

Molecular Force Fields with Gradient-Domain Machine Learning: Construction and Application to Dynamics of Small Molecules with Coupled Cluster Forces

Huziel E. Sauceda,¹ Stefan Chmiela,² Igor Poltavsky,³ Klaus-Robert Müller,^{2,4,5,*} and Alexandre Tkatchenko^{3,†}

¹*Fritz-Haber-Institut der Max-Planck-Gesellschaft, 14195 Berlin, Germany*

²*Machine Learning Group, Technische Universität Berlin, 10587 Berlin, Germany*

³*Physics and Materials Science Research Unit, University of Luxembourg, L-1511 Luxembourg, Luxembourg*

⁴*Department of Brain and Cognitive Engineering,*

Korea University, Anam-dong, Seongbuk-gu, Seoul 02841, Korea

⁵*Max Planck Institute for Informatics, Stuhlsatzenhausweg, 66123 Saarbrücken, Germany*

(Dated: February 4, 2019)

We present the construction of molecular force fields for small molecules (less than 25 atoms) using the recently developed symmetrized gradient-domain machine learning (sGDML) approach [Chmiela *et al.*, Nat. Commun. **9**, 3887 (2018); Sci. Adv. **3**, e1603015 (2017)]. This approach is able to accurately reconstruct complex high-dimensional potential-energy surfaces from just a few 100s of molecular conformations extracted from *ab initio* molecular dynamics trajectories. The data efficiency of the sGDML approach implies that atomic forces for these conformations can be computed with high-level wavefunction-based approaches, such as the “gold standard” CCSD(T) method. We demonstrate that the flexible nature of the sGDML model recovers local and non-local electronic interactions (e.g. H-bonding, proton transfer, lone pairs, changes in hybridization states, steric repulsion and $n \rightarrow \pi^*$ interactions) without imposing any restriction on the nature of interatomic potentials. The analysis of sGDML molecular dynamics trajectories yields new qualitative insights into dynamics and spectroscopy of small molecules close to spectroscopic accuracy.

I. INTRODUCTION

Molecular force fields (FFs) constitute one of the most important tools in chemistry, biology, and materials modelling due to their remarkable value in understanding systems that range from small molecules, e.g. ethanol with 9 atoms, up to large proteins, aiding in the exploration and the discovery of new materials and drugs. Creating physically inspired and handcrafted interatomic potentials with parameters fitted to experimental data or quantum-mechanical calculations has been a common practice since the early works on molecular dynamics [1–4]. The complexity of creating reliable interatomic interaction models using prior physical knowledge led to the development of dedicated specialized FFs for different material classes, including tight-binding potentials for semiconductors and metals [5], Tersoff potential for covalent materials [6], polarizable FFs [7], the TIP n P FFs for water [8, 9], and a wide variety of biomolecular FFs such as AMBER, CHARMM, MMFF, and GROMOS that often lead to reliable results for folded protein structures under ambient conditions [10–13]. The wealth of available interatomic potentials illustrates the vast amount of fundamentally different material classes, exposing that treating different types of interactions (metallic bonding, covalent chemistry, hydrogen bonding, non-covalent interactions, etc.) in a unified and seamless fashion is a complex challenge for handcrafted mechanistic FFs.

To resolve some of these challenges, a number of recently developed machine learned (ML) FFs exploit the

redundant information contained in datasets of *ab initio* calculations or molecular dynamics trajectories to reconstruct the underlying potential energy surface (PES) without imposing any particular handcrafted analytical form for the interatomic potential.

In particular, a vast amount of work has been done in molecular representations [14–31], neural networks architecture development [32–40], data sampling [41–45] and inference methods [46–60], as well as software development [61–63] and explanation methods [64, 65].

Based on rigorous statistical learning theory [66, 67], machine learning provides a powerful and general framework for constructing force fields, since ML approaches can reconstruct complex high-dimensional objects with arbitrary accuracy, provided that sufficient data samples (molecular energies and atomic forces) are used for training. Obviously, the computational cost of evaluating ML FFs lies in between empirical FFs and *ab initio* reference calculations. In particular, the sGDML approach employed here is 5-10 orders of magnitude faster than *ab initio* calculations and 2-3 orders of magnitude slower than classical FFs, this depends of course on the molecular system. For example, in the case of a CCSD(T)/cc-pVTZ calculation, the sGDML model can be up to 10^7 and 10^9 times faster for malondialdehyde and aspirin, respectively[61].

In the broadest sense, the challenge of accurately learning FFs is currently being addressed using two methods: Neural Networks (NN) [32–35, 37, 38, 68–70] and kernel-based models [17, 19, 23, 26, 41, 42, 44, 46, 48, 58–60]. Both approaches can be constructed to employ energy and/or force information. Learning forces is advantageous for several reasons: (i) FFs reconstruction in the force domain yield smoother PESs, eliminating artifacts

* klaus-robert.mueller@tu-berlin.de

† alexandre.tkatchenko@uni.lu

due to somewhat conflicting requirement of simultaneously reproducing accurate energies and forces [37, 61], the inherent uncertainty of the learning process, and using biased models that introduce unphysical approximations, e.g. atomic partitioning of the energy, (ii) obtaining energies from force models tends to diminish the noise in the prediction as a result of the integral operator, contrasting the behaviour of the forces generated by the gradient operator on energy models, and (iii) force models require smaller amounts of reference calculations to reach a desired accuracy [59]. Such data efficiency arises not only due to the fact that each force data point carries $3N$ components (where N is the number of atoms) per reference calculation, but also because those components are orthogonal, thus providing complete information about the immediate local environment [71].

Both NN and kernel-based methods can achieve formally any desired accuracy of predictions whenever a sufficiently large amount of training data is available. In contrast, when only 100s of data points are available, as in the case of high-level *ab initio* data, the kernel methods usually offer a better reconstruction efficiency (with a unique and well-defined solution) as they make greater use of prior information. Finally, it is important to emphasize the mandatory requirement of generating conservative ML-FFs, i.e. $\mathbf{F} = -\nabla E$, to guarantee stable simulations.

The symmetric Gradient Domain Machine Learning (sGDML) FF [60] retains all the advantages of kernel-based ML models which directly learn forces. In fact, training the sGDML model solely using forces, besides the availability of molecular energies, improves the learning process. Given that there is evidence that combining energies and forces in the loss function degrades the quality of the force prediction [37, 61]. The robustness of the method is explained by the fact that all atomic interactions are modelled globally, without resorting to an inherently non-unique partitioning into atom-wise, pairwise or many-body contributions. In the sGDML model: (i) Each FF model is explicitly constrained to be energy conserving, and (ii) The model complexity is further reduced through the incorporation of molecular symmetries (i.e. rigid and fluxional) that are automatically extracted from the reference dataset. All these important properties contribute to the ability of sGDML to reconstruct complex PES for molecules of intermediate size from modest amounts of reference data, an unfeasible task for non-dedicated molecular FFs. In particular, the sGDML model enables the reconstruction of CCSD(T)-quality FFs from a limited amount (~ 100 s) of reference molecular configurations [60].

In this article, we analyze some of the relevant quantum effects captured by the sGDML model while reconstructing the PES of small molecules. First, in section II, we give a short introduction to the GDML framework and its symmetrized version, the sGDML model. In section III a discussion regarding the advantages of gradient domain-based FFs is presented, as well as an analysis of

dedicated vs. transferable FFs. Then, in section IV, we describe some of the quantum interactions described by the sGDML’s reconstructed PES. In particular, we focus on three ubiquitous phenomena of general interest: section IV-A) lone-pairs and electrostatic interactions, section IV-B) intramolecular hydrogen bonds and proton transfer, and section IV-C) changes in atomic hybridization state and $n \rightarrow \pi^*$ interactions [72, 73]. Note that a qualitative description of these complex interactions by regular FF would require highly specialized models, while the sGDML model captures every interaction encoded in $-\mathbf{F}_i = \langle \Psi^* | \partial \mathcal{H} / \partial \mathbf{x}_i | \Psi \rangle$ with high accuracy. In the last section, we summarize our findings.

