



HHS Public Access

Author manuscript

J Chem Theory Comput. Author manuscript; available in PMC 2021 August 11.

Published in final edited form as:

J Chem Theory Comput. 2020 August 11; 16(8): 4757–4775. doi:10.1021/acs.jctc.0c00355.

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

Paraskevi Gkeka¹, Gabriel Stoltz^{2,3}, Amir Barati Farimani⁴, Zineb Belkacemi¹, Michele Ceriotti⁵, John Chodera⁶, Aaron R. Dinner⁷, Andrew Ferguson⁸, Jean-Bernard Maillet⁹, Hervé Minoux¹⁰, Christine Peter¹¹, Fabio Pietrucci¹², Ana Silveira⁶, Alexandre Tkatchenko¹³, Zofia Trstanova¹⁴, Rafal Wiewiora⁶, Tony Lelièvre^{2,3}

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

Abstract

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.

Keywords

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

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 meeting¹ brought together a diverse 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:

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

- **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 (non-linear) 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 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 PESTo 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 determination 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 directly 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 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 sampling 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 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 committer 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 AIQuraishi (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 where 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 proteins using their metastable states as target conformations, while in Section B.2 we give some examples where 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 α , 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 α 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 α for instance shows the importance of how subtle changes in protein conformation can translate into a distorted trimeric assembly of TNF α , impacting downstream signaling of TNFR1. Small compounds stabilizing this asymmetrical TNF α trimer can then be designed to treat or prevent TNF α -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 trajectories 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.

ACKNOWLEDGEMENTS

This review paper was written following a CECAM (Centre Européen de Calcul Atomique et Moléculaire) discussion meeting, hosted at the Sanofi Campus of Gentilly. The authors thank the CECAM as well as Sanofi for making this event possible. Moreover, the PG, GS and TL thank Dr. Marc Bianciotto for his proof reading and feedback.

6. Bibliography

1. Zhang Yue-Yu, Niu Haiyang, Piccini GiovanniMaria, Mendels Dan, and Parrinello Michele. Improving collective variables: The case of crystallization. *The Journal of Chemical Physics*, 150(9):094509, 2019. [PubMed: 30849916]
2. Behler J Atom-centered symmetry functions for constructing high-dimensional neural network potentials. *The Journal of Chemical Physics*, 134(7):074106, 2011. [PubMed: 21341827]
3. Bartók Albert P., Kondor Risi, and Csányi Gábor. On representing chemical environments. *Phys. Rev. B*, 87:184115, 2013.
4. Wales DJ Perspective: Insight into reaction coordinates and dynamics from the potential energy landscape. *J. Chem. Phys.*, 142(13):130901, 2015. doi: 10.1063/1.4916307. [PubMed: 25854218]
5. Peters Baron. Reaction coordinates and mechanistic hypothesis tests. *Annu. Rev. Phys. Chem.*, 67(1):669–690, 2016. [PubMed: 27090846]
6. McGibbon Robert T, Husic Brooke E, and Pande Vijay S. Identification of simple reaction coordinates from complex dynamics. *J. Chem. Phys.*, 146(4):044109, 2016. doi: 10.1063/1.4974306.
7. Chmiela Stefan, Sauceda Huziel E., Müller Klaus-Robert, and Tkatchenko Alexandre. Towards exact molecular dynamics simulations with machine-learned force fields. *Nat. Commun.*, 9:3887, 2018. doi: 10.1038/s41467-018-06169-2. [PubMed: 30250077]
8. Wang Jiang, Olsson Simon, Wehmeyer Christoph, Pérez Adrià, Charron Nicholas E., Fabritiis Gianni de, Noé Frank, and Clementi Cecilia. Machine learning of coarse-grained molecular dynamics force fields. *ACS Cent. Sci.*, 5(5):755–767, 2019. doi: 10.1021/acscentsci.8b00913. [PubMed: 31139712]
9. Häse Florian, Galván Ignacio Fernández, Aspuru-Guzik Alán, Lindh Roland, and Vacher Morgane. 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. [PubMed: 30881655]
10. Jolliffe Ian. *Principal Component Analysis*. Wiley Online Library, 2002.
11. Coifman Ronald R and Lafon Stéphane. Diffusion maps. *Applied and Computational Harmonic Analysis*, 21(1):5–30, 2006.
12. Coifman Ronald R, Kevrekidis Ioannis G, Lafon Stéphane, Maggioni Mauro, and Nadler Boaz. Diffusion maps, reduction coordinates, and low dimensional representation of stochastic systems. *Multiscale Model. Simul.*, 7(2):842–864, 2008.

13. Chen Wei and Ferguson Andrew L. 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. [PubMed: 30368832]
14. Chen Wei, Tan Aik Rui, and Ferguson Andrew L. Collective variable discovery and enhanced sampling using autoencoders: Innovations in network architecture and error function design. *J. Chem. Phys.*, 149(7):072312, 2018. [PubMed: 30134681]
15. Lamim Ribeiro João Marcelo, Bravo Pablo, Wang Yihang, and Tiwary Pratyush. Reweighted autoencoded variational bayes for enhanced sampling (RAVE). *J. Chem. Phys.*, 149(7):072301, 2018. [PubMed: 30134694]
16. Wehmeyer Christoph and Noé Frank. Time-lagged autoencoders: Deep learning of slow collective variables for molecular kinetics. *J. Chem. Phys.*, 148(24):241703, 6 2018. [PubMed: 29960344]
17. Ma Ao and Dinner Aaron R. Automatic method for identifying reaction coordinates in complex systems. *J. Phys. Chem. B*, 109(14):6769–6779, 2005. [PubMed: 16851762]
18. Shapeev Alexander V. Moment tensor potentials: A class of systematically improvable interatomic potentials. *Multiscale Modeling & Simulation*, 14(3):1153–1173, 1 2016. doi: 10.1137/15m1054183.
19. Chen H, Lu J, and Ortner C. Thermodynamic limit of crystal defects with finite temperature tight binding. *Arch. Ration. Mech. Anal.*, 230:701–733, 2018.
20. Lunghi Alessandro and Sanvito Stefano. 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. [PubMed: 31172029]
21. Artrith N and Behler J High-dimensional neural network potentials for metal surfaces: A prototype study for copper. *Phys. Rev. B*, 85:045439, 2012.
22. Podryabinkin Evgeny V., Tikhonov Evgeny V., Shapeev Alexander V., and Oganov Artem R. Accelerating crystal structure prediction by machine-learning interatomic potentials with active learning. *Phys. Rev. B*, 99:064114, 2019.
23. Gubaev Konstantin, Podryabinkin Evgeny V., Hart Gus L.W., and Shapeev Alexander V. Accelerating high-throughput searches for new alloys with active learning of interatomic potentials. *Computational Materials Science*, 156:148–156, 2019.
24. Huan Tran Doan, Batra Rohit, Chapman James, Kim Chiho, Chandrasekaran Anand, and Ramprasad Rampi. Iterative-learning strategy for the development of application-specific atomistic force fields. *The Journal of Physical Chemistry C*, 123(34):20715–20722, 2019.
25. Jinnouchi Ryosuke, Karsai Ferenc, and Kresse Georg. On-the-fly machine learning force field generation: Application to melting points. *Phys. Rev. B*, 100:014105, 2019.
26. Deringer Volker L., Proserpio Davide M., Csányi Gábor, and Pickard Chris J. Data-driven learning and prediction of inorganic crystal structures. *Faraday Discuss*, 211:45–59, 2018. [PubMed: 30043006]
27. Eickenberg M, Exarchakis G, Hirn M, Mallat S, and Thiry L Solid harmonic wavelet scattering for predictions of molecule properties. *J. Chem. Phys.*, 148(24):241732, 2018. [PubMed: 29960365]
28. Ferré G, Haut T, and Barros K Learning molecular energies using localized graph kernels. *The Journal of Chemical Physics*, 146(11):114107, 2017. [PubMed: 28330348]
29. Weininger David. 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. Steinhardt Paul J., Nelson David R., and Ronchetti Marco. Bond-orientational order in liquids and glasses. *Physical Review B*, 28(2):784–805, 7 1983. doi: 10.1103/physrevb.28.784.
31. Pietrucci Fabio and Laio Alessandro. 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, 8 2009. doi: 10.1021/ct900202f. [PubMed: 26616604]
32. Bartók Albert P., Payne Mike C., Kondor Risi, and Csányi Gábor. Gaussian approximation potentials: The accuracy of quantum mechanics, without the electrons. *Phys. Rev. Lett*, 104(13):136403, 4 2010. doi: 10.1103/PhysRevLett.104.136403. [PubMed: 20481899]

