approximate the slow CVs. Time-lagged autoencoders (TAEs) employ time-delayed autoencoding neural networks to learn slow CVs into which the molecular trajectory can be projected (i.e., encoded) and also used to predict the system state at the next time increment (i.e., decoded) (16). Variational dynamics encoders (VDEs) are similar to TAEs but employ a variational as opposed to traditional autoencoding architecture that introduces stochasticity into the decoding of the learned CVs (135, 136).

Enhanced sampling can be conducted in the learned slow CVs in a similar manner to that in the high-variance CVs, but the application of artificial biasing potentials perturbs the true system dynamics and subsequent applications of slow CV estimation techniques to the biased data must compensate for this effect (137–139).

## C. Enhanced sampling using local and global diffusion maps.
Using the illustrative example of diffusions maps, we discuss in this section how to use the proposed

reaction coordinate to enhance sampling and somehow perform some extrapolation procedure. Diffusion maps are a dimensionality reduction technique which allows for identifying the slowly-evolving principal modes of high-dimensional molecular systems (11, 12). It does so by computing an approximation of a Fokker-Planck operator on the trajectory point-cloud sampled from a probability distribution (typically the Boltzmann-Gibbs distribution corresponding to prescribed temperature). The construction is based on a normalized graph Laplacian matrix. In an appropriate limit, the matrix converges to the generator of overdamped Langevin dynamics. The spectral decomposition of the diffusion map matrix thus yields an approximation of the continuous spectral problem on the point-cloud (140) and leads to natural CVs.

Since the first appearance of diffusion maps (11), several improvements have been proposed including local scaling (141), variable bandwidth kernels (142) and target measure maps (TMDmap) (143). The latter scheme extends diffusion maps on point-clouds obtained from a surrogate distribution, ideally one that is easier to sample from. Based on the idea of importance sampling, it can be used on biased trajectories, and improves the accuracy and application of diffusion maps in high dimensions (143).

Several algorithms have used diffusion maps to learn the CVs adaptively and thus enhance the dynamics in the learned slowest dynamics (13, 105, 113, 114). These methods are based on iterative procedures whereby diffusion maps are employed as a tool to gradually uncover the intrinsic geometry of the local states and drive the sampling toward unexplored domains of the state space, either through sequential restarting (114) or pushing (105) the trajectory from the border of the point-cloud in the direction given by the reduced coordinates. All these methods try to gather local information about the metastable states to drive global sampling. In (73), the authors focused on the construction of diffusion maps within a metastable state by formalizing the concept of a local equilibrium based on the *quasi-stationary distribution* (144). This local equilibrium guarantees the convergence of the diffusion map within the metastable state. Moreover, the work provides the analytic form of the operator obtained when metastable trajectories are used within diffusion maps.

Finally, since the collective variables provided by diffusion maps are only defined on the sampled point cloud, one must apply extrapolation approaches. These might be very noisy and, more importantly, lose their meaning outside the convex hull of the point cloud. As a remedy, diffusion maps could be used as a tool to select collective variables from a database of physical reaction coordinates, similarly to (17), providing more physical insight into the abstract collective variables. This approach would allow to evaluate the CV outside the point cloud and provide more physical meaning into the abstract collective variables.

The local-global perspective has motivated a method allowing on-the-fly identification of metastable states as an ensemble of configurations along a trajectory, for which the diffusion map spectrum converges. Secondly, an enhanced sam-

pling algorithm based on QSD and diffusion maps has been proposed. For the latter, the main idea is a sample from the QSD allowing to build high-quality local CVs (within the metastable state) by considering the most correlated physical CVs to the diffusion coordinates. Once the best local CVs have been identified, one can use existing methods as metadynamics to enhance the sampling, effectively driving the dynamics to exit the metastable state. The authors in (73) demonstrate this idea on a toy-model example showing improved sampling over the standard approach.

Diffusion maps can also be used to a compute the committor function (145), which provides dynamical information about the connection between two metastable states and can be used as a reaction coordinate. Markov state models (MSM) can in principle be used to compute committor probabilities (146), but high dimensionality makes grid-based methods intractable. Similar work in this direction was done by (145, 147, 148). Diffusion-maps, especially the TMDmap (143), can be used for committor computations in high dimensions. The low computational complexity aids in the analysis of molecular trajectories and helps to unravel the dynamical behaviour at various temperatures.

As a future work, the quality of the diffusion map approximation could be improved by introducing more sophisticated kernels or point-cloud approximations similarly to (145). Also, diffusion maps could be extended to the approximation of generators of the underdamped Langevin dynamics.

### D. Extracting dynamical information from trajectory data.
Once good CVs or metastable states have been identified, these can be used to extract dynamical information. Let us describe in this section the approach followed by Thiede *et al.* (147), which is based on a Galerkin projection of the infinitesimal generator.

The approach in (147) builds on the MSM and related frameworks (115, 117, 128, 149–154). Dynamical statistics of interest are cast as solutions to equations involving the generator, i.e., the operator that describes the evolution of functions of the dynamics over infinitesimal times. Although the full generator cannot be determined in general, the equations can be solved by a Galerkin approximation. In this approximation, the dynamical statistic of interest is expanded in terms of a basis, and its generator equation is reduced to a linear form. The contributing matrix elements (inner products of basis elements and the generator) can be estimated from short MD trajectories. A key challenge is to generate basis sets consistent with the boundary conditions. Thiede *et al.* (147) considered two basis sets: indicator functions that reprise MSMs and diffusion maps (11). The latter showed promise for capturing smoothly varying dynamical statistics, such as committors and mean first-passage times with fewer basis functions, but the efficiency of a given basis is likely to be problem specific. Because the dynamical Galerkin approximation framework generalizes the notion of transition between states, the sampled configurations can be replaced by short trajectory segments. This allows treating memory that arises from incomplete description of the system by delay embedding (155, 156). This is an appealing alternative to

extending the lag time in an MSM because it does not sacrifice time resolution. Going forward, it will be interesting to investigate whether variational methods akin to those for elucidating time scales (115, 133) can be developed to permit representation of the dynamical statistics in terms of nonlinear functions.

### E. Tackling both Markovian and non-Markovian cases: Free energy, friction and mass profiles extracted from short MD trajectories using Langevin models.
In principle, the high-dimensional dynamics of a system composed by many atoms, when projected onto one (or a few) CV, can be modeled by a generalized Langevin equation (157, 158). Such stochastic differential equations contain several ingredients: a mass, a drift term corresponding to the mean force (gradient of the free energy landscape), a friction and a noise. Projecting on a low-dimensional space yields, in general, non-Markovian dynamics, except in the presence of time scale separation between CVs and bath coordinates and at coarse time resolution (157).

Clearly, the construction of optimal Langevin models along meaningful reaction coordinates is appealing from several viewpoints (159). On one side, the complex many-body dynamics is approximated by an equation that preserves physical intuition and is cheap to integrate. On the other side, exact kinetic rates - free from transition state theory approximations - between metastable states can be accessed more easily, by exploiting brute-force Langevin simulations or more elaborate methods (160). Compared to Markov state models, Langevin models are not restricted to Markovian dynamics and do not require the discretization of configuration space and the choice of a lag time, which are customary sources of errors.

For all these reasons, several algorithms have been developed to recast MD data into low-dimensional Langevin models (161–172). Usually, with these techniques, the terms of the Langevin equation are estimated employing very long equilibrium MD trajectories that ergodically sample the whole relevant free energy landscape. Of course such data are seldom available in complex applications featuring rare events, strongly limiting the scope to the case of barriers smaller than a few $k_B T$. Tackling the more general case of limited sampling and non-equilibrium MD trajectories is much more involved (173).

A possible and simple solution to this challenge - especially in the context of rare events - has been proposed in Ref. 174: the parameters of a generalized Langevin equation are optimized by minimizing the error between MD and Langevin probability distributions $P(x, \dot{x}, t)$ along the reaction coordinate $x$. Such out-of-equilibrium distributions are estimated from a set of short unbiased trajectories initiated close to a barrier top (with random thermal velocities) and allowed to relax into the adjacent free energy minima, in the spirit of committor analysis (a preliminary exploration of putative transition state structures can be nowadays performed at a moderate cost using, e.g., the prejudice-free techniques of Ref. 175–177).

Employing both benchmark models and solvated proline

dipeptide as a test case, numerical evidence indicates that ~100 short trajectories (of few picoseconds in the typical case of a small solute in water) encode all the information needed to reconstruct free energy, friction, and mass profiles (174). This approach, suitable also for high barriers of tens of $k_B T$ and non-Markovian dynamics, provides the thermodynamics and kinetics of activated processes in a conceptually direct way, employing only standard unbiased MD, at a competitive cost with respect to existing enhanced sampling methods. Furthermore, the systematic construction of Langevin models for different choices of CVs starting from the same initial data could help in reaction coordinate optimization.

## 4. Application of machine learning techniques in biological systems and drug discovery

Two of biology's biggest challenges are the prediction of protein structure based on its amino acid sequence, i.e., protein folding, as well as the dynamical conformational changes of the three-dimensional structure of proteins, i.e., protein dynamics. Beyond the actual problem of protein folding, which was recently set at a different basis after the breakthrough from AlphaFold and the impressive one million time faster Artificial Intelligence (AI) solution by AlQuraishi (178), the prediction of protein dynamics and mechanism of action is possible through the use of MD simulations.

Recent advances in computer hardware and algorithms have led to simulations of protein dynamics of size and time lengths that are intrinsic to biological processes. Dynamics of protein plasticity and drug binding/unbinding mechanisms are a few of the key processes that we would ideally like to capture through these large scale simulations. However, the analysis and interpretation of the large amount of data that are produced by these simulations is complex and should be carefully considered (179).

As discussed in Section B, despite the ever-growing time and length scales of simulations, unbiased MD is not able to explore the whole kinetic landscape of complex systems and carefully chosen, meaningful CVs can be used to represent the free energy surface of these systems in order to reveal the regions of low energy, i.e., stable and metastable states, as well as the barriers, i.e., transition states, between these regions (163, 169, 180). ML approaches have recently started being used for the discovery of meaningful CVs (14, 15, 133, 181, 182), while iterative schemes where CVs are being updated based on new simulation data provide promising results for challenging systems (181, 183, 184).