II. sGDML MODEL

The sGDML model enables the direct efficient construction of dedicated FFs for flexible molecules from high-level *ab initio* calculations. Unlike traditional FFs, it imposes no hypothesized analytical interaction models and thus in principle can model any physical interaction. Compared to other machine learning approaches, sGDML achieves high data efficiency through the incorporation of spatial and temporal physical symmetries. Global spatial symmetries include rotational and translational invariance of the energy, in addition to rigid and fluxional symmetries, which are recovered and enforced in an automatic data-driven manner. The homogeneity of time implies energy conservation, property that is enforced by learning in the gradient domain using as prior an analytically-integrable covariance function.

The latter is introduced as a linear operator constraint, by modeling the FF as the transformation of an underlying energy model. In particular, we train the gradient of a kernel ridge estimator on force labels \mathbf{F} , which – by construction – yields energy-conserving FFs that can be integrated to obtain the corresponding Born-Oppenheimer (BO) global potential-energy surface (PES) V_{BO} [59, 60]. Practically, this is achieved via the use of the Hessian matrix of a kernel κ as the covariance structure to solve the normal equation of the ridge estimator in the gradient domain [59, 60]:

$$(\mathbf{K}_{\text{Hess}(\kappa)} + \lambda \mathbb{I}) \vec{\alpha} = \nabla V_{BO} = -\mathbf{F}, \quad (1)$$

where $\mathbf{K}_{\text{Hess}(\kappa)}$ is the kernel matrix, λ is the regularization parameter, \mathbb{I} is the identity matrix, and $\vec{\alpha}$ are the parameter-vectors to train.

The sGDML model imposes additional permutational symmetry constraints on $\mathbf{K}_{\text{Hess}(\kappa)}$ to take advantage of PES redundancies due to rigid space group and fluxional symmetries [60]. Typically, extracting those symmetries requires chemical and physical intuition about the system at hand, e.g. rotational barriers, which is impractical in a ML setting. Through a data-driven multi-partite matching approach (see Fig. 1), we automate the discovery of permutation matrices $\mathbf{P}(\tau)$ corresponding to

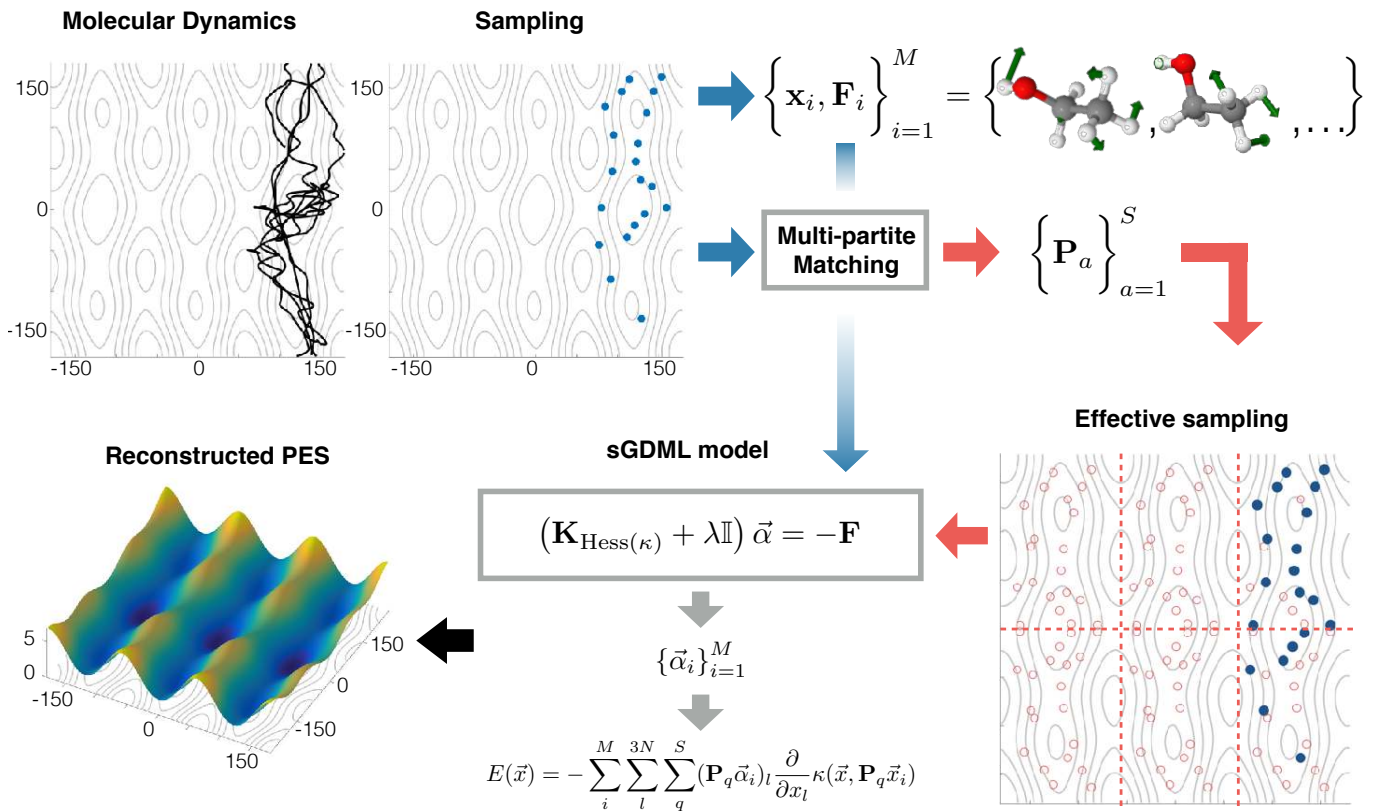


FIG. 1. Construction of the sGDML model. (1) The data used for training, $\{\mathbf{x}_i, \mathbf{F}_i\}_{i=1}^M$, is generated by random sampling of molecular dynamics trajectories (blue dots). The force on each atom is represented by a green arrow. (2) From the training set the permutational set of symmetries, $\{\mathbf{P}_a\}_{a=1}^S$, are computed by the multi-partite matching approach. This effectively enhances the size of the training set by a factor S . (3) The force field is trained by solving the linear system for the parameters $\{\alpha_j\}$. The reconstructed potential-energy surface is obtained by analytically integrating the force model.

the index permutation τ from molecular dynamics simulations by realizing the assignment between adjacency matrices $(\mathbf{A})_{ij} = \|\vec{r}_i - \vec{r}_j\|$ of molecular graph pairs G and H in different energy states,

$$\arg \min_{\tau} \mathcal{L}(\tau) = \|\mathbf{P}(\tau)\mathbf{A}_G\mathbf{P}(\tau)^\top - \mathbf{A}_H\|^2. \quad (2)$$

The resulting approximate pairwise matchings are subsequently synchronized using transitivity within the training set as the consistency criterion. A particular advantage of our solution is its ability to limit the symmetry recovery to energetically feasible permutational configurations, given that unfeasible permutation, e.g. the permutation of two random atoms, would not contribute any valuable information the symmetrized kernel and should not be considered. This severely reduces computational efforts in evaluating the model. Finally, the FF estimator trained on M reference geometries, with $3N$ partial derivatives and S symmetry transformations each, takes the form

$$\hat{\mathbf{f}}_{\mathbf{F}}(\vec{x}) = \sum_i^M \sum_l^{3N} \sum_q^S (\mathbf{P}_q \vec{\alpha}_i)_l \frac{\partial}{\partial x_l} \nabla \kappa(\vec{x}, \mathbf{P}_q \vec{x}_i). \quad (3)$$

The corresponding energy predictor is obtained by simply integrating $\hat{\mathbf{f}}_{\mathbf{F}}$ with respect to the Cartesian coordinates,

$$-\hat{f}_E(\vec{x}) = \int \hat{\mathbf{f}}_{\mathbf{F}} \cdot d\mathbf{x} = \sum_i^M \sum_l^{3N} \sum_q^S (\mathbf{P}_q \vec{\alpha}_i)_l \frac{\partial}{\partial x_l} \kappa(\vec{x}, \mathbf{P}_q \vec{x}_i). \quad (4)$$

Due to linearity of integration, the expression for the energy predictor is identical up to second derivative operator on the kernel function (see Fig. 1). Figure 1 gives a general perspective of the sGDML model by summarizing the training process, from sampling the MD trajectory and extracting the permutational symmetries to solving the linear system and reconstructing the embedded PES in the data.