33. Behler Jörg and Parrinello Michele. Generalized neural-network representation of high-dimensional potential-energy surfaces. *Phys. Rev. Lett*, 98(14):146401, 4 2007. doi: 10.1103/PhysRevLett.98.146401. [PubMed: 17501293]
34. Faber Felix A., Christensen Anders S., Huang Bing, and Lilienfeld O. Anatole von. Alchemical and structural distribution based representation for universal quantum machine learning. *The Journal of Chemical Physics*, 148(24):241717, 6 2018. doi: 10.1063/1.5020710. [PubMed: 29960351]
35. Willatt Michael J., Musil Félix, and Ceriotti Michele. Atom-density representations for machine learning. *The Journal of Chemical Physics*, 150(15):154110, 4 2019. doi: 10.1063/1.5090481. [PubMed: 31005079]
36. Bartók Albert P., De Sandip, Poelking Carl, Bernstein Noam, Kermode ames R., Csányi Gábor, and Ceriotti Michele. Machine learning unifies the modeling of materials and molecules. *Science Advances*, 3(12):e1701816, 12 2017. doi: 10.1126/sciadv.1701816. [PubMed: 29242828]
37. Zuo Yunxing, Chen Chi, Li Xiangguo, Deng Zhi, Chen Yiming, Behler Jörg, Csányi Gábor, Shapeev Alexander V., Thompson Aidan P., Wood Mitchell A., and Ong Shyue Ping. Performance and cost assessment of machine learning interatomic potentials. *The Journal of Physical Chemistry A*, 124(4):731–745, 2020. [PubMed: 31916773]
38. Behler Jörg. Perspective: Machine learning potentials for atomistic simulations. *J. Chem. Phys*, 145(17):170901, 2016. doi: 10.1063/1.4966192. [PubMed: 27825224]
39. Braams Bastiaan J. and Bowman Joel M. Permutationally invariant potential energy surfaces in high dimensionality. *Int. Rev. Phys. Chem*, 28:577–606, 2009. doi: 10.1080/01442350903234923.
40. Chmiela Stefan, Tkatchenko Alexandre, Sauceda Huziel E., Poltavsky Igor, Schütt Kristof T., and Müller Klaus-Robert. Machine learning of accurate energy-conserving molecular force fields. *Sci. Adv*, 3(5):e1603015, 2017. doi: 10.1126/sciadv.1603015. [PubMed: 28508076]
41. Schütt Kristof T, Arbabzadah Farhad, Chmiela Stefan, Müller Klaus R, and Tkatchenko Alexandre. Quantum-chemical insights from deep tensor neural networks. *Nat. Commun*, 8:13890, 2017. [PubMed: 28067221]
42. Schütt KT, Sauceda HE, Kindermans P-J, Tkatchenko A, and Müller K-R SchNet - a deep learning architecture for molecules and materials. *J. Chem. Phys*, 148(24):241722, 2018. [PubMed: 29960322]
43. Qu Chen and Bowman Joel M. 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. [PubMed: 30981221]
44. Deringer Volker L. and Csányi Gábor. Machine learning based interatomic potential for amorphous carbon. *Phys. Rev. B*, 95:094203, 2017.
45. Ambrosetti Alberto, Ferri Nicola, DiStasio Robert A., and Tkatchenko Alexandre. Wavelike charge density fluctuations and van der Waals interactions at the nanoscale. *Science*, 351: 1171–1176, 2016. doi: 10.1126/science.aae0509. [PubMed: 26965622]
46. Hermann Jan, DiStasio Robert A., and Tkatchenko Alexandre. 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. [PubMed: 28272886]
47. Bereau Tristan, DiStasio Robert A. Jr, Tkatchenko Alexandre, and Von Lilienfeld O. Anatole. 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. [PubMed: 29960330]
48. Grisafi Andrea and Ceriotti Michele. Incorporating long-range physics in atomic-scale machine learning. *The Journal of Chemical Physics*, 151(20):204105, 11 2019. doi: 10.1063/1.5128375. [PubMed: 31779318]
49. Glielmo Aldo, Zeni Claudio, and Vita Alessandro De. Efficient nonparametricn-body force fields from machine learning. *Physical Review B*, 97(18), 5 2018. doi: 10.1103/physrevb.97.184307.
50. Veit Max, Jain Sandeep Kumar, Bonakala Satyanarayana, Rudra Indranil, Hohl Detlef, and Csányi Gábor. Equation of state of fluid methane from first principles with machine learning potentials. *Journal of Chemical Theory and Computation*, 15(4):2574–2586, 2 2019. doi: 10.1021/acs.jctc.8b01242. [PubMed: 30794393]