In this section, we first present an example of dimensionality reduction for building a Markov State Model for the study of lysine methyltransferase SETD8 (see Section A). We next present some biological examples were adaptive MD/ML techniques can help gain access to non-crystallographic conformational states of disease-related proteins for drug discovery purposes (see Section B). In Section B.1, we discuss the possibility of conformational-specific targeting of pro-

teins using their metastable states as target conformations, while in Section B.2 we give some examples were ML techniques applied in MD simulations can provide information about potential allosteric binding sites or protein activation mechanisms upon ligand binding.

**A. Selection of efficient collective variables for MSMs: the example of SETD8.** Conformational changes in proteins span from thermal fluctuations of side chains and motions of active loops to major rearrangement of sub-domains, including unfolding and refolding processes (185). The ability to unveil the mechanisms underlying protein function requires quantifying the importance of these motions for the process of interest or, in other words, obtaining a representative ensemble of conformations.

Besides the relevance for devising enhanced sampling strategies, the discovery of CVs is decisive when analyzing simulation data sets by using, for instance, Markov State Models. In this context, the conformational study of the protein methyltransferase SETD8, an epigenetic enzyme essential in the regulation of the cell cycle, was discussed in (183).

SETD8 is characterized by a dynamically rich behavior, which has proven to be essential in enzymatic catalysis (186). In (183) the authors combined experiments and simulation in an attempt to span the up-to-that-time unexplored configurational space of SETD8. Several new X-ray structures were obtained by trapping conformations with small-molecule ligands (187). These, in turn, were used to build hypothetical structures by manually combining fragments observed in experiments.

The set of initial configurations was used to seed independent MD simulations in explicit solvent, resulting in an extensive simulation database. The search of reaction coordinates was done in different spaces of residue-residue distances, logistic distances, and backbone dihedrals. These CVs, usually referred to as "features" in the MSMs literature, are arbitrary choices, that have been traditionally based on human intuition and heuristics (188). This is arguably the "achilles heel" of MSMs and has prompted the development of ML approaches to bypass human intervention (16, 133).

Although a set of features is already a space with much fewer dimensions than the full atomic coordinates, it is still a high dimensional system that cannot be handled with MSMs. This requires further dimensionality reduction, which can be done using, for instance, the time-lagged independent component analysis (tICA), discussed in Section B.2. CVs obtained by tICA are linear combinations of features that, in principle, encompass the variance of the data while providing time scale separation. These are attributes of meaningful CVs (182), which explains the consensus regarding tICA as a suitable strategy for building MSMs (119, 124, 188, 189). The stage regarding data representation ends with clustering the conformational snapshots into discrete states using unsupervised ML protocols, such as the k-centers and k-means methods (190).

Given the multiple subjective decisions involved in selecting features and algorithms to represent the database, MSMs building must be allied with validation strategies. In this

context, Husic *et al.* (188) emphasize the importance of using a kinetically-motivated dimensionality reduction and cross-validation strategies to avoid over fitting. The study of SETD8 (183) uses both structural and kinetic criteria, and 50:50 shuffle-split cross-validation scheme with random divisions of the data into training and test sets (see Figure 3). As a result of such an extensive validation, the specific study successfully quantified an ensemble of kinetically relevant macrostates which, in addition, were validated with experiments.

**B. Machine learning-driven MD simulations in drug discovery.** The discovery of a new drug is a long, multi-step and expensive process. Any tool that can speed up any of the steps involved would have big implications down the entire drug discovery chain. Artificial intelligence is expected to significantly shape the future of many aspects of drug discovery during the forthcoming decades. It is already used to design evidence-based treatment plans for cancer patients, instantly analyze results from medical tests to escalate to the appropriate specialist immediately, and most recently to conduct scientific research for early-stage drug discovery.

Proteins, the most common drug targets, are dynamic molecular machineries whose function is intimately linked to their conformations. Destabilization of the subtle equilibrium of protein conformations can lead to severe pathologies, like in the well-known cases of KRAS G12X oncogenic mutations and prion disease. In this context, knowledge of the conformational landscape of targeted proteins would provide an outstanding advantage for the design of novel and original compounds stabilizing specific conformations of the protein (191).

Experimentally, the protein conformational space is often limited to few conformations that have been prone to crystallize. The use of GPUs and massive computational resources has enabled for the *in silico* alternative, MD simulations, to gain an important place in the first steps of drug discovery. Nevertheless, MD is limited to a few hundreds of microseconds of simulation, which limits the conformational space exploration.

New molecular modeling approaches combining MD simulations and ML techniques can help gain access to these non-crystallographic conformational states of a target protein. This knowledge would allow focusing on specific conformations of the protein in order to alter or restore its function. ML techniques can enable us to identify patterns in simulation data, build models that explain the different conformational states of a target and predict potential target-specific solutions for their druggability (13, 15, 181, 182, 184, 192–195).

As discussed in Section A, good CVs can guide enhanced sampling MD simulations in order to gain insights into long timescale dynamics of biomolecular systems. The difficulty of the identification of such CVs and in most cases the complexity of their definition has limited the number of available software for this purpose. PLUMED is an open-source, community-developed library that has been widely used in enhanced-sampling simulations of complex biological systems in combination with many MD engines, e.g., Amber, GROMACS, NAMD, and OpenMM (196–200). Most importantly, PLUMED can be interfaced with the host code using an API, accessible from multiple languages, including C++ and Python). This last functionality is important for adaptive protocols used for the identification of optimal CVs using iterative learning algorithms based on well developed ML libraries like Keras (201), TensorFlow (202), PyTorch (203) and Fastai (204). The MSM Builder package provides the user with software tools for predictive modeling of long timescale dynamics of biomolecular systems using statistical modeling to analyze physical simulations (205). Other tools that can be employed in MD/ML studies include among others MDTraj (206), ColVar module for VMD (192), OpenPathSampling (207).

***B.1. Conformational-specific targeting of proteins using cryptic binding sites.*** Drugs are traditionally designed to bind to the primary active site of their biological targets in order to induce a therapeutic effect. However, the high similarity between the orthosteric pockets among most of the protein families, leads in several cases to adverse effects. A new emerging direction in drug discovery is the use of alternative, transient, non-orthosteric binding sites that are not apparent in the protein's known crystallographic conformations and where small molecules can bind and modulate the biological target's function.

By binding to non-orthosteric sites of proteins, allosteric inhibitors can also exhibit a better selectivity vs proteins from the same family, as illustrated by SAR156497, a highly selective inhibitor of Aurora kinases (208). Well known drugs on the market work through this kind of mechanism of action (e.g., Lapatinib or Imatinib), but this mechanism was described *a posteriori*. Moreover, there are approved allosteric modulator drugs such as Cinacalcet for the treatment of hyperparathyroidism and Maraviroc for the treatment of AIDS, as well as many candidates at different stages of development (209, 210). Another aspect in targeting non-orthosteric pockets in drug discovery relies on the fact that allosteric inhibitors will not compete with endogenous ligands for binding, which can be critical when such endogenous ligands have very strong affinity for their protein.

One of the successful efforts in this direction is the example of PI3K$\alpha$, where a novel non-orthosteric pocket was identified using molecular dynamics (MD) simulations (211, 212). In (211), the authors used Functional Mode Analysis (213) and identified two dominant motions of PI3K$\alpha$ that influence both the active and allosteric pockets and are distinct between the wild-type protein and its oncogenic counterpart. Current work aims at extending this approach to other protein targets, where neural networks are employed in order to establish the link between oncogenic mutations and the protein's mode of action, with an ultimate goal to identify druggable mutant-specific conformations.