The addition of spatial, temporal and permutational symmetry constraints leads to a gain in data efficiency of more than two orders of magnitude [60]. Recently, we have systematically demonstrated that sGDML models trained on only few 100s of reference structures reconstruct molecular PESs with a mean average error of less than $0.06 \text{ kcal mol}^{-1}$ for small molecules with up to 15 atoms and less than $0.16 \text{ kcal mol}^{-1}$ for molecules as com-

plex as aspirin, paracetamol, and azobenzene [60] (See Table I, Tables S1 and S2). Hence, the explicit symmetrization incorporated in the GDML framework [59] results in robust learning models with the ability to preserve the complex subtleties encoded in the reference data.

The sGDML models for each molecule studied in this article were initially trained on DFT data at the generalized gradient approximation (GGA) level of theory with Perdew-Burke-Ernzerhof (PBE) [74] exchange-correlation functional and the Tkatchenko-Scheffler (TS) method [75] to account for van der Waals interactions. The training dataset was created by subsampling MD trajectories at constant temperature (500K) using the FHI-aims package [76]. In the case of keto-MDA, enol-MDA and ethanol we recomputed the training configurations using all-electron CCSD(T), while in the case of Aspirin we used all-electron CCSD [77–79] (see Supporting Information for further details).

III. COMPARISON OF sGDML TO OTHER ML-FF APPROACHES

A. Force vs. energy model

The unique approach used by the sGDML model contrasts other models that first develop an energy function and then get the forces by analytic differentiation [1–13, 17, 19, 23, 31–35, 37, 38, 46, 48, 68–70]. This is represented in the next diagram:

	Trained	Derived
sGDML :	$\hat{\mathbf{f}}_{\mathbf{F}}$	$\longrightarrow \hat{f}_E = -\int \hat{\mathbf{f}}_{\mathbf{F}} \cdot d\mathbf{x}$
E-ML :	\hat{f}_E	$\longrightarrow \hat{\mathbf{f}}_{\mathbf{F}} = -\nabla \hat{f}_E$

where the trained predictors and their post-training derived energy and forces are presented for the sGDML and an energy ML (E-ML) model, respectively. An interesting advantage of the sGDML over E-ML models is how the training error propagates to the derived quantities. Lets assume that the prediction errors associated with the models $\hat{\mathbf{f}}_{\mathbf{F}}$ and \hat{f}_E are γ_F and γ_E , respectively. Then, from the discrete approximation of the integral and the derivative operator, we obtain that the error in the derived energy, $-\int \hat{\mathbf{f}}_{\mathbf{F}} \cdot d\mathbf{x}$, is attenuated and given by $\sim \gamma_F \Delta x$ while the error in the derived forces, $-\nabla \hat{f}_E$, is amplified and given by $\sim \gamma_E / \Delta x$ (see Supporting Information for further details). A direct implication of these results is that, as a whole, FFs based on E-ML are potentially less stable than gradient based FFs. Empirical evidence supporting these results as well as a proof from signal processing theory were published in the original GDML paper [59].

Regarding the data efficiency of the sGDML, there is solid evidence from Gaussian processes (GPs) that learning linearizations of a function, e.g. gradients, is more

informative than learning single points [71]. Such data efficiency has been systematically shown in the GDML framework [59–61]. There is empirical evidence that more than $3N$ training data points would be needed in an E-ML model per each sample used in GDML to reach similar force accuracy [59]. Therefore, allowing to train molecular FFs using data from very accurate but computationally expensive reference methods, e.g. CCSD(T).

B. Performance of using forces vs. forces+energies for training

In the process of generating ML-FFs, the nature of the model, E-ML or gradient domain model, gives prior information regarding the problem to solve. This in the context of GPs would be equivalent to narrow the space of possible solutions. Then, a loss function is introduced in order to train the model by finding the best set of parameters that minimize such function (here presented without the regularization part):

Model	Loss function
sGDML :	$loss_F = \sum_{i=1}^M \ \hat{\mathbf{f}}_{\mathbf{F}}(\vec{x}_i) - \mathbf{F}_i\ ^2$
E-ML :	$loss_E = \sum_{i=1}^M \ \hat{f}_E(\vec{x}_i) - E_i\ ^2$

Using these loss functions would, in principle, give an optimal fitting respectively. There is the idea that such loss functions can be complemented by adding energy or force constraints, as they are often available in the reference data. In fact, several related works optimize a hybrid squared loss function of the form [37]:

$$loss = \sum_{i=1}^M \left\{ \|\hat{\mathbf{f}}_{\mathbf{F}}(\vec{x}_i) - \mathbf{F}_i\|^2 + \eta \|\hat{f}_E(\vec{x}_i) - E_i\|^2 \right\}, \quad (5)$$

where η is a linear trade-off hyper-parameter which absorbs the differences in units and weights the force and energy contribution. By training a model using this loss function a somewhat conflicting optimization race between energies and forces is introduced. Clear empirical evidence of this issue has been reported for NN-FFs [37] and for the sGDML+E model [61] where in both cases the quality of the forces degrades by introducing energy constraints.

C. Transferable and non-transferable models

Recently, the idea of transferable or across chemical compound space ML-FFs has been under discussion but due to its complexity less progress has been achieved compared to dedicated ML-FFs. Transferable models can generate qualitatively good predictions simultaneously for different molecular systems [55, 80–82], but

clearly cannot offer reliable results for PES reconstruction where energy prediction errors are often much larger than 1 kcal mol⁻¹ [55]. On the other hand, the accuracy achieved by state-of-the-art dedicated ML-FFs can even reach couple of *calories* per mole in energy predictions using only a few hundreds reference calculations for training [60]. In contrast, the above mentioned transferable model required ~ 13 million data configurations for training. Furthermore, any gradient field generated from these models is, at the moment, not reliable because of such energy errors, which makes prohibitive to accurately capture physical interaction [55]. From this discussion it is apparent that transferable models cannot be used to study dynamical properties of molecules, a task easily accomplished by dedicated FFs.

IV. MOLECULAR POTENTIAL-ENERGY SURFACES

In the framework of the BO approximation, V_{BO} contains all the information necessary to describe the dynamics of a molecular system. All electronic quantum interactions are encoded in V_{BO} , but in practice it is not possible to expand V_{BO} in different energetic contributions such as hydrogen bonding, electrostatics, dispersion interactions or other electronic effects. Therefore, the intricate form of V_{BO} resulting from an interplay between different quantum interactions, should be preserved in the reconstruction process. We will now demonstrate how sGDML is able to describe many complex features contained in the quantum-chemical conformational data.

TABLE I. Accuracy of total energies for sGDML@DFT (using PBE+TS functional) and sGDML@CCSD(T) models on various molecular dynamics datasets. Energy errors are in kcal mol⁻¹. These results, with the exception of enol-MDA, were originally published in Ref. [60]. All the models were trained using atomic forces for 1000 molecular conformations.