51. Ziegler James F. and Biersack Jochen P. The stopping and range of ions in matter. In Bromley D. Allan, editor, *Treatise on Heavy-Ion Science: Volume 6: Astrophysics, Chemistry, and Condensed Matter*, pages 93–129. Springer US, Boston, MA, 1985.
52. Lemke Tobias and Peter Christine. 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. [PubMed: 29120633]
53. Hunkler S, Lemke T, Peter C, and Kukharenko O Back-mapping based sampling: Coarse grained free energy landscapes as a guideline for atomistic exploration. *J. Chem. Phys*, 151(15):154102, 2019. [PubMed: 31640363]
54. Peter Christine and Kremer Kurt. Multiscale simulation of soft matter systems - from the atomistic to the coarse-grained level and back. *Soft Matter*, 5(22):4357–4366, 2009.
55. Rudzinski Joseph F and Noid WG. Coarse-graining entropy, forces, and structures. *J. Chem. Phys*, 135(21):214101, 2011. [PubMed: 22149773]
56. Noid WG. Perspective: Coarse-grained models for biomolecular systems. *J. Chem. Phys*, 139(9):090901, 2013. [PubMed: 24028092]
57. Potestio Raffaello, Peter Christine, and Kremer Kurt. Computer simulations of soft matter: Linking the scales. *Entropy*, 16(8):4199–4245, 2014.
58. Shell MS. Coarse-graining with the relative entropy. *Adv. Chem. Phys*, pages 395–441, 2016.
59. John ST and Csányi Gabor. Many-body coarse-grained interactions using Gaussian approximation potentials. *J. Phys. Chem. B*, 121(48):10934–10949, 2017. [PubMed: 29117675]
60. Zhang Linfeng, Han Jiequn, Wang Han, Car Roberto, and DeePCG Weinan E.: Constructing coarse-grained models via deep neural networks. *J. Chem. Phys*, 149(3):034101, 2018. [PubMed: 30037247]
61. Noid WG, Chu Jih-Wei, Ayton Gary S, Krishna Vinod, Izvekov Sergei, Voth Gregory A, Das Avisek, and Andersen Hans C. The multiscale coarse-graining method. I. A rigorous bridge between atomistic and coarse-grained models. *J. Chem. Phys*, 128(24):244114, 2008. [PubMed: 18601324]
62. Garrido L and Juste A On the determination of probability density functions by using neural networks. *Comput. Phys. Commun*, 115:25–31, 1998.
63. Ferguson Andrew L, Panagiotopoulos Athanassios Z, Debenedetti Pablo G, and Kevrekidis Ioannis G. Systematic determination of order parameters for chain dynamics using diffusion maps. *P. Natl. Acad. Sci. USA*, 107(31):13597–13602, 2010.
64. Ferguson Andrew L, Panagiotopoulos Athanassios Z, Kevrekidis Ioannis G, and Debenedetti Pablo G. Nonlinear dimensionality reduction in molecular simulation: The diffusion map approach. *Chem. Phys. Lett*, 509(1–3):1–11, 2011.
65. Wang J and Ferguson AL. Nonlinear machine learning in simulations of soft and biological materials. *Mol. Simul*, 44(13–14):1090–1107, 2018.
66. Pietrucci Fabio. 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. Ichiye Toshiko and Karplus Martin. 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. García Angel E. Large-amplitude nonlinear motions in proteins. *Phys. Rev. Lett*, 68(17): 2696, 1992. [PubMed: 10045464]
69. Amadei A, M Linssen AB, and Berendsen HJC Essential dynamics of proteins. *Proteins: Structure, Function, and Bioinformatics*, 17(4):412–425, 1993.
70. Ferguson Andrew L. Machine learning and data science in soft materials engineering. *J. Phys. Condens. Matter*, 30(4):043002, 2017.
71. Schölkopf B 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. Sittel F and Stock G Perspective: Identification of collective variables and metastable states of protein dynamics. *J. Chem. Phys*, 149(15):150901, 2018. [PubMed: 30342445]

73. Trstanova Zofia, Leimkuhler Ben, and Lelièvre Tony. Local and global perspectives on diffusion maps in the analysis of molecular systems. *Proc. R. Soc. A*, 476(2233):20190036, 2020. [PubMed: 32082050]
74. Jónsson Hannes, Mills G, and Jacobsen KW Nudged elastic band method for finding minimum energy paths of transitions. In Berne BJ, Ciccotti G, and Coker DF, editors, *Classical and Quantum Dynamics in Condensed Phase Simulations*, pages 385–404. World Scientific, 1998.
75. Weinan E, Weiqing Ren, and Vanden-Eijnden Eric. String method for the study of rare events. *Phys. Rev. B*, 66(5):52301, 2002.
76. Weinan E and Vanden-Eijnden Eric. Transition-path theory and path-finding algorithms for the study of rare events. *Annu. Rev. Phys. Chem*, 61:391–420, 2010. [PubMed: 18999998]
77. Lelièvre T and Stoltz G Partial differential equations and stochastic methods in molecular dynamics. *Acta Numerica*, 25:681–880, 2016.
78. Bolhuis Peter G, Chandler David, Dellago Christoph, and Geissler Phillip L. Transition path sampling: Throwing ropes over rough mountain passes, in the dark. *Annu. Rev. Phys. Chem*, 53(1):291–318, 2002. [PubMed: 11972010]
79. Banushkina Polina V and Krivov Sergei V. Optimal reaction coordinates. *WIREs: Comput. Mol. Sci*, 6(6):748–763, 2016.
80. Jung Hendrik, Covino Roberto, and Hummer Gerhard. Artificial intelligence assists discovery of reaction coordinates and mechanisms from molecular dynamics simulations. arXiv preprint, 1901.04595, 2019.
81. Loève Michel. *Probability Theory: Foundations, Random Sequences*. Van Nostrand, 1955.
82. Sirovich Lawrence. Turbulence and the dynamics of coherent structures. I. Coherent structures. *Q. Appl. Math*, 45(3):561–571, 1987.
83. Sirovich Lawrence. Turbulence and the dynamics of coherent structures. II. Symmetries and transformations. *Q. Appl. Math*, 45(3):573–582, 1987.
84. Park HM and Cho DH. The use of the Karhunen-Loeve decomposition for the modeling of distributed parameter systems. *Chem. Eng. Sci*, 51(1):81–98, 1996.
85. Chatterjee Anindya. An introduction to the proper orthogonal decomposition. *Current Science*, 78:808–817, 2000.
86. Liang YC, Lee HP, Lim SP, Lin WZ, Lee KH, and Wu CG. Proper orthogonal decomposition and its applications—Part I: Theory. *J. Sound Vib*, 252(3):527–544, 2002.
87. Schölkopf Bernhard, Smola Alexander, and Müller Klaus-Robert. Kernel principal component analysis. In Gerstner Wulfram, Germond Alain, Hasler Martin, and Nicoud Jean-Daniel, 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. Kramer Mark A. Nonlinear principal component analysis using autoassociative neural networks. *AIChE Journal*, 37(2):233–243, 1991.
89. Nguyen Phuong H. Complexity of free energy landscapes of peptides revealed by non-linear principal component analysis. *Proteins: Structure, Function, and Bioinformatics*, 65 (4):898–913, 2006. ISSN 1097–0134. doi: 10.1002/prot.21185.
90. Scholz Matthias, Fraunholz Martin, and Selbig Joachim. Nonlinear principal component analysis: neural network models and applications. In Gorban Alexander N., Kégl Balázs, Wunsch Donald C., and Zinovyev Andrei Y., editors, *Principal Manifolds for Data Visualization and Dimension Reduction*, pages 44–67. Springer, Berlin Heidelberg, 2008.
91. Comon Pierre. Independent component analysis, a new concept? *Signal Processing*, 36 (3):287–314, 1994.
92. Borg Ingwer and Groenen Patrick JF. *Modern Multidimensional Scaling: Theory and Applications*. Springer Science & Business Media, 2005.
93. Ceriotti Michele, Tribello Gareth A., and Parrinello Michele. Simplifying the representation of complex free-energy landscapes using sketch-map. *Proceedings of the National Academy of Sciences*, 108(32):13023–13028, 7 2011. doi: 10.1073/pnas.1108486108.
94. Roweis Sam T and Saul Lawrence K. Nonlinear dimensionality reduction by locally linear embedding. *Science*, 290(5500):2323–2326, 2000. [PubMed: 11125150]