Beyond single protein conformations, multimeric protein assembly also appears as a challenging area where ML could play a role in drug discovery. The recent example on TNF$\alpha$ for instance shows the importance of how subtle
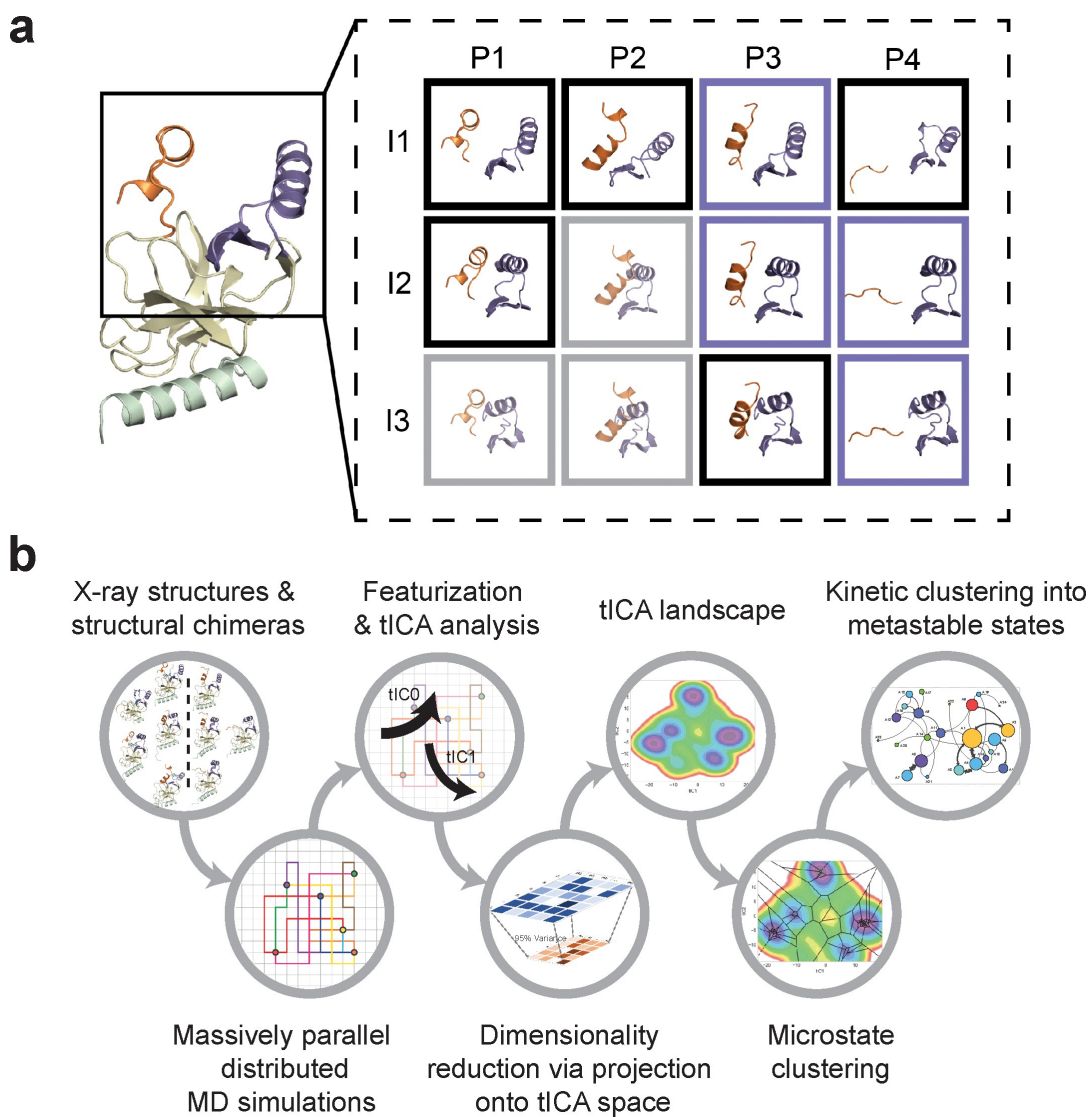
**Fig. 3.** Construction of conformational landscapes of apo- and SAM-bound SETD8 through diversely seeded, parallel molecular dynamics simulations and Markov state models.(a) Combinatorial construction of structural chimeras using crystallographically-derived conformations. (b) Workflow for dynamic conformational landscapes construction using MSM. For more information we refer the reader to the original publication 183. (Image source: Ref. 183. Use permitted under the Creative Commons Attribution License CC BY 4.0., https://creativecommons.org/licenses/by/4.0/).

changes in protein conformation can translate into a distorted trimeric assembly of TNF$\alpha$, impacting downstream signaling of TNFR1. Small compounds stabilizing this asymmetrical TNF$\alpha$ trimer can then be designed to treat or prevent TNF$\alpha$-related diseases (214).

***B.2. Compound-specific effect of binding.*** Another promising direction in the drug discovery process is the compound-specific effect of protein binding (215, 216). For example, a small organic compound can be used to boost the enzymatic activity of a protein enzyme or evaluate allosteric binders by the stabilization of its active conformation. In finding allosteric binding sites, ML algorithms such as k-means and Markov Models can significantly help in reducing the dimensions of drug binding events. The connections between statistical mechanics principles, such as Boltzmann Machines, and the discovery of the binding sites in proteins can be insightful. As an example, one can run thousands of small tra-

jectories of drug binding and unbinding events and learn the reaction coordinates using tICA (time-independent Component Analysis) in order to find the possible allosteric binding sites (215). These trajectories can be generated using different initial seeds (both different locations and orientations) and may range from 50 ns to 500 ns.

In the activation pathway of many proteins such as G Protein Coupled Receptors (GPCRs), the conformational changes are subtle and are limited to the sequential motion of residue switches triggering a signal from ligand to intracellular motifs. Finding these intricate motions in high dimensional space requires ML techniques to reduce the system's dimensions (216). Among these methods, variational autoencoders (VAE) and tICA (sparse or kernel) can be used to achieve learning and finding the reaction coordinates for such complex proteins.

## 5. Concluding remarks and perspective

Let us conclude this review by presenting some global perspectives on the interactions between machine learning approaches and molecular simulation, which are common to all the situations we discussed – from devising numerical potentials based on ab-initio reference data to the identification of collective variables in actual simulation of biological proteins.

First, we have seen that the aims of the coarse-graining procedures may be very different in nature. From the material presented in this review, one can identify three major purposes: (1) *a modeling objective*: using machine learning techniques to improve models, for instance by better representing force fields and potential energy surfaces; (2) *a numerical objective*: improving the efficiency of numerical methods, for instance by devising good collective variables to be used in conjunction with enhanced sampling techniques, such as free energy biased sampling techniques; (3) *a data analysis objective*: providing an efficient post-processing tool, as for instance a Markov state model to interpret the raw simulation data from molecular dynamics and identify states of interest. Concerning the choice of the learning methods, some common trends are shared by all methods, namely ensuring that one has access to a sufficiently rich database (sufficient variability of configurations for force fields, long reactive trajectories to identify CVs) and representing correctly the data (starting possibly with some putative CVs/descriptors, and then using some regression from there to sparsify/optimally combine these initial guesses). The precise choice of the learning method and the reduced model to work with, however, depend very much on the goal and priority of the user, and the system under consideration. The priority can be *the accuracy* (being as precise and as close as possible to some reference model, e.g., all-atom results when coarse-graining, or reproducing DFT energies when constructing numerical potentials), *the transferability* (learning how to coarse-grain small systems and extending the method to larger ones, learning energies at a given temperature and using the potential at another one) or the CPU/GPU *computational cost*.

When using black box learning techniques, based for example on neural networks, a problem which is often raised is the *interpretability* of the result. This is discussed for example in (80) which attempts to reconcile machine learning models (specifically a neural network approach to optimal reaction coordinates) with physical insight by means of symbolic regression techniques, also known as genetic programming. Such techniques appear very promising for the future, being able to distill fundamental natural laws from numerical data (217).

Another important element is the *reproducibility* of the results: one should favor approaches which are easy enough to cross-check and to repeat on various architectures. This also requires the researchers to ensure that the coarse-graining technique they propose yield robust results. For example, the results should not depend on the initial weights in a neural network, or on the sampled point used as inputs. Finally, this includes considering well established databases, or making

databases available to other users/developers; and also relying on standard and well maintained packages when using external libraries.

One idea which would help setting up common benchmarks and/or agreeing on common aims/priorities would be to organize some competition or prediction contest, which should ideally be simple enough so that even small groups can participate since this requires agreeing on common goals. Setting up the rules of such a competition would already be quite an achievement. Another important idea would be to emphasize transferability in all approaches, and more systematically work with some databases of some sort and then test on different databases.

## 6. Bibliography

1. Yue-Yu Zhang, Haiyang Niu, GiovanniMaria Piccini, Dan Mendels, and Michele Parrinello. Improving collective variables: The case of crystallization. *The Journal of Chemical Physics*, 150(9):094509, 2019.
2. J. Behler. Atom-centered symmetry functions for constructing high-dimensional neural network potentials. *The Journal of Chemical Physics*, 134(7):074106, 2011.
3. Albert P. Bartók, Risi Kondor, and Gábor Csányi. On representing chemical environments. *Phys. Rev. B*, 87:184115, 2013.
4. D. J. Wales. Perspective: Insight into reaction coordinates and dynamics from the potential energy landscape. *J. Chem. Phys.*, 142(13):130901, 2015. doi: 10.1063/1.4916307.
5. Baron Peters. Reaction coordinates and mechanistic hypothesis tests. *Annu. Rev. Phys. Chem.*, 67(1):669–690, 2016.
6. Robert T McGibbon, Brooke E Husic, and Vijay S Pande. Identification of simple reaction coordinates from complex dynamics. *J. Chem. Phys.*, 146(4):044109, 2016. doi: 10.1063/1.4974306.
7. Stefan Chmiela, Huziel E. Sauceda, Klaus-Robert Müller, and Alexandre Tkatchenko. Towards exact molecular dynamics simulations with machine-learned force fields. *Nat. Commun.*, 9:3887, 2018. doi: 10.1038/s41467-018-06169-2.
8. Jiang Wang, Simon Olsson, Christoph Wehmeyer, Adrià Pérez, Nicholas E. Charron, Gianni de Fabritiis, Frank Noé, and Cecilia Clementi. Machine learning of coarse-grained molecular dynamics force fields. *ACS Cent. Sci.*, 5(5):755–767, 2019. doi: 10.1021/acscentsci.8b00913.
9. Florian Häse, Ignacio Fernández Galván, Alán Aspuru-Guzik, Roland Lindh, and Morgane Vacher. How machine learning can assist the interpretation of ab initio molecular dynamics simulations and conceptual understanding of chemistry. *Chem. Sci.*, 10:2298–2307, 2019. doi: 10.1039/C8SC04516J.
10. Ian Jolliffe. *Principal Component Analysis*. Wiley Online Library, 2002.
11. Ronald R Coifman and Stéphane Lafon. Diffusion maps. *Applied and Computational Harmonic Analysis*, 21(1):5–30, 2006.
12. Ronald R Coifman, Ioannis G Kevrekidis, Stéphane Lafon, Mauro Maggioni, and Boaz Nadler. Diffusion maps, reduction coordinates, and low dimensional representation of stochastic systems. *Multiscale Model. Simul.*, 7(2):842–864, 2008.
13. Wei Chen and Andrew L Ferguson. Molecular enhanced sampling with autoencoders: On-the-fly collective variable discovery and accelerated free energy landscape exploration. *J. Comput. Chem.*, 39(25):2079–2102, 2018.
14. Wei Chen, Aik Rui Tan, and Andrew L Ferguson. Collective variable discovery and enhanced sampling using autoencoders: Innovations in network architecture and error function design. *J. Chem. Phys.*, 149(7):072312, 2018.
15. João Marcelo Lamim Ribeiro, Pablo Bravo, Yihang Wang, and Pratyush Tiwary. Reweighted autoencoded variational bayes for enhanced sampling (RAVE). *J. Chem. Phys.*, 149(7):072301, 2018.
16. Christoph Wehmeyer and Frank Noé. Time-lagged autoencoders: Deep learning of slow collective variables for molecular kinetics. *J. Chem. Phys.*, 148(24):241703, June 2018.
17. Ao Ma and Aaron R Dinner. Automatic method for identifying reaction coordinates in complex systems. *J. Phys. Chem. B*, 109(14):6769–6779, 2005.
18. Alexander V. Shapeev. Moment tensor potentials: A class of systematically improvable interatomic potentials. *Multiscale Modeling & Simulation*, 14(3):1153–1173, jan 2016. doi: 10.1137/15m1054183.
19. H. Chen, J. Lu, and C. Ortner. Thermodynamic limit of crystal defects with finite temperature tight binding. *Arch. Ration. Mech. Anal.*, 230:701–733, 2018.
20. Alessandro Lunghi and Stefano Sanvito. A unified picture of the covalent bond within quantum-accurate force fields: From organic molecules to metallic complexes' reactivity. *Science Advances*, 5(5):eaaw2210, 2019.
21. N. Artrith and J. Behler. High-dimensional neural network potentials for metal surfaces: A prototype study for copper. *Phys. Rev. B*, 85:045439, 2012.
22. Evgeny V. Podryabinkin, Evgeny V. Tikhonov, Alexander V. Shapeev, and Artem R.