Dataset	Energy Prediction			
	DFT		CCSD(T)	
Molecule	MAE	RMSE	MAE	RMSE
keto-MDA	0.10	0.13	0.06	0.08
Ethanol	0.07	0.09	0.05	0.07
enol-MDA	0.07	0.09	0.07	0.07
Aspirin	0.19	0.25	0.16 ^a	0.21 ^a

^a CCSD

In practice, these important features or interactions (e.g. energy barriers or H-bond interactions) often come in the form of subtle variations in the V_{BO} of less than 0.1 kcal mol⁻¹, which is one order of magnitude lower than so-called chemical accuracy [60]. For example, the relative stability of *trans* and *gauche* conformers of ethanol is within 0.1 kcal mol⁻¹. Any model with an expected

error above that threshold runs the risk of misrepresenting or even inverting this subtle energy difference, which will lead to incorrect occupation probabilities and hence qualitatively wrong dynamical properties. The sGDML model has been shown to satisfy the stringent accuracy requirement of 0.1–0.2 kcal mol⁻¹ for molecules with up to 15 atoms [60]. Moreover, we found that using coupled-cluster reference data not only generates a more accurate description of the quantum system, but also improves the learning errors as shown in Table I. In this section, we exemplify the accuracy and insights obtained with sGDML with ubiquitous and challenging features of general interest in chemical physics: electron lone pairs, electrostatic interactions, intramolecular hydrogen bonds, proton tunneling effect, and other electronic effects (e.g. steric repulsion, change in the bond nature, and bonding–antibonding orbital interaction). Figure 2 shows an overview of the different types of molecules and their reconstructed PES chosen to highlight the mentioned effects in this study.

A. Electrostatic interactions and electron lone pairs

First, we focus our attention on electrostatic interactions, including atom–atom and lone-pair–atom interaction. Here, the concept of electron lone pairs play a central role; these are ubiquitous molecular features responsible for a wide variety of physical and chemical phenomena. Lone pairs are valence electrons of an atom that are not shared with any other atom in a molecule. Some examples of atoms in molecules that often present lone pairs are nitrogen and oxygen. To illustrate the interactions induced by lone pairs, we will use the keto tautomer of malondialdehyde (keto-MDA) and ethanol molecule shown in the first two rows of Fig. 2. These fluxional molecules have complex PES with a rich variety of physical phenomena (e.g. electrostatics and steric repulsion) for which the reconstruction process is not a trivial task (see Fig. 3).

1. Oxygen–oxygen atom repulsion in keto-MDA

To illustrate the interatomic repulsion interaction we use the keto-MDA molecule, whose PES complexity is evident despite its small size (see Fig. 3-A). For example, the PES contains flat regions, which correspond to the global minima of the molecule (depicted in dark blue in Fig. 3-A), but also display intricate pathways to move to local minima. Also, one can notice the sudden energy increase when the two oxygen atoms are in the closest configuration (structure (1) in Fig. 3-A). As mentioned before, the PES is the result of many complex interactions but certainly there are parts of the PES which can be mainly ascribed to a particular phenomenon. This is the case for the yellow region in Fig. 3-A, where the


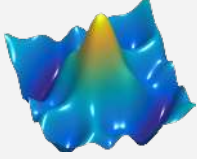
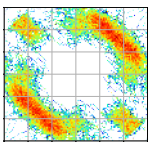
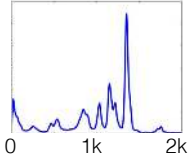
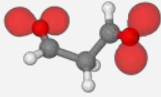

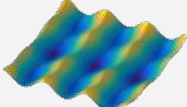
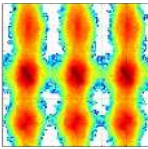
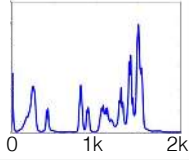
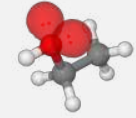

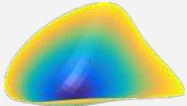
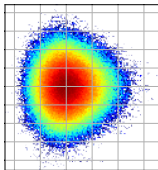
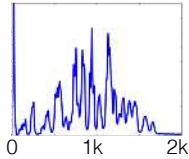
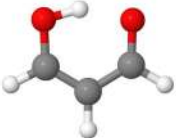
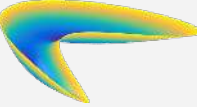
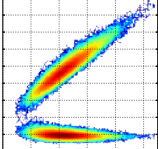
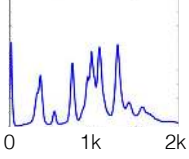
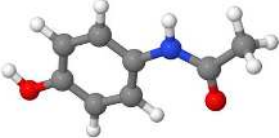
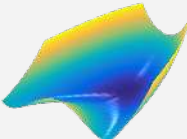
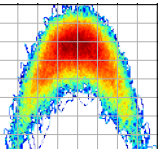
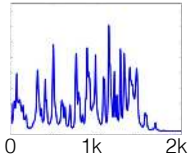
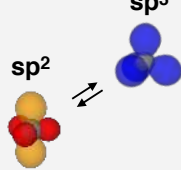
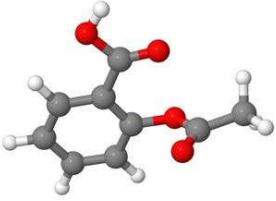
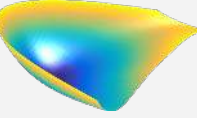
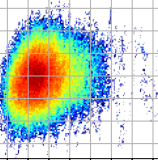
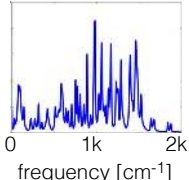
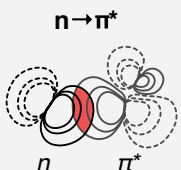
Molecules	Properties				
	PES	FES	H-bond	Spectrum	electronic effects
 keto Malondialdehyde					
 Ethanol					
 Salicylic acid			O—H...O		
 enol Malondialdehyde			O—H...O O...H—O		
 Paracetamol					
 Aspirin				 frequency [cm ⁻¹]	

FIG. 2. Molecules under study and their properties. From left to right: List of molecules and their molecular structure, potential-energy surface along two relevant torsional degrees of freedom, free energy surface at 300 K, type of intramolecular hydrogen bonds (if applicable), vibrational spectrum at 300 K, and type of electronic effect to study (if applicable). The free energy scale from lowest (red) to higher (blue) is in $k_B T$. The last column shows the electron lone pairs in keto-MDA and ethanol, $sp^2 \rightarrow sp^3$ hybridization transition present in paracetamol and $n \rightarrow \pi^*$ interaction in aspirin[72].

closeness between the two oxygen atoms suggests that the steep increase in the energy could be primarily attributed to the electrostatic repulsion between the lone pairs in each atom. Additionally, we know that electron lone pair clouds have large spatial extent compared to shared electrons, therefore steric effects caused by elec-

tron cloud overlap could also be playing an important role in this region due to the close proximity between the two oxygen atoms ($r_{OO} \sim 2.6 \text{ \AA}$). From these two interactions, only the electrostatic contribution is roughly incorporated in regular FFs as constant point charges located on each atom. This greatly constrains their flex-