95. Zhang Zhenyue and Wang Jing. MLLE: Modified locally linear embedding using multiple weights. In Schölkopf Bernhard, Platt John, and Hofmann Thomas, editors, *Advances in Neural Information Processing Systems 19: Proceedings of the 2006 Conference*, pages 1593–1600, Cambridge, December 2007. MIT Press.
96. Das Payel, Moll Mark, Stamati Hernan, Kavradi Lydia E, and Clementi Cecilia. Low-dimensional, free-energy landscapes of protein-folding reactions by nonlinear dimensionality reduction. *P. Natl. Acad. Sci. USA*, 103(26):9885–9890, 2006.
97. Tenenbaum Joshua B, De Silva Vin, and Langford John C. A global geometric framework for nonlinear dimensionality reduction. *Science*, 290(5500):2319–2323, 2000. [PubMed: 11125149]
98. Silva Vin D. and Tenenbaum Joshua B. Global versus local methods in nonlinear dimensionality reduction. In Thrun S and Obermayer K, editors, *Advances in Neural Information Processing Systems 15*, pages 705–712. MIT Press, Cambridge, MA, 2002.
99. Wang Jianzhong. Local tangent space alignment. In *Geometric Structure of High-Dimensional Data and Dimensionality Reduction*, pages 221–234. Springer, Berlin Heidelberg, 2012.
100. Weinberger Kilian Q and Saul Lawrence K. Unsupervised learning of image manifolds by semidefinite programming. *Int. J. Comput. Vision*, 70(1):77–90, 2006.
101. Belkin Mikhail and Niyogi Partha. Laplacian eigenmaps for dimensionality reduction and data representation. *Neural Computation*, 15(6):1373–1396, 2003.
102. Donoho David L and Grimes Carrie. Hessian eigenmaps: Locally linear embedding techniques for high-dimensional data. *P. Natl. Acad. Sci. USA*, 100(10):5591–5596, 2003.
103. Coifman Ronald R, Lafon Stephane, Lee Ann B, Maggioni Mauro, Nadler Boaz, Warner Frederick, and Zucker Steven W. 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. Shamsi Zahra, Cheng Kevin J, and Shukla Diwakar. REinforcement learning based Adaptive samPLing: REAPing Rewards by Exploring Protein Conformational Landscapes. *J. Phys. Chem. B*, 122:8386–8395, 2018. [PubMed: 30126271]
105. Chiavazzo Eliodoro, Covino Roberto, Coifman Ronald R, Gear C William, Georgiou Anastasia S, Hummer Gerhard, and Kevrekidis Ioannis G. Intrinsic map dynamics exploration for uncharted effective free-energy landscapes. *Proc. Natl. Acad. Sci. USA*, 114 (28):E5494–E5503, 2017. [PubMed: 28634293]
106. Ferguson Andrew L, Panagiotopoulos Athanassios Z, Debenedetti Pablo G, and Kevrekidis Ioannis G. Integrating diffusion maps with umbrella sampling: Application to alanine dipeptide. *J. Chem. Phys.* 134(13):135103, 2011. [PubMed: 21476776]
107. Tribello GA, Ceriotti M, and Parrinello M Using sketch-map coordinates to analyze and bias molecular dynamics simulations. *Proceedings of the National Academy of Sciences*, 109(14):5196–5201, 3 2012. doi: 10.1073/pnas.1201152109.
108. Abrams Cameron F and Vanden-Eijnden Eric. On-the-fly free energy parameterization via temperature accelerated molecular dynamics. *Chem. Phys. Lett*, 547:114–119, 2012. [PubMed: 23226688]
109. Hashemian Behrooz, Millán Daniel, and Arroyo Marino. Modeling and enhanced sampling of molecular systems with smooth and nonlinear data-driven collective variables. *J. Chem. Phys.* 139(21):214101, 2013. [PubMed: 24320358]
110. Li Chun-Guang, Guo Jun, Chen Guang, Nie Xiang-Fei, and Yang Zhen. A version of Isomap with explicit mapping. In *2006 International Conference on Machine Learning and Cybernetics*, pages 3201–3206. IEEE, 2006.
111. Spiwok Vojtěch and Králová Blanka. Metadynamics in the conformational space nonlinearly dimensionally reduced by Isomap. *J. Chem. Phys.* 135(22):224504, 2011. [PubMed: 22168700]
112. Branduardi Davide, Gervasio Francesco Luigi, and Parrinello Michele. From A to B in free energy space. *J. Chem. Phys.* 126(5):054103, 2007. [PubMed: 17302470]
113. Preto Jordane and Clementi Cecilia. Fast recovery of free energy landscapes via diffusion-map-directed molecular dynamics. *Phys. Chem. Chem. Phys.* 16(36):19181–19191, 2014. [PubMed: 24955434]

114. Zheng Wenwei, Rohrdanz Mary A, and Clementi Cecilia. Rapid exploration of configuration space with diffusion-map-directed molecular dynamics. *J. Phys. Chem. B*, 117(42):12769–12776, 2013. [PubMed: 23865517]
115. Noé Frank and Nüske Feliks. A variational approach to modeling slow processes in stochastic dynamical systems. *Multiscale Model. Sim*, 11(2):635–655, 2013.
116. Mezi Igor. Spectral properties of dynamical systems, model reduction and decompositions. *Nonlinear Dynamics*, 41(1–3):309–325, 2005.
117. Williams Matthew O, Kevrekidis Ioannis G, and Rowley Clarence W. A data-driven approximation of the Koopman operator: Extending dynamic mode decomposition. *J. Nonlinear Sci*, 25(6):1307–1346, 2015.
118. Wu Hao, Nüske Feliks, Paul Fabian, Klus Stefan, Koltai Péter, and Noé Frank. Variational Koopman models: slow collective variables and molecular kinetics from short offequilibrium simulations. *J. Chem. Phys*, 146(15):154104, 2017. [PubMed: 28433026]
119. Pérez-Hernández Guillermo, Paul Fabian, Giorgino Toni, Fabritiis Gianni De, and Noé Frank. Identification of slow molecular order parameters for Markov model construction. *J. Chem. Phys*, 139(1):015102, 2013. [PubMed: 23822324]
120. Nüske Feliks, Keller Bettina G, Pérez-Hernández Guillermo, Mey Antonia S J S, and Noé Frank. Variational approach to molecular kinetics. *J. Chem. Theory Comput*, 10(4):1739–1752, 2014. [PubMed: 26580382]
121. Noé Frank and Clementi Cecilia. Kinetic distance and kinetic maps from molecular dynamics simulation. *J. Chem. Theory Comput*, 11(10):5002–5011, 2015. [PubMed: 26574285]
122. Noé Frank, Banisch Ralf, and Clementi Cecilia. Commute maps: Separating slowly mixing molecular configurations for kinetic modeling. *J. Chem. Theory Comput*, 12(11):5620–5630, 2016. [PubMed: 27696838]
123. Pérez-Hernández Guillermo and Noé Frank. 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. [PubMed: 27792332]
124. Schwantes Christian R. and Pande Vijay S. 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, 3 2013. [PubMed: 23750122]
125. Klus Stefan, Nüske Feliks, Koltai Péter, Wu Hao, Kevrekidis Ioannis, Christof Schütte, and Frank Noé. Data-driven model reduction and transfer operator approximation. *J. Nonlinear Sci*, 28(3):985–1010, 2018.
126. Pande Vijay S, Beauchamp Kyle, and Bowman Gregory R. Everything you wanted to know about Markov state models but were afraid to ask. *Methods*, 52(1):99–105, 2010. [PubMed: 20570730]
127. Noé Frank and Rosta Edina. Markov models of molecular kinetics. *The Journal of Chemical Physics*, 151(19):190401, 2019. [PubMed: 31757166]
128. Schütte Christof, Fischer Alexander, Huisinga Wilhelm, and Deuffhard Peter. A direct approach to conformational dynamics based on hybrid Monte Carlo. *J. Comput. Phys*, 151(1):146–168, 1999.
129. Deuffhard P and Weber M Robust Perron cluster analysis in conformation dynamics. *Linear algebra and its applications*, 398:161–184, 2005.
130. Röblitz S and Weber M 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. Schwantes Christian R and Pande Vijay S. Modeling molecular kinetics with tICA and the kernel trick. *J. Chem. Theory Comput*, 11(2):600–608, 2015. [PubMed: 26528090]
132. Andrew Galen, Arora Raman, Bilmes Jeff, and Livescu Karen. Deep canonical correlation analysis. In Dasgupta Sanjoy and McAllester David, 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. Mardt Andreas, Pasquali Luca, Wu Hao, and Noé Frank. VAMPnets for deep learning of molecular kinetics. *Nature Communications*, 9:5, 2018.