Oganov. Accelerating crystal structure prediction by machine-learning interatomic potentials with active learning. *Phys. Rev. B*, 99:064114, 2019.

23. Konstantin Gubaev, Evgeny V. Podryabinkin, Gus L.W. Hart, and Alexander V. Shapeev. Accelerating high-throughput searches for new alloys with active learning of interatomic potentials. *Computational Materials Science*, 156:148–156, 2019.

24. Tran Doan Huan, Rohit Batra, James Chapman, Chiho Kim, Anand Chandrasekaran, and Rampi Ramprasad. Iterative-learning strategy for the development of application-specific atomistic force fields. *The Journal of Physical Chemistry C*, 123(34):20715–20722, 2019.

25. Ryosuke Jinnouchi, Ferenc Karsai, and Georg Kresse. On-the-fly machine learning force field generation: Application to melting points. *Phys. Rev. B*, 100:014105, 2019.

26. Volker L. Deringer, Davide M. Proserpio, Gábor Csányi, and Chris J. Pickard. Data-driven learning and prediction of inorganic crystal structures. *Faraday Discuss.*, 211:45–59, 2018.

27. M. Eickenberg, G. Exarchakis, M. Hirn, S. Mallat, and L. Thiry. Solid harmonic wavelet scattering for predictions of molecule properties. *J. Chem. Phys.*, 148(24):241732, 2018.

28. G. Ferré, T. Haut, and K. Barros. Learning molecular energies using localized graph kernels. *The Journal of Chemical Physics*, 146(11):114107, 2017.

29. David Weininger. SMILES, a chemical language and information system. 1. Introduction to methodology and encoding rules. *Journal of Chemical Information and Modeling*, 28(1):31–36, 1988. doi: 10.1021/ci00057a005.

30. Paul J. Steinhardt, David R. Nelson, and Marco Ronchetti. Bond-orientational order in liquids and glasses. *Physical Review B*, 28(2):784–805, jul 1983. doi: 10.1103/physrevb.28.784.

31. Fabio Pietrucci and Alessandro Laio. A collective variable for the efficient exploration of protein beta-sheet structures: Application to SH3 and GB1. *Journal of Chemical Theory and Computation*, 5(9):2197–2201, aug 2009. doi: 10.1021/ct900202f.

32. Albert P. Bartók, Mike C. Payne, Risi Kondor, and Gábor Csányi. Gaussian approximation potentials: The accuracy of quantum mechanics, without the electrons. *Phys. Rev. Lett.*, 104(13):136403, apr 2010. doi: 10.1103/PhysRevLett.104.136403.

33. Jörg Behler and Michele Parrinello. Generalized neural-network representation of high-dimensional potential-energy surfaces. *Phys. Rev. Lett.*, 98(14):146401, apr 2007. doi: 10.1103/PhysRevLett.98.146401.

34. Felix A. Faber, Anders S. Christensen, Bing Huang, and O. Anatole von Lilienfeld. Alchemical and structural distribution based representation for universal quantum machine learning. *The Journal of Chemical Physics*, 148(24):241717, jun 2018. doi: 10.1063/1.5020710.

35. Michael J. Willatt, Félix Musil, and Michele Ceriotti. Atom-density representations for machine learning. *The Journal of Chemical Physics*, 150(15):154110, apr 2019. doi: 10.1063/1.5090481.

36. Albert P. Bartók, Sandip De, Carl Poelking, Noam Bernstein, James R. Kermode, Gábor Csányi, and Michele Ceriotti. Machine learning unifies the modeling of materials and molecules. *Science Advances*, 3(12):e1701816, dec 2017. doi: 10.1126/sciadv.1701816.

37. Yunxing Zuo, Chi Chen, Xiangguo Li, Zhi Deng, Yiming Chen, Jörg Behler, Gábor Csányi, Alexander V. Shapeev, Aidan P. Thompson, Mitchell A. Wood, and Shyue Ping Ong. Performance and cost assessment of machine learning interatomic potentials. *The Journal of Physical Chemistry A*, 124(4):731–745, 2020.

38. Jörg Behler. Perspective: Machine learning potentials for atomistic simulations. *J. Chem. Phys.*, 145(17):170901, 2016. doi: 10.1063/1.4966192.

39. Bastiaan J. Braams and Joel M. Bowman. Permutationally invariant potential energy surfaces in high dimensionality. *Int. Rev. Phys. Chem.*, 28:577–606, 2009. doi: 10.1080/01442350903234923.

40. Stefan Chmiela, Alexandre Tkatchenko, Huziel E. Sauceda, Igor Poltavsky, Kristof T. Schütt, and Klaus-Robert Müller. Machine learning of accurate energy-conserving molecular force fields. *Sci. Adv.*, 3(5):e1603015, 2017. doi: 10.1126/sciadv.1603015.

41. Kristof T Schütt, Farhad Arbabzadah, Stefan Chmiela, Klaus R Müller, and Alexandre Tkatchenko. Quantum-chemical insights from deep tensor neural networks. *Nat. Commun.*, 8:13890, 2017.

42. K. T. Schütt, H. E. Sauceda, P.-J. Kindermans, A. Tkatchenko, and K.-R. Müller. SchNet - a deep learning architecture for molecules and materials. *J. Chem. Phys.*, 148(24):241722, 2018.

43. Chen Qu and Joel M. Bowman. A fragmented, permutationally invariant polynomial approach for potential energy surfaces of large molecules: Application to n-methyl acetamide. *J. Chem. Phys.*, 150(14):141101, 2019. doi: 10.1063/1.5092794.

44. Volker L. Deringer and Gábor Csányi. Machine learning based interatomic potential for amorphous carbon. *Phys. Rev. B*, 95:094203, 2017.

45. Alberto Ambrosetti, Nicola Ferri, Robert A. DiStasio, and Alexandre Tkatchenko. Wavelike charge density fluctuations and van der Waals interactions at the nanoscale. *Science*, 351:1171–1176, 2016. doi: 10.1126/science.aae0509.

46. Jan Hermann, Robert A. DiStasio, and Alexandre Tkatchenko. First-principles models for van der Waals interactions in molecules and materials: Concepts, theory, and applications. *Chem. Rev.*, 117:4714–4758, 2017. doi: 10.1021/acs.chemrev.6b00446.

47. Tristan Bereau, Robert A. DiStasio Jr, Alexandre Tkatchenko, and O. Anatole Von Lilienfeld. Non-covalent interactions across organic and biological subsets of chemical space: Physics-based potentials parametrized from machine learning. *J. Chem. Phys.*, 148(24):241706, 2018.

48. Andrea Grisafi and Michele Ceriotti. Incorporating long-range physics in atomic-scale machine learning. *The Journal of Chemical Physics*, 151(20):204105, nov 2019. doi: 10.1063/1.5128375.

49. Aldo Glielmo, Claudio Zeni, and Alessandro De Vita. Efficient nonparametric n-body force fields from machine learning. *Physical Review B*, 97(18), may 2018. doi: 10.1103/physrevb.97.184307.

50. Max Veit, Sandeep Kumar Jain, Satyanarayana Bonakala, Indranil Rudra, Detlef Hohl, and Gábor Csányi. Equation of state of fluid methane from first principles with machine learning potentials. *Journal of Chemical Theory and Computation*, 15(4):2574–2586, feb 2019. doi: 10.1021/acs.jctc.8b01242.

51. James F. Ziegler and Jochen P. Biersack. The stopping and range of ions in matter. In D. Allan Bromley, editor, *Treatise on Heavy-Ion Science: Volume 6: Astrophysics, Chemistry, and Condensed Matter*, pages 93–129. Springer US, Boston, MA, 1985.

52. Tobias Lemke and Christine Peter. Neural network based prediction of conformational free energies - A new route toward coarse-grained simulation models. *J. Chem. Theory Comput.*, 13(12):6213–6221, 2017.

53. S. Hunkler, T. Lemke, C. Peter, and O. Kukharenko. Back-mapping based sampling: Coarse grained free energy landscapes as a guideline for atomistic exploration. *J. Chem. Phys.*, 151(15):154102, 2019.

54. Christine Peter and Kurt Kremer. Multiscale simulation of soft matter systems - from the atomistic to the coarse-grained level and back. *Soft Matter*, 5(22):4357–4366, 2009.

55. Joseph F Rudzinski and W G Noid. Coarse-graining entropy, forces, and structures. *J. Chem. Phys.*, 135(21):214101, 2011.

56. W G Noid. Perspective: Coarse-grained models for biomolecular systems. *J. Chem. Phys.*, 139(9):090901, 2013.

57. Raffaello Potestio, Christine Peter, and Kurt Kremer. Computer simulations of soft matter: Linking the scales. *Entropy*, 16(8):4199–4245, 2014.