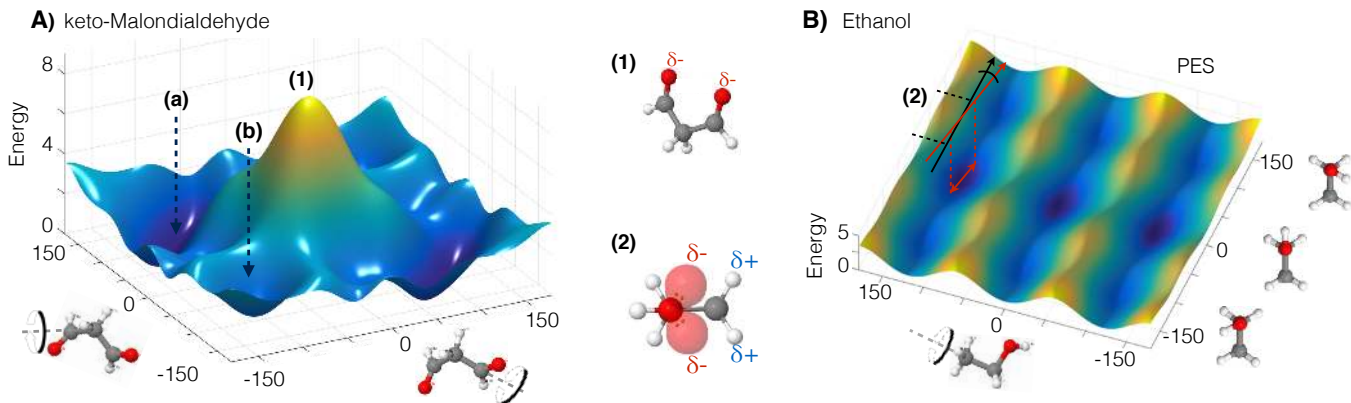


FIG. 3. Features of the PES. A) PES for keto-MDA. The structure (1) leads to a steep increase in energy due to the close distance between two negatively charged oxygen atoms. The regions (a) and (b) are the global and a local minima of keto-MDA. B) Ethanol’s PES. Structure (2) shows the effect of oxygen’s lone pair and the partial positive charges in methyl’s hydrogen atoms and their coupling is represented by a red arrow in the PES of ethanol. Both PES were predicted using sGDML@CCSD(T) models [60].

ibility and reliability to describe complex interactions. Nonetheless, systematic studies of such interactions using the sGDML model could spawn new ideas regarding their integration into regular FFs and ultimately increase the predictive power of empirical FFs.

2. Electron lone pairs in ethanol

A particularly interesting case of strong effects of electron lone pairs is the ethanol molecule, where the lone pairs of the oxygen atom interact with the partially positive hydrogen atoms of the methyl group (structure (2) in Fig. 3). This molecule has two rotors – the hydroxyl and methyl groups – as its main degrees of freedom. The PES for ethanol Fig. 3-B exhibits a very subtle quasi-linear dependence between the dihedral angles of the methyl and hydroxyl functional groups in the *trans* configuration (angle zero of the hydroxyl group in Fig. 3-B) as shown by the red arrows in Fig. 3-B. Such coupling is evident when analyzing the normal modes for this configuration, where the lowest two vibrations correspond to the aforementioned coupled motion (see also Fig. 4-C). The origin of this coupling can be understood by the electrostatic attraction between the lone pairs in the oxygen atom and the partially positively charged hydrogen atoms in the methyl rotor (structure (2) in Fig. 3). The correct description of this phenomenon within the reference data is crucial to obtain accurate physical properties of any molecular system with lone pairs [60]. It is important to stress that an accurate and general description of lone pairs is a characteristic that goes beyond the capabilities of regular FFs [60]. The handmade FFs that attempt to include lone pairs introduce explicit extra point charges [83], which results in highly specialized models [84].

3. Dynamics of coupled rotors in ethanol

The coupling between the hydroxyl and methyl rotors in ethanol manifests in the two lowest normal modes (1) and (2) shown in Fig. 4-C. The normal mode (1) corresponds to the direction indicated by the red arrow in the PES (Fig. 3-B), while (2) moves perpendicularly. Here we analyze the dynamical implications of the coupling between lone pairs and the methyl rotor. By performing molecular dynamics simulations at 300K using the NVE and NVT ensembles, we have obtained different probability distributions and therefore different free energy surfaces (FES) as shown in Fig. 4-A. In both cases the FES considerably differs from its underlying PES in the *trans* region. The FES_{NVE} develops a deep minimum which traps the system into a state that boosts the occupation of the lowest vibrational mode, illustrated by (1) in Fig. 4-C. In contrast, FES_{NVT} reverts the direction of the coupling as depicted in Fig. 4-A-(2), which in general promotes the occupation of both vibrational modes (1) and (2) in Fig. 4-C. In both cases, the FES shows an interesting and contrasting behavior that highlights the importance of the coupling between the two rotors in ethanol. A possible explanation for the difference between the two ensembles is that the NVT distributes the energy between all the molecular degrees of freedom in a more efficient way compared to NVE. The NVE ensemble relies only on the anharmonicities of the PES to redistribute the energy. Furthermore, the strong coupling between the two rotors suppresses the energy redistribution between vibrational normal modes in NVE. It is clear that only FFs with an explicit description of lone pairs can show this coupling.

As shown here for ethanol and keto-MDA, the interactions involving electron lone pairs (e.g. electrostatic interactions, steric repulsion and $n \rightarrow \pi^*$ interactions) encoded in the PES play a fundamental role in the dynam-

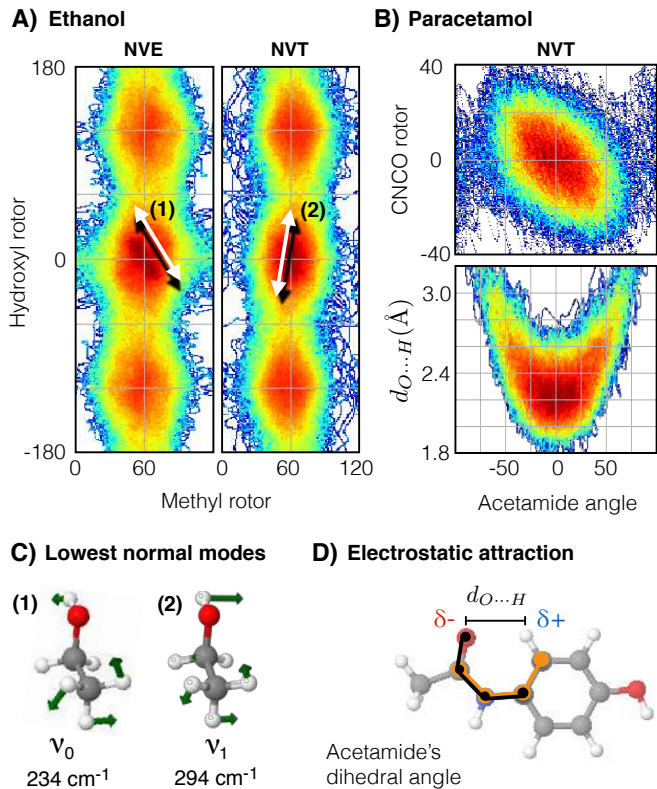


FIG. 4. Free energy for A) ethanol and B) paracetamol molecules at 300K using $F(T; x, y) = -k_B T \ln[P(x, y)]$ where $P(x, y)$ is the sampling obtained from classical molecular dynamics. For ethanol two different ensembles are presented: NVE and NVT. C) Two lowest vibrational normal modes and their frequencies for *trans* ethanol. The normal modes (1) and (2) are represented in the free energies A-NVE and A-NVT, respectively. D) Qualitative representation of the partial charges in paracetamol. The dihedral angle of the acetamide group respect to the benzene ring is represented in orange and the inner rotational degree of freedom CNCO is represented in black.

ics of the molecule and defining the free energy (Fig. 2). This has direct implications on its thermodynamics and spectroscopic properties given that a different sampling of the PES re-weights and shifts the peaks in the vibrational spectrum. Therefore, since electron lone pairs participate in many other more complex interactions, their appropriate description is the first steps to generate reliable force fields.