134. Chen Wei, Sidky Hythem, and Ferguson Andrew L. Nonlinear discovery of slow molecular modes using state-free reversible VAMPnets. *J. Chem. Phys.*, 150:214114, 2019. [PubMed: 31176319]
135. Hernández Carlos X, Wayment-Steele Hannah K, Sultan Mohammad M, Husic Brooke E, and Pande Vijay S. Variational encoding of complex dynamics. *Phys. Rev. E*, 97(6): 062412, 2018. [PubMed: 30011547]
136. Wayment-Steele Hannah K and Pande Vijay S. Note: Variational encoding of protein dynamics benefits from maximizing latent autocorrelation. *J. Chem. Phys.*, 149:216101, 2018. [PubMed: 30525733]
137. Quer Jannes, Donati Luca, Keller Bettina G, and Weber Marcus. An automatic adaptive importance sampling algorithm for molecular dynamics in reaction coordinates. *SIAM J. Sci. Comput.*, 40(2):A653–A670, 2018.
138. Donati Luca and Keller Bettina G. Girsanov reweighting for metadynamics simulations. *J. Chem. Phys.*, 149(7):072335, 2018. [PubMed: 30134671]
139. Donati Lorenzo, Hartmann Carsten, and Keller Bettina G. Girsanov reweighting for path ensembles and Markov state models. *J. Chem. Phys.*, 146(24):244112, 2017. [PubMed: 28668056]
140. Nadler Boaz, Lafon Stephane, Coifman Ronald, and Kevrekidis Ioannis G. Diffusion Maps a Probabilistic Interpretation for Spectral Embedding and Clustering Algorithms. In Gorban Alexander N, Kégl Balázs, Wunsch Donald C, and Zinovyev Andrei Y, editors, *Principal Manifolds for Data Visualization and Dimension Reduction*, pages 238–260, Berlin, Heidelberg, 2008. Springer. ISBN 978–3-540–73750-6.
141. Rohrdanz Mary A., Zheng Wenwei, Maggioni Mauro, and Clementi Cecilia. Determination of reaction coordinates via locally scaled diffusion map. *J. Chem. Phys.*, 134(12):124116, 2011. [PubMed: 21456654]
142. Berry Tyrus and Harlim John. Variable bandwidth diffusion kernels. *Appl. Comput. Harmon. Anal.*, 40(1):68–96, 2016. ISSN 1096603X. doi: 10.1016/j.acha.2015.01.001.
143. Banisch Ralf, Trstanova Zofia, Bittracher Andreas, Klus Stefan, and Koltai Péter. Diffusion maps tailored to arbitrary non-degenerate Itô processes. *Appl. Comput. Harmon. Anal.*, 48:242–265, 2018. ISSN 1063–5203.
144. Collet Pierre, Martinez Servet, and Martin Jaime San. *Quasi-Stationary Distributions: Markov Chains, Diffusions and Dynamical Systems*. Springer Science & Business Media, 2012. ISBN 9783642331305.
145. Lai Rongjie and Lu Jianfeng. Point cloud discretization of Fokker–Planck operators for committor functions. *Multiscale Model. Simul.*, 16(2):710–726, 2018.
146. Prinz Jan-Hendrik, Held Martin, Smith Jeremy C, and Noé Frank. Efficient computation, sensitivity, and error analysis of committor probabilities for complex dynamical processes. *Multiscale Model. Simul.*, 9(2):545–567, 2011.
147. Thiede EH, Giannakis D, Dinner AR, and Weare J Galerkin approximation of dynamical quantities using trajectory data. *J Chem Phys*, 150(24):244111, 6 2019. [PubMed: 31255053]
148. Khoo Yuehaw, Lu Jianfeng, and Ying Lexing. Solving for high dimensional committor functions using artificial neural networks. *Research in the Mathematical Sciences*, 6:1, 2019.
149. Molgedey Lutz and Schuster Heinz Georg. Separation of a mixture of independent signals using time delayed correlations. *Phys. Rev. Lett.*, 72(23):3634–3637, 1994. [PubMed: 10056251]
150. Takano Hiroshi and Miyashita Seiji. Relaxation modes in random spin systems. *J. Phys. Soc. Jpn.*, 64(10):3688–3698, 1995.
151. Hirao Hidetomo, Koseki Sachiko, and Takano Hiroshi. Molecular dynamics study of relaxation modes of a single polymer chain. *J. Phys. Soc. Jpn.*, 66(11):3399–3405, 1997.
152. Swope William C, Pitera Jed W, and Suits Frank. Describing protein folding kinetics by molecular dynamics simulations. 1. Theory. *J. Phys. Chem. B*, 108(21):6571–6581, 2004.
153. Prinz Jan-Hendrik, Wu Hao, Sarich Marco, Keller Bettina, Senne Martin, Held Martin, Chodera John D, Schütte Christof, and Noé Frank. Markov models of molecular kinetics: Generation and validation. *J. Chem. Phys.*, 134(17):174105, 2011. [PubMed: 21548671]

154. Giannakis Dimitrios, Slawinska Joanna, and Zhao Zhizhen. Spatiotemporal feature extraction with data-driven Koopman operators. In Storcheus Dmitry, Rostamizadeh Afshin, and Kumar Sanjiv, 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. Takens Floris. Detecting strange attractors in turbulence. In Dynamical Systems and Turbulence, volume 898 of Lecture Notes in Mathematics, pages 366–381. Springer, 1981.
156. Aeyels Dirk. Generic observability of differentiable systems. *SIAM J. Control Optim*, 19 (5):595–603, 1981.
157. Zwanzig Robert. Nonequilibrium Statistical Mechanics. Oxford University Press, 2001.
158. Łuczka J. Non-markovian stochastic processes: Colored noise. *Chaos*, 15(2):026107, 2005.
159. Camilloni Carlo and Pietrucci Fabio. Advanced simulation techniques for the thermodynamic and kinetic characterization of biological systems. *Adv. Phys.:*X, 3(1):1477531, 2018.
160. Hänggi Peter, Talkner Peter, and Borkovec Michal. Reaction-rate theory: Fifty years after Kramers. *Rev. Mod. Phys.*, 62:251–341, 1990. doi: 10.1103/RevModPhys.62.251.
161. Straub John E, Borkovec Michal, and Berne Bruce J. Calculation of dynamic friction on intramolecular degrees of freedom. *J. Phys. Chem*, 91(19):4995–4998, 1987.
162. Hummer Gerhard and Kevrekidis Ioannis G. Coarse molecular dynamics of a peptide fragment: Free energy, kinetics, and long-time dynamics computations. *J. Chem. Phys.*, 118(23):10762–10773, 2003.
163. Lange Oliver F and Grubmüller Helmut. Collective Langevin dynamics of conformational motions in proteins. *J. Chem. Phys.*, 124(21):214903, 2006. [PubMed: 16774438]
164. Hummer Gerhard. 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. Horenko Illia, Hartmann Carsten, Schütte Christof, and Noé Frank. Data-based parameter estimation of generalized multidimensional Langevin processes. *Phys. Rev. E*, 76(1): 016706, 2007.
166. Micheletti Cristian, Bussi Giovanni, and Laio Alessandro. Optimal Langevin modeling of out-of-equilibrium molecular dynamics simulations. *J. Chem. Phys.*, 129(7):074105, 2008. [PubMed: 19044758]
167. Darve Eric, Solomon Jose, and Kia Amirali. Computing generalized Langevin equations and generalized Fokker–Planck equations. *P. Natl. Acad. Sci. USA*, 106(27):10884–10889, 2009.
168. Legoll Frédéric and Lelièvre Tony. Effective dynamics using conditional expectations. *Nonlinearity*, 23(9):2131, 2010.
169. Schaudinnus Norbert, Bastian Björn, Hegger Rainer, and Stock Gerhard. Multidimensional Langevin modeling of nonoverdamped dynamics. *Phys. Rev. Lett.*, 115(5):050602, 2015. [PubMed: 26274405]
170. Meloni Roberto, Camilloni Carlo, and Tiana Guido. Properties of low-dimensional collective variables in the molecular dynamics of biopolymers. *Phys. Rev. E*, 94(5):052406, 2016. [PubMed: 27967023]
171. Lesnicki Dominika, Vuilleumier Rodolphe, Carof Antoine, and Rotenberg Benjamin. Molecular hydrodynamics from memory kernels. *Phys. Rev. Lett.*, 116(14):147804, 2016. [PubMed: 27104730]
172. Daldrop Jan O, Kappler Julian, Brünig Florian N, and Netz Roland R. Butane dihedral angle dynamics in water is dominated by internal friction. *P. Natl. Acad. Sci. USA*, 115(20): 5169–5174, 2018.
173. Zhang Qi, Bruji Jasna, and Vanden-Eijnden Eric. Reconstructing free energy profiles from nonequilibrium relaxation trajectories. *J. Stat. Phys.*, 144(2):344–366, 2011.
174. Pérez-Villa Andrea and Pietrucci Fabio. Free energy, friction, and mass profiles from short molecular dynamics trajectories. arXiv preprint, 1810.00713, 2018.
175. Samanta Amit, Chen Ming, Yu Tang-Qing, Tuckerman Mark, and E Weinan. Sampling saddle points on a free energy surface. *J. Chem. Phys.*, 140(16):164109, 2014. [PubMed: 24784255]