58. M S Shell. Coarse-graining with the relative entropy. *Adv. Chem. Phys.*, pages 395–441, 2016.

59. S. T. John and Gabor Csányi. Many-body coarse-grained interactions using Gaussian approximation potentials. *J. Phys. Chem. B*, 121(48):10934–10949, 2017.

60. Linfeng Zhang, Jiequn Han, Han Wang, Roberto Car, and Weinan E. DeePCG: Constructing coarse-grained models via deep neural networks. *J. Chem. Phys.*, 149(3):034101, 2018.

61. W G Noid, Jhih-Wei Chu, Gary S Ayton, Vinod Krishna, Sergei Izvekov, Gregory A Voth, Avisek Das, and Hans C Andersen. The multiscale coarse-graining method. I. A rigorous bridge between atomistic and coarse-grained models. *J. Chem. Phys.*, 128(24):244114, 2008.

62. L. Garrido and A. Juste. On the determination of probability density functions by using neural networks. *Comput. Phys. Commun.*, 115:25–31, 1998.

63. Andrew L Ferguson, Athanassios Z Panagiotopoulos, Pablo G Debenedetti, and Ioannis G Kevrekidis. Systematic determination of order parameters for chain dynamics using diffusion maps. *P. Natl. Acad. Sci. USA*, 107(31):13597–13602, 2010.

64. Andrew L Ferguson, Athanassios Z Panagiotopoulos, Ioannis G Kevrekidis, and Pablo G Debenedetti. Nonlinear dimensionality reduction in molecular simulation: The diffusion map approach. *Chem. Phys. Lett.*, 509(1-3):1–11, 2011.

65. J Wang and AL Ferguson. Nonlinear machine learning in simulations of soft and biological materials. *Mol. Simul.*, 44(13-14):1090–1107, 2018.

66. Fabio Pietrucci. Strategies for the exploration of free energy landscapes: unity in diversity and challenges ahead. *Reviews in Physics*, 2:32–45, 2017. doi: 10.1016/j.revip.2017.05.001.

67. Toshiko Ichiye and Martin Karplus. Collective motions in proteins: a covariance analysis of atomic fluctuations in molecular dynamics and normal mode simulations. *Proteins: Struct., Funct., Bioinf.*, 11(3):205–217, 1991.

68. Angel E García. Large-amplitude nonlinear motions in proteins. *Phys. Rev. Lett.*, 68(17):2696, 1992.

69. A. Amadei, A.B.M Linssen, and H.J.C. Berendsen. Essential dynamics of proteins. *Proteins: Structure, Function, and Bioinformatics*, 17(4):412–425, 1993.

70. Andrew L Ferguson. Machine learning and data science in soft materials engineering. *J. Phys. Condens. Matter*, 30(4):043002, 2017.

71. B. Schölkopf. The kernel trick for distances. In *Proceedings of the 13th International Conference on Neural Information Processing Systems*, NIPS'00, page 283–289, Cambridge, MA, USA, 2000. MIT Press.

72. F. Sittel and G. Stock. Perspective: Identification of collective variables and metastable states of protein dynamics. *J. Chem. Phys.*, 149(15):150901, 2018.

73. Zofia Trstanova, Ben Leimkuhler, and Tony Lelièvre. Local and global perspectives on diffusion maps in the analysis of molecular systems. *Proc. R. Soc. A*, 476(2233):20190036, 2020.

74. Hannes Jónsson, G. Mills, and K.W. Jacobsen. Nudged elastic band method for finding minimum energy paths of transitions. In B.J. Berne, G. Ciccotti, and D.F. Coker, editors, *Classical and Quantum Dynamics in Condensed Phase Simulations*, pages 385–404. World Scientific, 1998.

75. Weinan E, Ren Weiqing, and Eric Vanden-Eijnden. String method for the study of rare events. *Phys. Rev. B*, 66(5):52301, 2002.

76. Weinan E and Eric Vanden-Eijnden. Transition-path theory and path-finding algorithms for the study of rare events. *Annu. Rev. Phys. Chem.*, 61:391–420, 2010.

77. T. Lelièvre and G. Stoltz. Partial differential equations and stochastic methods in molecular dynamics. *Acta Numerica*, 25:681–880, 2016.

78. Peter G Bolhuis, David Chandler, Christoph Dellago, and Phillip L Geissler. Transition path sampling: Throwing ropes over rough mountain passes, in the dark. *Annu. Rev. Phys. Chem.*, 53(1):291–318, 2002.

79. Polina V Banushkina and Sergei V Krivov. Optimal reaction coordinates. *WIREs: Comput. Mol. Sci.*, 6(6):748–763, 2016.

80. Hendrik Jung, Roberto Covino, and Gerhard Hummer. Artificial intelligence assists discovery of reaction coordinates and mechanisms from molecular dynamics simulations. *arXiv preprint*, 1901.04595, 2019.

81. Michel Loève. *Probability Theory: Foundations, Random Sequences*. Van Nostrand, 1955.

82. Lawrence Sirovich. Turbulence and the dynamics of coherent structures. I. Coherent structures. *Q. Appl. Math.*, 45(3):561–571, 1987.

83. Lawrence Sirovich. Turbulence and the dynamics of coherent structures. II. Symmetries and transformations. *Q. Appl. Math.*, 45(3):573–582, 1987.

84. HM Park and DH Cho. The use of the Karhunen-Loeve decomposition for the modeling of distributed parameter systems. *Chem. Eng. Sci.*, 51(1):81–98, 1996.

85. Anindya Chatterjee. An introduction to the proper orthogonal decomposition. *Current Science*, 78:808–817, 2000.

86. YC Liang, HP Lee, SP Lim, WZ Lin, KH Lee, and CG Wu. Proper orthogonal decomposition and its applications—Part I: Theory. *J. Sound Vib.*, 252(3):527–544, 2002.

87. Bernhard Schölkopf, Alexander Smola, and Klaus-Robert Müller. Kernel principal component analysis. In Wulfram Gerstner, Alain Germond, Martin Hasler, and Jean-Daniel

Nicoud, editors, *Artificial Neural Networks — ICANN'97: 7th International Conference Lausanne, Switzerland, October 8–10, 1997 Proceedings*, pages 583–588, Berlin Heidelberg, October 1997. Springer.

88. Mark A Kramer. Nonlinear principal component analysis using autoassociative neural networks. *AIChE Journal*, 37(2):233–243, 1991.

89. Phuong H. Nguyen. Complexity of free energy landscapes of peptides revealed by nonlinear principal component analysis. *Proteins: Structure, Function, and Bioinformatics*, 65 (4):898–913, 2006. ISSN 1097-0134. doi: 10.1002/prot.21185.

90. Matthias Scholz, Martin Fraunholz, and Joachim Selbig. Nonlinear principal component analysis: neural network models and applications. In Alexander N. Gorban, Balázs Kégl, Donald C. Wunsch, and Andrei Y. Zinovyev, editors, *Principal Manifolds for Data Visualization and Dimension Reduction*, pages 44–67. Springer, Berlin Heidelberg, 2008.

91. Pierre Comon. Independent component analysis, a new concept? *Signal Processing*, 36 (3):287–314, 1994.

92. Ingwer Borg and Patrick JF Groenen. *Modern Multidimensional Scaling: Theory and Applications*. Springer Science & Business Media, 2005.

93. Michele Ceriotti, Gareth A. Tribello, and Michele Parrinello. Simplifying the representation of complex free-energy landscapes using sketch-map. *Proceedings of the National Academy of Sciences*, 108(32):13023–13028, jul 2011. doi: 10.1073/pnas.1108486108.

94. Sam T Roweis and Lawrence K Saul. Nonlinear dimensionality reduction by locally linear embedding. *Science*, 290(5500):2323–2326, 2000.

95. Zhenyue Zhang and Jing Wang. MLLE: Modified locally linear embedding using multiple weights. In Bernhard Schölkopf, John Platt, and Thomas Hofmann, editors, *Advances in Neural Information Processing Systems 19: Proceedings of the 2006 Conference*, pages 1593–1600, Cambridge, December 2007. MIT Press.

96. Payel Das, Mark Moll, Hernan Stamati, Lydia E Kavraki, and Cecilia Clementi. Low-dimensional, free-energy landscapes of protein-folding reactions by nonlinear dimensionality reduction. *P. Natl. Acad. Sci. USA*, 103(26):9885–9890, 2006.

97. Joshua B Tenenbaum, Vin De Silva, and John C Langford. A global geometric framework for nonlinear dimensionality reduction. *Science*, 290(5500):2319–2323, 2000.

98. Vin D. Silva and Joshua B. Tenenbaum. Global versus local methods in nonlinear dimensionality reduction. In S. Thrun and K. Obermayer, editors, *Advances in Neural Information Processing Systems 15*, pages 705–712. MIT Press, Cambridge, MA, 2002.

99. Jianzhong Wang. Local tangent space alignment. In *Geometric Structure of High-Dimensional Data and Dimensionality Reduction*, pages 221–234. Springer, Berlin Heidelberg, 2012.

100. Kilian Q Weinberger and Lawrence K Saul. Unsupervised learning of image manifolds by semidefinite programming. *Int. J. Comput. Vision*, 70(1):77–90, 2006.

101. Mikhail Belkin and Partha Niyogi. Laplacian eigenmaps for dimensionality reduction and data representation. *Neural Computation*, 15(6):1373–1396, 2003.

102. David L Donoho and Carrie Grimes. Hessian eigenmaps: Locally linear embedding techniques for high-dimensional data. *P. Natl. Acad. Sci. USA*, 100(10):5591–5596, 2003.

103. Ronald R Coifman, Stephane Lafon, Ann B Lee, Mauro Maggioni, Boaz Nadler, Frederick Warner, and Steven W Zucker. Geometric diffusions as a tool for harmonic analysis and structure definition of data: Diffusion maps. *P. Natl. Acad. Sci. USA*, 102(21):7426–7431, 2005.