B. Intramolecular H-bond and proton transfer

Another complex phenomenon accurately captured with the sGDML method is hydrogen bonding (H-bond). The *intramolecular* H-bond is a subtle interaction, which dictates the dynamical behavior of many molecules [85]. It is responsible for very important molecular features such as the molecular structure and vibrational spectra,

which result in macroscopic properties e.g. solubility and permeability [86]. Here, we will study two different types of H-bonds: standard donor-acceptor H-bond and the symmetric H-bond present in salicylic acid and the enol form of malondialdehyde (enol-MDA), respectively (Fig. 5). The symmetric H-bond allows proton tunneling to occur due to thermal fluctuations assisted by quantum nuclear effects to overcome the energetic barrier. In the case of a standard donor-acceptor H-bond, the proton is fixed to the proton donor (PD) while its dynamical behavior is strongly affected by the electron lone pair belonging to a neighboring atom (proton acceptor, PA). These two kinds of intramolecular H-bond appear very often in molecules and their presence can drastically change the physical properties of any molecule, as we show in this section for salicylic acid and enol-MDA molecules.

1. Intramolecular H-bond

We start by analyzing the salicylic acid molecule (Fig. 5-A). This molecule presents a standard donor-acceptor kind of H-bond between the hydroxyl and carboxylic acid groups. From the schematic representation in Fig. 5-A, we can see that the effect of the H-bond in this molecule consists in allowing the proton to stretch from the PD oxygen towards the PA oxygen. The middle point (transition state) between PD and PA is represented by a red plane in salicylic acid's PES in Fig. 5-A bottom. The energy necessary to reach this point is ~ 10 kcal mol $^{-1}$, barely accessible at room temperature. We would also like to highlight the narrow structure of the reduced PES in the transition state (red plane in Fig. 5-A bottom), which gives us an idea of how directional the H-bond is. This directionality of the interaction changes the dynamics of the participating functional groups, which results in a characteristic red-shift in the stretching frequency of O-H induced by the H-bond [87-89]. From a vibrational normal mode analysis on the sGDML reconstructed PES (see Fig. 6), we observe the red-shift of the O-H stretching frequency in the participating hydroxyl group. Furthermore, the H-bond also creates a blue shift in normal modes perpendicular to the H-bond (see Fig. 6), which is a direct evidence of a O-H...O bond. The proper description of these molecular features will be directly displayed in spectroscopic properties such as IR and Raman spectrum which is often used to characterize molecular structures and their interactions.

2. Proton transfer

The second type of H-bond we analyze is the symmetric H-bond in enol-MDA, which exhibits a symmetric double-well reduced PES as schematically represented in Fig. 5-B. The energetic barrier separating the two minima is ~ 4 kcal mol $^{-1}$ which occurs when the interatomic

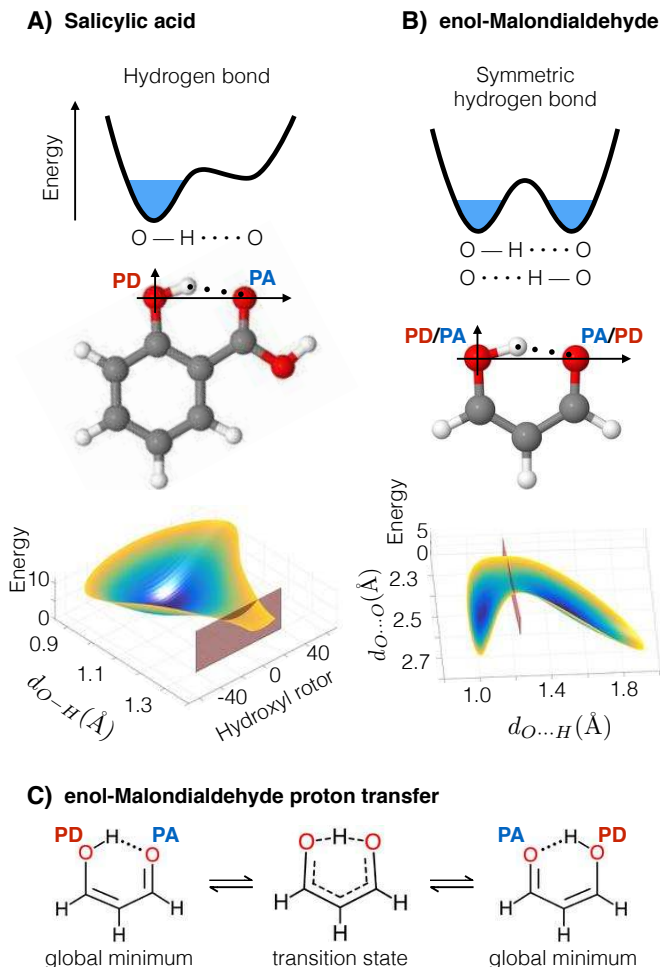


FIG. 5. Intramolecular hydrogen bond in A) salicylic acid and B) enol-Malondialdehyde. The top row shows schematically the type of H-bond in each case. In the middle row are the molecular structures and their respective proton reaction directions (from PD to PA). In the case of enol-MDA the PES is symmetric, then PD and PA are interchangeable. In the bottom row the PESs are shown where the red plane indicates the transition state of the proton. The energy at the transition states are ~ 10 kcal mol $^{-1}$ for salicylic acid and ~ 4 kcal mol $^{-1}$ for enol-MDA. In the transition state, the enol-MDA molecule has a C_{2v} point symmetry, which the sGDML model exploits to increase the reconstruction accuracy.

distance between oxygen atoms is $d_{O...O} = 2.38$ Å, this allows proton transfer between the two oxygen atoms even at room temperature. The transition state in the PES is shown by the red plane in Fig. 5-B bottom.

The energetic barrier for the proton transfer has two possible contributions: Electron density rearrangement and quasi-aromaticity. From the schematic representation of the enol-MDA in Fig. 5-C, we see that going from the global minimum to the transition state entails a redistribution of the electron density. Therefore, the redistribution of π electrons in the molecule induces an energetic penalty [90], which contributes to the generation of a higher energy barrier. There is also evidence

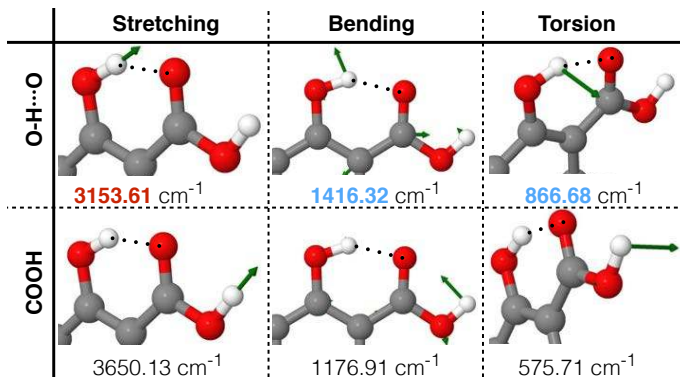


FIG. 6. Intramolecular H-bond in salicylic acid and its implications in the vibrational normal modes. The normal modes of the hydroxyl functional group involved in a H-bond (first row) and the reference ones in the carboxylic acid (COOH) functional group (second row).

that molecules like enol-MDA behave as quasi-aromatic systems in the transition state [91]. This phenomenon tends to stabilize the molecule in the transition state, which lowers the energetic barrier. Capturing these intricate and subtle quantum phenomena and their dynamics requires high-level quantum chemistry methods, with CCSD(T) being the only method that converges to the correct energetic barrier. We found that systematically increasing the amount of electron correlation energy in our calculations, the energy barrier decreases as $\sim 13 \rightarrow \sim 5 \rightarrow \sim 4$ kcal mol $^{-1}$ for HF \rightarrow CCSD \rightarrow CCSD(T), respectively. This result highlights the importance of the correlation energy in such complex phenomena as the H-bond interaction and proton transfer.