176. Pietrucci Fabio and Saitta Antonino Marco. 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. Pipolo Silvio, Salanne Mathieu, Ferlat Guillaume, Klotz Stefan, A Marco Saitta, and Fabio Pietrucci. Navigating at will on the water phase diagram. *Phys. Rev. Lett*, 119(24):245701, 2017. [PubMed: 29286747]
178. AlQuraishi Mohammed. End-to-end differentiable learning of protein structure. *Cell Systems*, 8(4):292–301.e3, 2019. ISSN 2405–4712. doi: 10.1016/j.cels.2019.03.006. [PubMed: 31005579]
179. Shaw David E., Maragakis Paul, Lindorff-Larsen Kresten, Piana Stefano, Dror Ron O., Eastwood Michael P., Bank Joseph A., Jumper John M., Salmon John K., Shan Yibing, and Wriggers Willy. Atomic-level characterization of the structural dynamics of proteins. *Science*, 330(6002):341–346, 2010. ISSN 0036–8075. doi: 10.1126/science.1187409. [PubMed: 20947758]
180. Krivov Sergei V. and Karplus Martin. 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.0800228105.
181. Sultan Mohammad M., Wayment-Steele Hannah K., and Pande Vijay S. 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. [PubMed: 29529369]
182. Brandt Simon, Sittel Florian, Ernst Matthias, and Stock Gerhard. Machine learning of biomolecular reaction coordinates. *J. Phys. Chem. Lett*, 9(9):2144–2150, 4 2018. [PubMed: 29630378]
183. Chen Shi, Wiewiora Rafal P., Meng Fanwang, Babault Nicolas, Ma Anqi, Yu Wenyu, Qian Kun, Hu Hao, Zou Hua, Wang Junyi, Fan Shijie, Blum Gil, Pittella-Silva Fabio, Beauchamp Kyle A., Tempel Wolfman, Jiang Hualiang, Chen Kaixian, Skene Robert, Zheng Y. George, Brown Peter J., Jin Jian, Luo Cheng, Chodera John D., and Luo Minkui. The dynamic conformational landscapes of the protein methyltransferase SETD8. *eLife*, 8: e45403, 2019. [PubMed: 31081496]
184. Trapl D, Horvacanin I, Mareska V, Ozelik F, Unal G, and Spiwok V Anncolvar: Approximation of Complex Collective Variables by Artificial Neural Networks for Analysis and Biasing of Molecular Simulations. *Front. Mol. Biosci*, 6:25, 2019. [PubMed: 31058167]
185. Henzler-Wildman Katherine and Kern Dorothee. Dynamic personalities of proteins. *Nature*, 450(7172):964–972, 12 2007. [PubMed: 18075575]
186. Schramm Vern L. Enzymatic transition states, transition-state analogs, dynamics, thermodynamics, and lifetimes. *Annu. Rev. Biochem*, 80(1):703–732, 7 2011. [PubMed: 21675920]
187. Lee GM and Craik CS Trapping moving targets with small molecules. *Science*, 324 (5924):213–215, 4 2009. [PubMed: 19359579]
188. Husic Brooke E., McGibbon Robert T., Sultan Mohammad M., and Pande Vijay S. Optimized parameter selection reveals trends in Markov state models for protein folding. *J. Chem. Phys*, 145(19):194103, 11 2016. [PubMed: 27875868]
189. Noé Frank and Clementi Cecilia. Collective variables for the study of long-time kinetics from molecular trajectories: Theory and methods. *Curr. Opin. Struc. Biol*, 43:141–147, 4 2017.
190. Bowman Gregory R. 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. Wodak Shoshana J., Paci Emanuele, Dokholyan Nikola V, Berezovsky Igo N, Horovitz Amnon, Li Jing, Hilser Vincent J., Bahar Ivet, Karanicolas John, Stock Gerhard, Hamm Peter, Stote Roland H., Jerome Eberhardt, Chebaro Yassmine, Dejaegere Annick, Marco Cecchini, Changeux Jean-Pierre, Bolhuis Peter G., Vreede Jocelyne, Faccioli Pietro, Orioli Simone, Ravasio Riccardo, Yan Le, Brito Carolina, Wyart Mathieu, Gkeka Paraskevi, Rivalta Ivan, Palermo Giulia, McCammon J. Andrew, Panecka-Hofman Joanna, Wade Rebecca C., Pizio Antonella Di, Niv Masha Y., Nussinov Ruth, Tsai Chung-Jung, Jang Hyunbum, Padhorny Dzmitry, Kozakov Dima, and McLeish Tom. Allostery in its many disguises: From theory to applications. *Structure*, 27(4):566–578, 2019. doi: 10.1016/j.str.2019.01.003. [PubMed: 30744993]