104. Zahra Shamsi, Kevin J Cheng, and Diwakar Shukla. REinforcement learning based Adaptive samPling: REAPing Rewards by Exploring Protein Conformational Landscapes. *J. Phys. Chem. B*, 122:8386–8395, 2018.

105. Eliodoro Chiavazzo, Roberto Covino, Ronald R Coifman, C William Gear, Anastasia S Georgiou, Gerhard Hummer, and Ioannis G Kevrekidis. Intrinsic map dynamics exploration for uncharted effective free-energy landscapes. *Proc. Natl. Acad. Sci. USA*, 114 (28):E5494–E5503, 2017.

106. Andrew L Ferguson, Athanassios Z Panagiotopoulos, Pablo G Debenedetti, and Ioannis G Kevrekidis. Integrating diffusion maps with umbrella sampling: Application to alanine dipeptide. *J. Chem. Phys.*, 134(13):135103, 2011.

107. G. A. Tribello, M. Ceriotti, and M. Parrinello. Using sketch-map coordinates to analyze and bias molecular dynamics simulations. *Proceedings of the National Academy of Sciences*, 109(14):5196–5201, mar 2012. doi: 10.1073/pnas.1201152109.

108. Cameron F Abrams and Eric Vanden-Eijnden. On-the-fly free energy parameterization via temperature accelerated molecular dynamics. *Chem. Phys. Lett.*, 547:114–119, 2012.

109. Behrooz Hashemian, Daniel Millán, and Marino Arroyo. Modeling and enhanced sampling of molecular systems with smooth and nonlinear data-driven collective variables. *J. Chem. Phys.*, 139(21):214101, 2013.

110. Chun-Guang Li, Jun Guo, Guang Chen, Xiang-Fei Nie, and Zhen Yang. A version of Isomap with explicit mapping. In *2006 International Conference on Machine Learning and Cybernetics*, pages 3201–3206. IEEE, 2006.

111. Vojtěch Spiwok and Blanka Králová. Metadynamics in the conformational space nonlinearly dimensionally reduced by Isomap. *J. Chem. Phys.*, 135(22):224504, 2011.

112. Davide Branduardi, Francesco Luigi Gervasio, and Michele Parrinello. From A to B in free energy space. *J. Chem. Phys.*, 126(5):054103, 2007.

113. Jordane Preto and Cecilia Clementi. Fast recovery of free energy landscapes via diffusion-map-directed molecular dynamics. *Phys. Chem. Chem. Phys.*, 16(36):19181–19191, 2014.

114. Wenwei Zheng, Mary A Rohrdanz, and Cecilia Clementi. Rapid exploration of configuration space with diffusion-map-directed molecular dynamics. *J. Phys. Chem. B*, 117(42):12769–12776, 2013.

115. Frank Noé and Feliks Nüske. A variational approach to modeling slow processes in stochastic dynamical systems. *Multiscale Model. Sim.*, 11(2):635–655, 2013.

116. Igor Mezić. Spectral properties of dynamical systems, model reduction and decompositions. *Nonlinear Dynamics*, 41(1-3):309–325, 2005.

117. Matthew O Williams, Ioannis G Kevrekidis, and Clarence W Rowley. A data-driven approximation of the Koopman operator: Extending dynamic mode decomposition. *J. Nonlinear Sci.*, 25(6):1307–1346, 2015.

118. Hao Wu, Feliks Nüske, Fabian Paul, Stefan Klus, Péter Koltai, and Frank Noé. Variational Koopman models: slow collective variables and molecular kinetics from short off-equilibrium simulations. *J. Chem. Phys.*, 146(15):154104, 2017.

119. Guillermo Pérez-Hernández, Fabian Paul, Toni Giorgino, Gianni De Fabritiis, and Frank Noé. Identification of slow molecular order parameters for Markov model construction. *J. Chem. Phys.*, 139(1):015102, 2013.

120. Feliks Nüske, Bettina G Keller, Guillermo Pérez-Hernández, Antonia S J S Mey, and Frank Noé. Variational approach to molecular kinetics. *J. Chem. Theory Comput.*, 10(4):1739–1752, 2014.

121. Frank Noé and Cecilia Clementi. Kinetic distance and kinetic maps from molecular dynamics simulation. *J. Chem. Theory Comput*, 11(10):5002–5011, 2015.

122. Frank Noé, Ralf Banisch, and Cecilia Clementi. Commute maps: Separating slowly mixing molecular configurations for kinetic modeling. *J. Chem. Theory Comput.*, 12(11):5620–5630, 2016.

123. Guillermo Pérez-Hernández and Frank Noé. Hierarchical time-lagged independent component analysis: Computing slow modes and reaction coordinates for large molecular systems. *J. Chem. Theory Comput.*, 12(12):6118–6129, 2016.

124. Christian R. Schwantes and Vijay S. Pande. Improvements in Markov state model construction reveal many non-native interactions in the folding of NTL9. *Journal of Chemical Theory and Computation*, 9(4):2000–2009, March 2013.

125. Stefan Klus, Feliks Nüske, Péter Koltai, Hao Wu, Ioannis Kevrekidis, Christof Schütte, and Frank Noé. Data-driven model reduction and transfer operator approximation. *J. Nonlinear Sci.*, 28(3):985–1010, 2018.

126. Vijay S Pande, Kyle Beauchamp, and Gregory R Bowman. Everything you wanted to know about Markov state models but were afraid to ask. *Methods*, 52(1):99–105, 2010.

127. Frank Noé and Edina Rosta. Markov models of molecular kinetics. *The Journal of Chemical Physics*, 151(19):190401, 2019.

128. Christof Schütte, Alexander Fischer, Wilhelm Huisinga, and Peter Deuflhard. A direct approach to conformational dynamics based on hybrid Monte Carlo. *J. Comput. Phys.*, 151(1):146–168, 1999.

129. P. Deuflhard and M. Weber. Robust Perron cluster analysis in conformation dynamics. *Linear algebra and its applications*, 398:161–184, 2005.

130. S. Röblitz and M. Weber. Fuzzy spectral clustering by PCCA+: Application to Markov state models and data classification. *Advances in Data Analysis and Classification*, 7(2): 147–179, 2013.

131. Christian R Schwantes and Vijay S Pande. Modeling molecular kinetics with tICA and the kernel trick. *J. Chem. Theory Comput.*, 11(2):600–608, 2015.

132. Galen Andrew, Raman Arora, Jeff Bilmes, and Karen Livescu. Deep canonical correlation analysis. In Sanjoy Dasgupta and David McAllester, editors, *Proceedings of the 30th International Conference on Machine Learning*, volume 28 of *Proceedings of Machine Learning Research*, pages 1247–1255, Atlanta, Georgia, USA, 17–19 Jun 2013. PMLR.

133. Andreas Mardt, Luca Pasquali, Hao Wu, and Frank Noé. VAMPnets for deep learning of molecular kinetics. *Nature Communications*, 9:5, 2018.

134. Wei Chen, Hythem Sidky, and Andrew L Ferguson. Nonlinear discovery of slow molecular modes using state-free reversible VAMPnets. *J. Chem. Phys.*, 150:214114, 2019.

135. Carlos X Hernández, Hannah K Wayment-Steele, Mohammad M Sultan, Brooke E Husic, and Vijay S Pande. Variational encoding of complex dynamics. *Phys. Rev. E*, 97(6): 062412, 2018.

136. Hannah K Wayment-Steele and Vijay S Pande. Note: Variational encoding of protein dynamics benefits from maximizing latent autocorrelation. *J. Chem. Phys.*, 149:216101, 2018.

137. Jannes Quer, Luca Donati, BETTINA G Keller, and Marcus Weber. An automatic adaptive importance sampling algorithm for molecular dynamics in reaction coordinates. *SIAM J. Sci. Comput.*, 40(2):A653–A670, 2018.

138. Luca Donati and Bettina G Keller. Girsanov reweighting for metadynamics simulations. *J. Chem. Phys.*, 149(7):072335, 2018.

139. Lorenzo Donati, Carsten Hartmann, and Bettina G Keller. Girsanov reweighting for path ensembles and Markov state models. *J. Chem. Phys.*, 146(24):244112, 2017.

140. Boaz Nadler, Stephane Lafon, Ronald Coifman, and Ioannis G Kevrekidis. Diffusion Maps - a Probabilistic Interpretation for Spectral Embedding and Clustering Algorithms. In Alexander N Gorban, Balázs Kégl, Donald C Wunsch, and Andrei Y Zinovyev, editors, *Principal Manifolds for Data Visualization and Dimension Reduction*, pages 238–260, Berlin, Heidelberg, 2008. Springer. ISBN 978-3-540-73750-6.

141. Mary A. Rohrdanz, Wenwei Zheng, Mauro Maggioni, and Cecilia Clementi. Determination of reaction coordinates via locally scaled diffusion map. *J. Chem. Phys.*, 134(12):124116, 2011.

142. Tyrus Berry and John Harlim. Variable bandwidth diffusion kernels. *Appl. Comput. Harmon. Anal.*, 40(1):68–96, 2016. ISSN 1096603X. doi: 10.1016/j.acha.2015.01.001.

143. Ralf Banisch, Zofia Trstanova, Andreas Bittracher, Stefan Klus, and Péter Koltai. Diffusion maps tailored to arbitrary non-degenerate Itô processes. *Appl. Comput. Harmon. Anal.*, 48:242–265, 2018. ISSN 1063-5203.

144. Pierre Collet, Servet Martinez, and Jaime San Martin. *Quasi-Stationary Distributions: Markov Chains, Diffusions and Dynamical Systems*. Springer Science & Business Media, 2012. ISBN 9783642331305.

145. Rongjie Lai and Jianfeng Lu. Point cloud discretization of Fokker–Planck operators for committor functions. *Multiscale Model. Simul.*, 16(2):710–726, 2018.

146. Jan-Hendrik Prinz, Martin Held, Jeremy C Smith, and Frank Noé. Efficient computation, sensitivity, and error analysis of committor probabilities for complex dynamical processes. *Multiscale Model. Simul.*, 9(2):545–567, 2011.