The resulting free energy of the system obtained from MD simulations using sGDML@CCSD(T) at room temperature (see Fig. 2) displays a very low proton transition rate between the two oxygen atoms but it is still accessible at room temperature. This suggests that nuclear quantum effects would considerably increase the transition rate due to tunneling effects reshaping the FES and consequently its vibrational spectrum and thermodynamics. In general, the local electron density delocalization induced by intramolecular H-bonds influence macroscopic properties such as solubility and permeability [86], but a fundamentally different macroscopic implication of the symmetric H-bond is proton transport in extended systems like water. Therefore, the need of creating FFs capable of handling H-bonds in all of their flavors to accurately describe complex biological systems becomes obvious, and data-driven models enable a robust solution to this problem.

C. Hybridization change $sp^2 \leftrightarrow sp^3$ and $n \rightarrow \pi^*$ interactions

The two type of interactions mentioned in the previous subsections, electrostatic and H-bonds, are often approximately implemented in empirical FFs. The advantages of flexible fully data-driven models are featured in describing any quantum interaction coming from $-\mathbf{F}_i = \langle \Psi^* | \partial H / \partial \mathbf{x}_i | \Psi \rangle$, without relying on prior knowledge of the phenomena or its connection to any classical electrodynamic or mechanical concepts. To exemplify this, we consider the paracetamol and aspirin molecules. Their dynamics are strongly influenced by delocalized π electrons and the $n \rightarrow \pi^*$ interaction [92], respectively. The quantitative description of these phenomena occurs naturally in our data-driven FF even when only a restricted amount of reference data is available.

Hence, it is important to highlight that only now with the development of accurate ML-FFs trained on ab-initio data it is feasible to study in full extent the dynamical implication of such electronic effects at finite temperatures.

1. Hybridization state change in paracetamol

Paracetamol is a molecule with a shallow global minimum consisting in a planar configuration (Fig. 2) stabilized by being a conjugated system. From Fig. 2 for paracetamol, a steep energy increase is evident as illustrated by yellow regions in the PES. This represents the breaking of the conjugated state, given that the nitrogen atom changes its hybridization state from $sp^2 \rightarrow sp^3$ producing an energetic penalty (see electronic effects column in Fig. 2). Such electronic effects raise the energy, leading to an effectively inaccessible region in this direction of the PES. In fact, this region is hardly visited by MD simulations at 300 K, therefore it is not represented in the FES in Fig. 4-B.

Another important contribution to the planar structure of paracetamol is the electrostatic interaction between the lone pair on the carbonyl oxygen and the positively charged nearest hydrogen atom (see Fig. 4-D). We find a linear coupling between the acetamide main dihedral angle and the carbonyl dihedral angle, depicted in orange and black in Fig. 4-D, respectively. The projected FES in these two variables shown in Fig. 4-B-top, reveals that near the global minimum the system moves without altering the $d_{O\dots H}$ distance since the internal dihedral angle CNCO flexes to follow the minimum free energy path. Certainly, paracetamol is a highly fluxional molecule containing four correlated rotors moving in a complex PES due to its electronic structure.

2. $n \rightarrow \pi^*$ interaction in Aspirin

Another important electronic effect is the overlap between occupied (lone pair n) and antibonding (π^*) orbitals, this electronic effect is depicted in Fig. 2. The aspirin molecule is a particularly interesting case in this regard given the dominant role of this interaction in its molecular behaviour. This crucial $n \rightarrow \pi^*$ attraction interaction is responsible for the binding between the ester and carbonyl groups, which dictates the structure of the global minimum. This effect is amplified at finite temperature given that thermal fluctuations enhance the overlap between the lone pair, n , in the carbonyl group and the antibonding orbital in the ester group, π^* [60]. We have recently shown that the energy functional form of conventional FFs put into close contact four negatively charged oxygen atoms in aspirin; such strong charge repulsion leads to a misrepresentation of its PES [60]. The sole incorporation of the missing lone pairs in all four closely interacting oxygen atoms and their directionality could greatly improve the results in regular FFs. In general, there are many other electronic effects (e.g. $n \rightarrow \sigma^*$ interactions [93], hyperconjugation, configuration dependent charge densities and Jahn–Teller effect [94]) that are not explicitly incorporated in conventional FFs, nor captured by less robust ML-FF frameworks, which limits the reliability and predictive power of the dynamics. This rigorous requirement is justified by the increasing demand of computationally inexpensive and highly accurate PESs to interpret and obtain further insights into state-of-the-art spectroscopic experimental results [94–100].

In summary, we have analyzed a wide variety of energy landscapes reconstructed with high fidelity (see Fig. 2). Being trained directly on molecular forces from *ab initio* calculations, the generated PES encodes the broad range of fundamental interactions coming from the solution of the quantum-mechanical problem. This indicates that our model is a general ML approach capable of describing arbitrary interatomic interactions contained in the reference data.

V. CONCLUSIONS

We have presented the construction of molecular force fields using the symmetrized Gradient Domain Machine Learning model. This framework reconstructs high-dimensional manifold embedded in the training data from a few 100s of samples, allowing the use of highly-accurate *ab initio* reference data such as the “gold standard” CCSD(T) method. The flexibility of the sGDML model comes from its intrinsic nature of being a fully data-driven universal approximator, which grants the adaptability to describe any kind of quantum interaction.

This was demonstrated here by describing H-bonds, proton transfer, lone pairs, changes in hybridization states, steric repulsion, and $n \rightarrow \pi^*$ interactions and

by obtaining insightful results from molecular dynamics simulations. From a careful analysis of the PES and MD simulations, we highlighted the importance of electron lone pairs in generating the strong coupling between the two rotors in ethanol and in the dynamics of keto-MDA. On the other hand, the proper description of H-bonds revealed the proton dynamics in salicylic acid and enol-MDA molecules, yielding further understanding about its implications on the vibrational spectrum. Regarding electronic effects, the main contribution is that MLFFs can be used as trustworthy tools able to describe non-trivial interactions. For example, from MD simulations we observed the $sp^2 \leftrightarrow sp^3$ hybridization change of the nitrogen atom in paracetamol which help us to understand better the strong consequences of breaking the conjugated molecular system, and that in aspirin the $n \rightarrow \pi^*$ interaction is enhanced at higher temperatures giving extra stability to the molecular global minima.

The main advantages of the sGDML model over other machine learning methods are copious: (i) it is highly data efficient, due to being trained in the gradient domain, (ii) it is robust due to modeling all atomic interactions globally, without any kind of inherent non-unique partitioning of the energy or force contributions, (iii) it uses energy conservation as a prior, therefore encoding this fundamental physical law in the core of every gradient-domain FF model, and (iv) it correctly represents spatial symmetries using explicit constraints that are automatically extracted from the data.

A number of challenges remain to be solved in order to extend the applicability of sGDML to larger systems. In spite of its many advantages, a global model imposes limits on the maximum molecule size, as well as the training set size. Overcoming this fundamental limitation without compromising its robustness calls for the introduction of a well-reasoned fragmentation scheme that divides the reconstruction problem into smaller independent subproblems without oversimplifying the nature of the interactions. A data-driven approach could achieve this task in

a way that is tailored to preserving the intricate phenomena studied in this article, as opposed to applying general coarse-graining techniques. Such an approach will benefit from the explicit knowledge about fluxional symmetries within the system, which our algorithm is already able to extract. In its current formulation, the sGDML model captures different interaction scales, with no need to separating them. Nevertheless, an explicit decoupling of long-range interactions could be a new avenue to further increase data efficiency on the way to increasingly larger and complex molecules.