192. Fiorin Giacomo, Klein Michael L., and Héning Jérôme. 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. Man-Un Ung Peter, Rahman Rayees, and Schlessinger Avner. 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. [PubMed: 29861272]
194. Degiacomi Matteo T. 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. [PubMed: 31031199]
195. Óscar Díaz, Dalton James A.R., and Giraldo Jesús. Artificial intelligence: A novel approach for drug discovery. *Trends in Pharmacological Sciences*, 40(8):550–551, 2019. doi: doi: 10.1016/j.tips.2019.06.005. [PubMed: 31279568]
196. Bonomi Massimiliano, Branduardi Davide, Bussi Giovanni, Camilloni Carlo, Provasi Davide, Raiteri Paolo, Donadio Davide, Marinelli Fabrizio, Pietrucci Fabio, Broglia Ricardo A., and Parrinello Michele. Plumed: A portable plugin for free-energy calculations with molecular dynamics. *Computer Physics Communications*, 180(10):1961–1972, 2009. ISSN 0010–4655. doi: 10.1016/j.cpc.2009.05.011.
197. Case David A., Cheatham Thomas E. III, Darden Tom, Gohlke Holger, Luo Ray, Merz Kenneth M. Jr., Onufriev Alexey, Simmerling Carlos, Wang Bing, and Woods Robert J. The Amber biomolecular simulation programs. *Journal of Computational Chemistry*, 26 (16):1668–1688, 2005. doi: 10.1002/jcc.20290. [PubMed: 16200636]
198. Berendsen HJC, van der Spoel D, and van Drunen R GROMACS: A message-passing parallel molecular dynamics implementation. *Computer Physics Communications*, 91(1): 43–56, 1995. ISSN 0010–4655. doi: 10.1016/0010-4655(95)00042-E.
199. Phillips James C., Braud Rosemary, Wang Wei, Gumbart James, Tajkhorshid Emad, Villa Elizabeth, Chipot Christophe, Skeel Robert D., Laxmikant KalÃ©, and Klaus Schulten. Scalable molecular dynamics with NAMD. *J. Comput. Chem*, 26(16):1781–1802, 2005. doi: 10.1002/jcc.20289. [PubMed: 16222654]
200. Eastman P, Swails J, Chodera JD, McGibbon RT, Zhao Y, Beauchamp KA, Wang LP, Simmonett AC, Harrigan MP, Stern CD, Wiewiora RP, Brooks BR, and Pande VS OpenMM 7: Rapid development of high performance algorithms for molecular dynamics. *PLoS Comput. Biol*, 13(7):e1005659, 2017. [PubMed: 28746339]
201. Chollet François et al. Keras, 2015. <https://keras.io>.
202. Abadi Martín, Agarwal Ashish, Barham Paul, Brevdo Eugene, Chen Zhifeng, Citro Craig, Corrado Greg S., Davis Andy, Dean Jeffrey, Devin Matthieu, Ghemawat Sanjay, Goodfellow Ian, Harp Andrew, Irving Geoffrey, Isard Michael, Jia Yangqing, Jozefowicz Rafal, Kaiser Lukasz, Kudlur Manjunath, Levenberg Josh, Mané Dan, Monga Rajat, Moore Sherry, Murray Derek, Olah Chris, Schuster Mike, Shlens Jonathon, Steiner Benoit, Sutskever Ilya, Talwar Kunal, Tucker Paul, Vanhoucke Vincent, Vasudevan Vijay, Viégas Fernanda, Vinyals Oriol, Warden Pete, Wattenberg Martin, Wicke Martin, Yu Yuan, and Zheng Xiaoqiang. TensorFlow: Large-scale machine learning on heterogeneous systems, 2015. Software available from tensorflow.org.
203. Paszke Adam, Gross Sam, Chintala Soumith, Chanan Gregory, Yang Edward, DeVito Zachary, Lin Zeming, Desmaison Alban, Antiga Luca, and Lerer Adam. Automatic differentiation in pytorch. In NIPS 2017 Workshop on Autodiff, 2017.
204. Howard Jeremy et al. fastai, 2018. <https://github.com/fastai/fastai>.
205. Harrigan Matthew P., Sultan Mohammad M., Hernández Carlos X., Husic Brooke E., Eastman Peter, Schwantes Christian R., Beauchamp Kyle A., McGibbon Robert T., and Pande Vijay S. MSMBuiler: Statistical models for biomolecular dynamics. *Biophysical Journal*, 112(1):10–15, 2017. ISSN 0006–3495. doi: 10.1016/j.bpj.2016.10.042. [PubMed: 28076801]
206. McGibbon Robert T., Beauchamp Kyle A., Harrigan Matthew P., Klein Christoph, Swails Jason M., Hernández Carlos X., Schwantes Christian R., Wang Lee-Ping, Lane Thomas J., and Pande Vijay S. 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. [PubMed: 26488642]

207. Swenson David W. H., Prinz Jan-Hendrik, Noe Frank, Chodera John D., and Bolhuis Peter G. 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. [PubMed: 30359525]
208. Carry Jean-Christophe, Clerc François, Minoux Hervé, Schio Laurent, Mauger Jacques, Nair Anil, Parmantier Eric, Moigne Ronan Le, Delorme Cecile, Nicolas Jean-Paul, Krick Alain, Abecassis Pierre-Yves, Crocq-Stuerga Veronique, Pouzieux Stephanie, Delarbre Laure, Maignan Sebastien, Bertrand Thomas, Bjergarde Kirsten, Ma Nina, Lachaud Sylvette, Guizani Houlf, Lebel Rémi, Doerflinger Gilles, Monget Sylvie, Perron Sebastien, Gasse Francis, Angouillant-Boniface Odile, Filoche-Romme Bruno, Murer Michel, Gontier Sylvie, Prevost Celine, Monteiro Marie-Line, and Combeau Cecile. Sar156497, an exquisitely selective inhibitor of aurora kinases. *J. Med. Chem.*, 58(1):362–375, 2015. [PubMed: 25369539]
209. DrugBank. <https://www.drugbank.ca>, 2020.
210. Clinical Trials. <https://clinicaltrials.gov>, 2020.
211. Gkeka Paraskevi, Evangelidis Thomas, Pavlaki Maria, Lazani Vasiliki, Christoforidis Savvas, Agianian Bogos, and Cournia Zoe. 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. Gkeka Paraskevi, Papafotika Alexandra, Christoforidis Savvas, and Cournia Zoe. Exploring a non-ATP pocket for potential allosteric modulation of PI3K α . *J. Phys. Chem. B*, 119 (3):1002–1016, 2015. doi: 10.1021/jp506423e. [PubMed: 25299356]
213. Hub Jochen S. and de Groot Bert L. Detection of functional modes in protein dynamics. *PLOS Comput. Biol.*, 5(8):1–13, 08 2009. doi: 10.1371/journal.pcbi.1000480.
214. O’Connell James Philip, Porter John Robert, Rapecki Stephen Edward, Norman Timothy John, Warrelow Graham John, Arakaki Tracy Lynn, Burgin Alex Buntin, Pitt William Ross, Calmiano Mark Daniel, Schubert David Andreas, Lightwood Daniel John, and Wootton Rebecca Jayne. Novel TNF α structure for use in therapy, 2015. PCT / E P2015 / 074491.
215. Farimani Amir Barati, Feinberg Evan N., and Pande Vijay. Binding pathway of opiates to μ -opioid receptors revealed by machine learning. *Biophys. J.*, 114:62a–63a, 02 2018. doi: 10.1016/j.bpj.2017.11.390.
216. Feinberg Evan N., Farimani Amir Barati, Uprety Rajendra, Hunkele Amanda, Pasternak Gavril, Majumdar Susruta, and Pande Vijay. Machine learning harnesses molecular dynamics to discover new μ -opioid chemotypes. arXiv preprint, 1803.04479, 03 2018.
217. Schmidt Michael and Lipson Hod. Distilling free-form natural laws from experimental data. *Science*, 324(5923):81–85, 2009. [PubMed: 19342586]

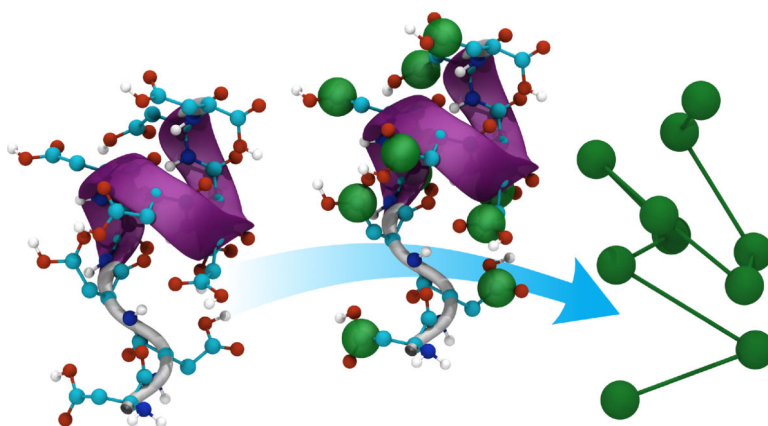


Fig. 1. Particle-based coarse-graining: high dimensional free energy surfaces (FES) can be extract from atomistic data and used as a basis for CG models (52, 53).