147. E. H. Thiede, D. Giannakis, A. R. Dinner, and J. Weare. Galerkin approximation of dynamical quantities using trajectory data. *J Chem Phys*, 150(24):244111, Jun 2019.

148. Yuehaw Khoo, Jianfeng Lu, and Lexing Ying. Solving for high dimensional committor functions using artificial neural networks. *Research in the Mathematical Sciences*, 6:1, 2019.

149. Lutz Molgedey and Heinz Georg Schuster. Separation of a mixture of independent signals using time delayed correlations. *Phys. Rev. Lett.*, 72(23):3634–3637, 1994.

150. Hiroshi Takano and Seiji Miyashita. Relaxation modes in random spin systems. *J. Phys. Soc. Jpn.*, 64(10):3688–3698, 1995.

151. Hidetomo Hirao, Sachiko Koseki, and Hiroshi Takano. Molecular dynamics study of relax-

ation modes of a single polymer chain. *J. Phys. Soc. Jpn*, 66(11):3399–3405, 1997.

152. William C Swope, Jed W Pitera, and Frank Suits. Describing protein folding kinetics by molecular dynamics simulations. 1. Theory. *J. Phys. Chem. B*, 108(21):6571–6581, 2004.

153. Jan-Hendrik Prinz, Hao Wu, Marco Sarich, Bettina Keller, Martin Senne, Martin Held, John D Chodera, Christof Schütte, and Frank Noé. Markov models of molecular kinetics: Generation and validation. *J. Chem. Phys.*, 134(17):174105, 2011.

154. Dimitrios Giannakis, Joanna Slawinska, and Zhizhen Zhao. Spatiotemporal feature extraction with data-driven Koopman operators. In Dmitry Storcheus, Afshin Rostamizadeh, and Sanjiv Kumar, editors, *Proceedings of the 1st International Workshop on Feature Extraction: Modern Questions and Challenges at NIPS 2015*, volume 44 of *Proceedings of Machine Learning Research*, pages 103–115, Montreal, Canada, 11 Dec 2015. PMLR.

155. Floris Takens. Detecting strange attractors in turbulence. In *Dynamical Systems and Turbulence*, volume 898 of *Lecture Notes in Mathematics*, pages 366–381. Springer, 1981.

156. Dirk Aeyels. Generic observability of differentiable systems. *SIAM J. Control Optim.*, 19 (5):595–603, 1981.

157. Robert Zwanzig. *Nonequilibrium Statistical Mechanics*. Oxford University Press, 2001.

158. J Łuczka. Non-markovian stochastic processes: Colored noise. *Chaos*, 15(2):026107, 2005.

159. Carlo Camilloni and Fabio Pietrucci. Advanced simulation techniques for the thermodynamic and kinetic characterization of biological systems. *Adv. Phys.:X.*, 3(1):1477531, 2018.

160. Peter Hänggi, Peter Talkner, and Michal Borkovec. Reaction-rate theory: Fifty years after Kramers. *Rev. Mod. Phys.*, 62:251–341, 1990. doi: 10.1103/RevModPhys.62.251.

161. John E Straub, Michal Borkovec, and Bruce J Berne. Calculation of dynamic friction on intramolecular degrees of freedom. *J. Phys. Chem.*, 91(19):4995–4998, 1987.

162. Gerhard Hummer and Ioannis G Kevrekidis. Coarse molecular dynamics of a peptide fragment: Free energy, kinetics, and long-time dynamics computations. *J. Chem. Phys.*, 118(23):10762–10773, 2003.

163. Oliver F Lange and Helmut Grubmüller. Collective Langevin dynamics of conformational motions in proteins. *J. Chem. Phys.*, 124(21):214903, 2006.

164. Gerhard Hummer. Position-dependent diffusion coefficients and free energies from Bayesian analysis of equilibrium and replica molecular dynamics simulations. *New J. Phys.*, 7(1):34, 2005.

165. Illia Horenko, Carsten Hartmann, Christof Schütte, and Frank Noé. Data-based parameter estimation of generalized multidimensional Langevin processes. *Phys. Rev. E*, 76(1): 016706, 2007.

166. Cristian Micheletti, Giovanni Bussi, and Alessandro Laio. Optimal Langevin modeling of out-of-equilibrium molecular dynamics simulations. *J. Chem. Phys.*, 129(7):074105, 2008.

167. Eric Darve, Jose Solomon, and Amirali Kia. Computing generalized Langevin equations and generalized Fokker–Planck equations. *P. Natl. Acad. Sci. USA*, 106(27):10884–10889, 2009.

168. Frédéric Legoll and Tony Lelièvre. Effective dynamics using conditional expectations. *Nonlinearity*, 23(9):2131, 2010.

169. Norbert Schaudinnus, Björn Bastian, Rainer Hegger, and Gerhard Stock. Multidimensional Langevin modeling of nonoverdamped dynamics. *Phys. Rev. Lett.*, 115(5):050602, 2015.

170. Roberto Meloni, Carlo Camilloni, and Guido Tiana. Properties of low-dimensional collective variables in the molecular dynamics of biopolymers. *Phys. Rev. E*, 94(5):052406, 2016.

171. Dominika Lesnicki, Rodolphe Vuilleumier, Antoine Carof, and Benjamin Rotenberg. Molecular hydrodynamics from memory kernels. *Phys. Rev. Lett.*, 116(14):147804, 2016.

172. Jan O Daldrop, Julian Kappler, Florian N Brünig, and Roland R Netz. Butane dihedral angle dynamics in water is dominated by internal friction. *P. Natl. Acad. Sci. USA*, 115(20): 5169–5174, 2018.

173. Qi Zhang, Jasna Brujić, and Eric Vanden-Eijnden. Reconstructing free energy profiles from nonequilibrium relaxation trajectories. *J. Stat. Phys.*, 144(2):344–366, 2011.

174. Andrea Pérez-Villa and Fabio Pietrucci. Free energy, friction, and mass profiles from short molecular dynamics trajectories. *arXiv preprint*, 1810.00713, 2018.

175. Amit Samanta, Ming Chen, Tang-Qing Yu, Mark Tuckerman, and Weinan E. Sampling saddle points on a free energy surface. *J. Chem. Phys.*, 140(16):164109, 2014.

176. Fabio Pietrucci and Antonino Marco Saitta. Formamide reaction network in gas phase and solution via a unified theoretical approach: Toward a reconciliation of different prebiotic scenarios. *P. Natl. Acad. Sci. USA*, 112(49):15030–15035, 2015.

177. Silvio Pipolo, Mathieu Salanne, Guillaume Ferlat, Stefan Klotz, A Marco Saitta, and Fabio Pietrucci. Navigating at will on the water phase diagram. *Phys. Rev. Lett.*, 119(24):245701, 2017.

178. Mohammed AlQuraishi. End-to-end differentiable learning of protein structure. *Cell Systems*, 8(4):292–301.e3, 2019. ISSN 2405-4712. doi: https://doi.org/10.1016/j.cels.2019.03.006.

179. David E. Shaw, Paul Maragakis, Kresten Lindorff-Larsen, Stefano Piana, Ron O. Dror, Michael P. Eastwood, Joseph A. Bank, John M. Jumper, John K. Salmon, Yibing Shan, and Willy Wriggers. Atomic-level characterization of the structural dynamics of proteins. *Science*, 330(6002):341–346, 2010. ISSN 0036-8075. doi: 10.1126/science.1187409.

180. Sergei V. Krivov and Martin Karplus. Diffusive reaction dynamics on invariant free energy profiles. *P. Natl. Acad. Sci. USA*, 105(37):13841–13846, 2008. ISSN 0027-8424. doi: 10.1073/pnas.0800224105.

181. Mohammad M. Sultan, Hannah K. Wayment-Steele, and Vijay S. Pande. Transferable neural networks for enhanced sampling of protein dynamics. *Journal of Chemical Theory and Computation*, 14(4):1887–1894, 2018. doi: 10.1021/acs.jctc.8b00025.

182. Simon Brandt, Florian Sittel, Matthias Ernst, and Gerhard Stock. Machine learning of biomolecular reaction coordinates. *J. Phys. Chem. Lett.*, 9(9):2144–2150, apr 2018.

183. Shi Chen, Rafal P. Wiewiora, Fanwang Meng, Nicolas Babault, Anqi Ma, Wenyu Yu, Kun Qian, Hao Hu, Hua Zou, Junyi Wang, Shijie Fan, Gil Blum, Fabio Pittella-Silva, Kyle A. Beauchamp, Wolfram Tempel, Hualiang Jiang, Kaixian Chen, Robert Skene, Y. George Zheng, Peter J. Brown, Jian Jin, Cheng Luo, John D. Chodera, and Minkui Luo. The dynamic conformational landscapes of the protein methyltransferase SETD8. *eLife*, 8: e45403, 2019.

184. D. Trapl, I. Horvacanin, V. Mareska, F. Ozcelik, G. Unal, and V. Spiwok. Anncolvar: Approximation of Complex Collective Variables by Artificial Neural Networks for Analysis and Biasing of Molecular Simulations. *Front. Mol. Biosci.*, 6:25, 2019.

185. Katherine Henzler-Wildman and Dorothee Kern. Dynamic personalities of proteins. *Nature*, 450(7172):964–972, dec 2007.

186. Vern L. Schramm. Enzymatic transition states, transition-state analogs, dynamics, thermodynamics, and lifetimes. *Annu. Rev. Biochem.*, 80(1):703–732, July 2011.

187. G. M. Lee and C. S. Craik. Trapping moving targets with small molecules. *Science*, 324 (5924):213–215, apr 2009.

188. Brooke E. Husic, Robert T. McGibbon, Mohammad M. Sultan, and Vijay S. Pande. Optimized parameter selection reveals trends in Markov state models for protein folding. *J. Chem. Phys.*, 145(19):194103, nov 2016.

189. Frank Noé and Cecilia Clementi. Collective variables for the study of long-time kinetics from molecular trajectories: Theory and methods. *Curr. Opin. Struc. Biol.*, 43:141–147, April 2017.