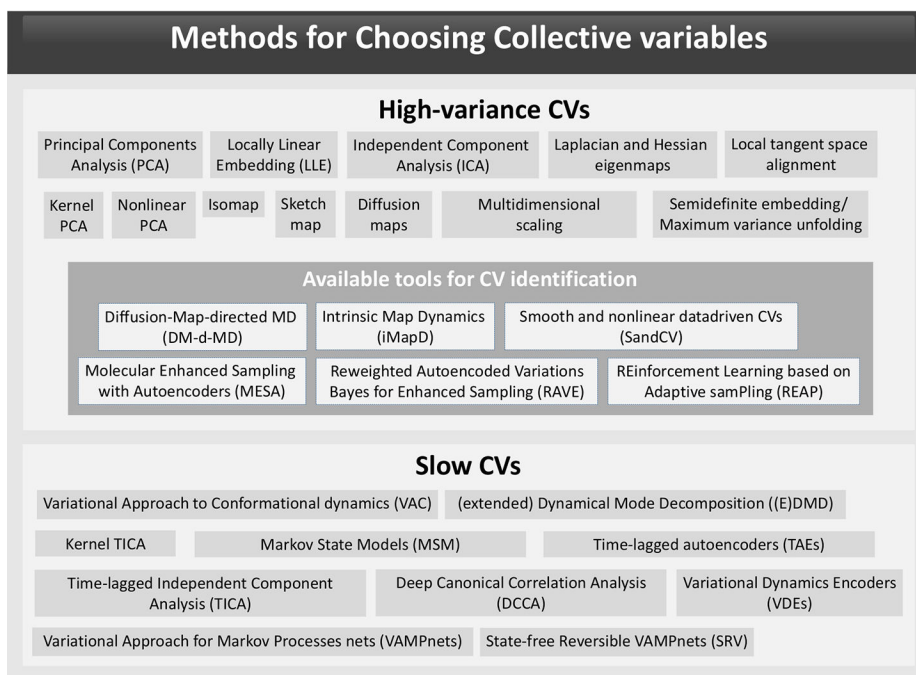


Fig. 2. Representative methods for CV identification. All related citations are in the main text.

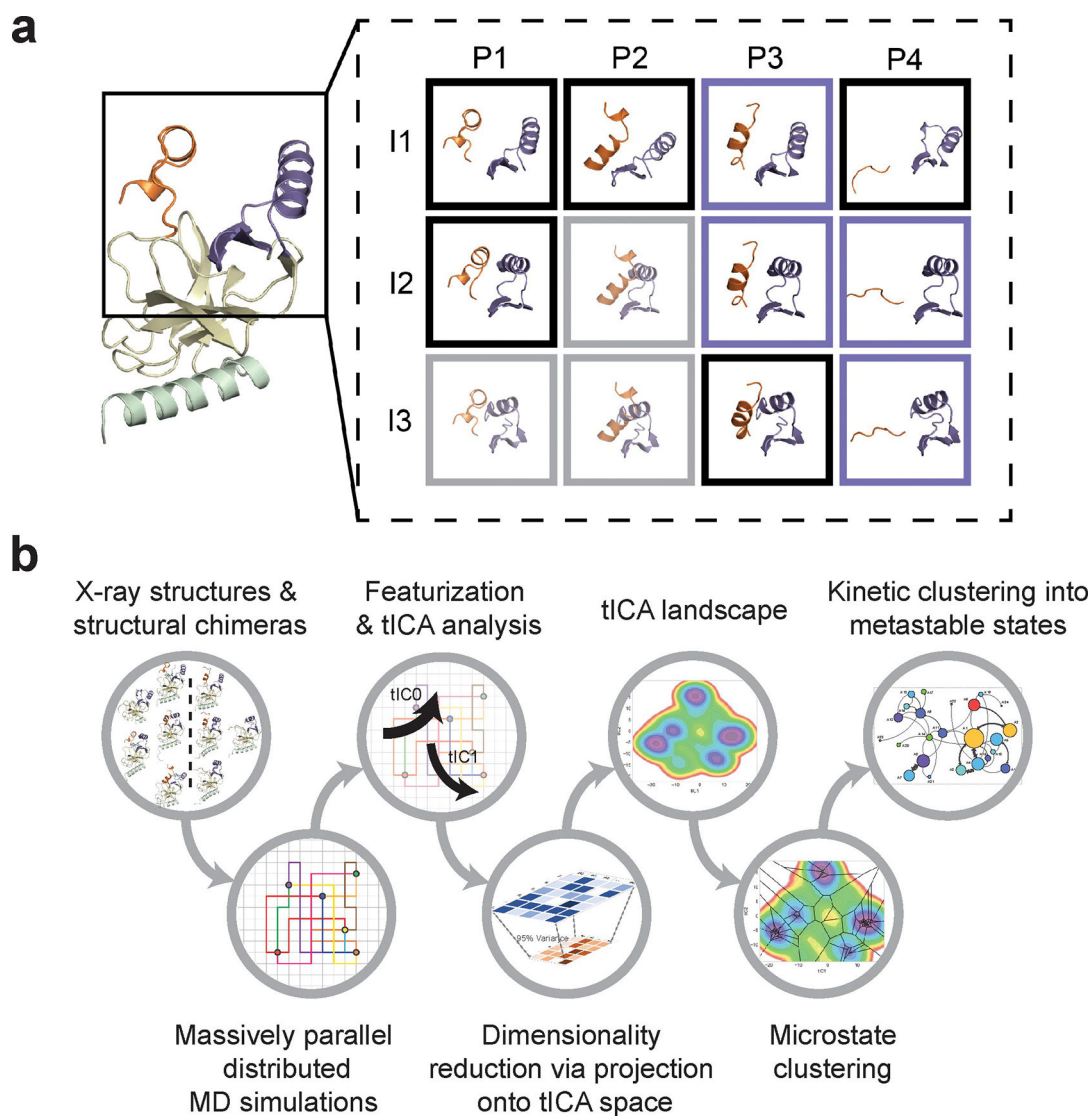


Fig. 3. Construction of conformational landscapes of apo- and SAM-bound SETD8 through diversely seeded, parallel molecular dynamics simulations and Markov state models. (a) Combinatorial construction of structural chimeras using crystallographically-derived conformations. (b) Workflow for dynamic conformational 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/>).

Table 1.

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

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

Author Manuscript

Author Manuscript

Author Manuscript

Author Manuscript