190. Gregory R. Bowman. An overview and practical guide to building Markov state models. In *An Introduction to Markov State Models and Their Application to Long Timescale Molecular Simulation*, Advances in Experimental Medicine and Biology, pages 7–22. Springer Netherlands, 2014.

191. Shoshana J. Wodak, Emanuele Paci, Nikolay V. Dokholyan, Igor N. Berezovsky, Amnon Horovitz, Jing Li, Vincent J. Hilser, Ivet Bahar, John Karanicolas, Gerhard Stock, Peter Hamm, Roland H. Stote, Jerome Eberhardt, Yassmine Chebaro, Annick Dejaegere, Marco Cecchini, Jean-Pierre Changeux, Peter G. Bolhuis, Jocelyne Vreede, Pietro Faccioli, Simone Orioli, Riccardo Ravasio, Le Yan, Carolina Brito, Matthieu Wyart, Paraskevi Gkeka, Ivan Rivalta, Giulia Palermo, J. Andrew McCammon, Joanna Panecka-Hofman, Rebecca C. Wade, Antonella Di Pizio, Masha Y. Niv, Ruth Nussinov, Chung-Jung Tsai, Hyunbum Jang, Dzmitry Padhorny, Dima Kozakov, and Tom McLeish. Allostery in its many disguises: From theory to applications. *Structure*, 27(4):566–578, 2019. doi: https://doi.org/10.1016/j.str.2019.01.003.

192. Giacomo Fiorin, Michael L. Klein, and Jérôme Hénin. Using collective variables to drive molecular dynamics simulations. *Molecular Physics*, 111(22-23, SI):3345–3362, 2013. ISSN 0026-8976. doi: 10.1080/00268976.2013.813594.

193. Peter Man-Un Ung, Rayees Rahman, and Avner Schlessinger. Redefining the protein kinase conformational space with machine learning. *Cell Chemical Biology*, 25(7):916–924, 2018. doi: doi:10.1016/j.chembiol.2018.05.002.

194. Matteo T. Degiacomi. Coupling molecular dynamics and deep learning to mine protein conformational space. *Structure*, 27(6):1034–1040, 2019. doi: 10.1016/j.str.2019.03.018.

195. Díaz Óscar, James A.R. Dalton, and Jesús Giraldo. Artificial intelligence: A novel approach for drug discovery. *Trends in Pharmacological Sciences*, 40(8):550–551, 2019. doi: 10.1016/j.tips.2019.06.005.

196. Massimiliano Bonomi, Davide Branduardi, Giovanni Bussi, Carlo Camilloni, Davide Provasi, Paolo Raiteri, Davide Donadio, Fabrizio Marinelli, Fabio Pietrucci, Ricardo A. Broglia, and Michele Parrinello. Plumed: A portable plugin for free-energy calculations with molecular dynamics. *Computer Physics Communications*, 180(10):1961–1972, 2009. ISSN 0010-4655. doi: https://doi.org/10.1016/j.cpc.2009.05.011.

197. David A. Case, Thomas E. Cheatham III, Tom Darden, Holger Gohlke, Ray Luo, Kenneth M. Merz Jr., Alexey Onufriev, Carlos Simmerling, Bing Wang, and Robert J. Woods. The Amber biomolecular simulation programs. *Journal of Computational Chemistry*, 26 (16):1668–1688, 2005. doi: 10.1002/jcc.20290.

198. H.J.C. Berendsen, D. van der Spoel, and R. van Drunen. GROMACS: A message-passing parallel molecular dynamics implementation. *Computer Physics Communications*, 91(1): 43–56, 1995. ISSN 0010-4655. doi: https://doi.org/10.1016/0010-4655(95)00042-E.

199. James C. Phillips, Rosemary Braun, Wei Wang, James Gumbart, Emad Tajkhorshid, Elizabeth Villa, Christophe Chipot, Robert D. Skeel, Laxmikant Kalé, and Klaus Schulten. Scalable molecular dynamics with NAMD. *J. Comput. Chem.*, 26(16):1781–1802, 2005. doi: 10.1002/jcc.20289.

200. P. Eastman, J. Swails, J. D. Chodera, R. T. McGibbon, Y. Zhao, K. A. Beauchamp, L. P. Wang, A. C. Simmonett, M. P. Harrigan, C. D. Stern, R. P. Wiewiora, B. R. Brooks, and V. S. Pande. OpenMM 7: Rapid development of high performance algorithms for molecular dynamics. *PLoS Comput. Biol.*, 13(7):e1005659, 2017.

201. François Chollet et al. Keras, 2015. https://keras.io.

202. Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Greg S. Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Ian Goodfellow, Andrew Harp, Geoffrey Irving, Michael Isard, Yangqing Jia, Rafal Jozefowicz, Lukasz Kaiser, Manjunath Kudlur, Josh Levenberg, Dan Mané, Rajat Monga, Sherry Moore, Derek Murray, Chris Olah, Mike Schuster, Jonathon Shlens, Benoit Steiner, Ilya Sutskever, Kunal Talwar, Paul Tucker, Vincent Vanhoucke, Vijay Vasudevan, Fernanda Viégas, Oriol Vinyals, Pete Warden, Martin Wattenberg, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng. TensorFlow: Large-scale machine learning on heterogeneous systems, 2015. Software available from tensorflow.org.

203. Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in pytorch. In *NIPS 2017 Workshop on Autodiff*, 2017.

204. Jeremy Howard et al. fastai, 2018. https://github.com/fastai/fastai.

205. Matthew P. Harrigan, Mohammad M. Sultan, Carlos X. Hernández, Brooke E. Husic, Peter Eastman, Christian R. Schwantes, Kyle A. Beauchamp, Robert T. McGibbon, and Vijay S. Pande. MSMBuilder: Statistical models for biomolecular dynamics. *Biophysical Journal*, 112(1):10–15, 2017. ISSN 0006-3495. doi: https://doi.org/10.1016/j.bpj.2016.10.042.

206. Robert T. McGibbon, Kyle A. Beauchamp, Matthew P. Harrigan, Christoph Klein, Jason M. Swails, Carlos X. Hernández, Christian R. Schwantes, Lee-Ping Wang, Thomas J. Lane, and Vijay S. Pande. MDTraj: A modern open library for the analysis of molecular dynamics trajectories. *Biophysical Journal*, 109(8):1528–1532, 2015. doi: 10.1016/j.bpj.2015.08.015.

207. David W. H. Swenson, Jan-Hendrik Prinz, Frank Noe, John D. Chodera, and Peter G. Bolhuis. OpenPathSampling: A Python framework for path sampling simulations. 2. Building and customizing path ensembles and sample schemes. *Journal of Chemical Theory and Computation*, 15(2):837–856, 2019. doi: 10.1021/acs.jctc.8b00627.

208. Jean-Christophe Carry, François Clerc, Hervé Minoux, Laurent Schio, Jacques Mauger, Anil Nair, Eric Parmantier, Ronan Le Moigne, Cecile Delorme, Jean-Paul Nicolas, Alain Krick, Pierre-Yves Abecassis, Veronique Crocq-Stuerga, Stephanie Pouzieux, Laure Delarbre, Sebastien Maignan, Thomas Bertrand, Kirsten Bjergarde, Nina Ma, Sylvette Lachaud, Houlfa Guizani, Rémi Lebel, Gilles Doerflinger, Sylvie Monget, Sebastien Perron, Francis Gasse, Odile Angouillant-Boniface, Bruno Filoche-Romme, Michel Murer, Sylvie Gontier, Celine Prevost, Marie-Line Monteiro, and Cecile Combeau. Sar156497, an exquisitely selective inhibitor of aurora kinases. *J. Med. Chem.*, 58(1):362–375, 2015.

209. DrugBank. https://www.drugbank.ca, 2020.

210. Clinical Trials. https://clinicaltrials.gov, 2020.

211. Paraskevi Gkeka, Thomas Evangelidis, Maria Pavlaki, Vasiliki Lazani, Savvas Christoforidis, Bogos Agianian, and Zoe Cournia. Investigating the structure and dynamics of the Pik3ca wild-type and H1047R oncogenic mutant. *PLOS Comput. Biol.*, 10(10):1–12, 10 2014. doi: 10.1371/journal.pcbi.1003895.

212. Paraskevi Gkeka, Alexandra Papafotika, Savvas Christoforidis, and Zoe Cournia. Exploring a non-ATP pocket for potential allosteric modulation of PI3K$\alpha$. *J. Phys. Chem. B*, 119 (3):1002–1016, 2015. doi: 10.1021/jp506423e.

213. Jochen S. Hub and Bert L. de Groot. Detection of functional modes in protein dynamics. *PLOS Comput. Biol.*, 5(8):1–13, 08 2009. doi: 10.1371/journal.pcbi.1000480.

214. James Philip O'Connell, John Robert Porter, Alastair Lawson, Boris Kroeplien, Stephen Edward Rapecki, Timothy John Norman, Graham John Warrellow, Tracy Lynn Arakaki, Alex Buntin Burgin, William Ross Pitt, Mark Daniel Calmiano, David Andreas Schubert, Daniel John Lightwood, and Rebecca Jayne Wootton. Novel TNF$\alpha$ structure for use in therapy, 2015. PCT / E P2015 / 074491.

215. Amir Barati Farimani, Evan N. Feinberg, and Vijay Pande. Binding pathway of opiates to $\mu$-opioid receptors revealed by machine learning. *Biophys. J.*, 114:62a–63a, 02 2018. doi: 10.1016/j.bpj.2017.11.390.

216. Evan N. Feinberg, Amir Barati Farimani, Rajendra Uprety, Amanda Hunkele, Gavril Pasternak, Susruta Majumdar, and Vijay Pande. Machine learning harnesses molecular dynamics to discover new $\mu$-opioid chemotypes. *arXiv preprint*, 1803.04479, 03 2018.

217. Michael Schmidt and Hod Lipson. Distilling free-form natural laws from experimental data. *Science*, 324(5923):81–85, 2009.