VI. SUPPLEMENTARY MATERIAL

See supplementary material for additional tables with prediction accuracy for total energies and forces from DFT and CCSD(T) data. Also, some supplementary notes regarding reference data generation, molecular dynamics details, and error propagation.

The sGDML code and documentation is available at <http://quantum-machine.org/gdml/>.

VII. ACKNOWLEDGEMENTS

We thank Dr. M. Gastegger for helpful discussions. S.C., A.T., and K.-R.M. thank the Deutsche Forschungsgemeinschaft (project MU 987/20-1) for funding this work. A.T. is funded by the European Research Council with ERC-CoG grant BeStMo. KRM acknowledges partial support by BMBF (BZML and BBDC) as well as by the Institute of Information & Communications Technology Planning & Evaluation (IITP) grant funded by the Korea government (No. 2017-0-00451). Part of this research was performed while the authors were visiting the Institute for Pure and Applied Mathematics, which is supported by the NSF.

-
- [1] B. J. Alder and T. E. Wainwright, *J. Chem. Phys.* **31**, 459 (1959).
- [2] A. Rahman, *Phys. Rev.* **136**, A405 (1964).
- [3] L. Verlet, *Phys. Rev.* **159**, 98 (1967).
- [4] A. Rahman and F. H. Stillinger, *J. Chem. Phys.* **55**, 3336 (1971).
- [5] M. S. Daw and M. I. Baskes, *Phys. Rev. B* **29**, 6443 (1984).
- [6] J. Tersoff, *Phys. Rev. B* **37**, 6991 (1988).
- [7] A. Warshel, P. K. Sharma, M. Kato, and W. W. Parson, *Biochim. Biophys. Acta, Proteins Proteomics* **1764**, 1647 (2006).
- [8] W. L. Jorgensen, J. Chandrasekhar, J. D. Madura, R. W. Impey, and M. L. Klein, *J. Chem. Phys.* **79**, 926 (1983).
- [9] M. W. Mahoney and W. L. Jorgensen, *J. Chem. Phys.* **112**, 8910 (2000).
- [10] P. K. Weiner and P. A. Kollman, *J. Comput. Chem.* **2**, 287 (1981).
- [11] B. R. Brooks, R. E. Bruccoleri, B. D. Olafson, D. J. States, S. Swaminathan, and M. Karplus, *J. Comput. Chem.* **4**, 187 (1983).
- [12] T. A. Halgren, *J. Comput. Chem.* **17**, 490 (1996).
- [13] T. A. Soares, P. H. Hünenberger, M. A. Kastenholz, V. Kräutler, T. Lenz, R. D. Lins, C. Oostenbrink, and W. F. van Gunsteren, *J. Comput. Chem.* **26**, 725 (2005).
- [14] M. Rupp, A. Tkatchenko, K.-R. Müller, and O. A. von Lilienfeld, *Phys. Rev. Lett.* **108**, 58301 (2012).
- [15] A. P. Bartók, M. C. Payne, R. Kondor, and G. Csányi, *Phys. Rev. Lett.* **104**, 136403 (2010).
- [16] K. Hansen, G. Montavon, F. Biegler, S. Fazli, M. Rupp, M. Scheffler, O. A. von Lilienfeld, A. Tkatchenko, and K.-R. Müller, *J. Chem. Theory Comput.* **9**, 3404 (2013).

- [17] A. P. Bartók and G. Csányi, *Int. J. Quantum Chem.* **115**, 1051 (2015).
- [18] K. Hansen, F. Biegler, R. Ramakrishnan, W. Pronobis, O. A. von Lilienfeld, K.-R. Müller, and A. Tkatchenko, *J. Phys. Chem. Lett.* **6**, 2326 (2015).
- [19] M. Rupp, R. Ramakrishnan, and O. A. von Lilienfeld, *J. Phys. Chem. Lett.* **6**, 3309 (2015).
- [20] S. De, A. P. Bartók, G. Csányi, and M. Ceriotti, *Phys. Chem. Chem. Phys.* **18**, 13754 (2016).
- [21] N. Artrith, A. Urban, and G. Ceder, *Phys. Rev. B* **96**, 014112 (2017).
- [22] A. P. Bartók, S. De, C. Poelking, N. Bernstein, J. R. Kermode, G. Csányi, and M. Ceriotti, *Sci. Adv.* **3**, e1701816 (2017).
- [23] A. Glielmo, P. Sollich, and A. De Vita, *Phys. Rev. B* **95**, 214302 (2017).
- [24] K. Yao, J. E. Herr, and J. Parkhill, *J. Chem. Phys.* **146**, 014106 (2017).
- [25] F. A. Faber, L. Hutchison, B. Huang, J. Gilmer, S. S. Schoenholz, G. E. Dahl, O. Vinyals, S. Kearnes, P. F. Riley, and O. A. von Lilienfeld, *J. Chem. Theory Comput.* **13**, 5255 (2017).
- [26] M. Eickenberg, G. Exarchakis, M. Hirn, S. Mallat, and L. Thiry, *J. Chem. Phys.* **148**, 241732 (2018).
- [27] A. Glielmo, C. Zeni, and A. De Vita, *Phys. Rev. B* **97**, 184307 (2018).
- [28] A. Grisafi, D. M. Wilkins, G. Csányi, and M. Ceriotti, *Phys. Rev. Lett.* **120**, 036002 (2018).
- [29] Y.-H. Tang, D. Zhang, and G. E. Karniadakis, *J. Chem. Phys.* **148**, 034101 (2018).
- [30] W. Pronobis, A. Tkatchenko, and K.-R. Müller, *J. Chem. Theory Comput.* **14**, 2991 (2018).
- [31] F. A. Faber, A. S. Christensen, B. Huang, and O. A. von Lilienfeld, *J. Chem. Phys.* **148**, 241717 (2018).
- [32] J. Behler and M. Parrinello, *Phys. Rev. Lett.* **98**, 146401 (2007).
- [33] K. V. J. Jose, N. Artrith, and J. Behler, *J. Chem. Phys.* **136**, 194111 (2012).
- [34] J. Behler, *J. Chem. Phys.* **145**, 170901 (2016).
- [35] M. Gastegger, J. Behler, and P. Marquetand, *Chem. Sci.* **8**, 6924 (2017).
- [36] K. T. Schütt, F. Arbabzadah, S. Chmiela, K. R. Müller, and A. Tkatchenko, *Nat. Commun.* **8**, 13890 (2017).
- [37] K. T. Schütt, H. E. Sauceda, P.-J. Kindermans, A. Tkatchenko, and K.-R. Müller, *J. Chem. Phys.* **148**, 241722 (2018).
- [38] K. T. Schütt, P.-J. Kindermans, H. E. Sauceda, S. Chmiela, A. Tkatchenko, and K.-R. Müller, in *Advances in Neural Information Processing Systems 30* (Curran Associates, Inc., 2017) pp. 991–1001.
- [39] K. Ryczko, K. Mills, I. Luchak, C. Homenick, and I. Tamblyn, *Comput. Mater. Sci.* **149**, 134 (2018).
- [40] L. Zhang, J. Han, H. Wang, R. Car, and E. Weinan, *Phys. Rev. Lett.* **120**, 143001 (2018).
- [41] Z. Li, J. R. Kermode, and A. De Vita, *Phys. Rev. Lett.* **114**, 096405 (2015).
- [42] E. V. Podryabinkin and A. V. Shapeev, *Comput. Mater. Sci.* **140**, 171 (2017).
- [43] P. O. Dral, A. Owens, S. N. Yurchenko, and W. Thiel, *J. Chem. Phys.* **146**, 244108 (2017).
- [44] A. Mardt, L. Pasquali, H. Wu, and F. Noé, *Nat. Commun.* **9**, 5 (2018).
- [45] F. Noé and H. Wu, “Boltzmann generators - sampling equilibrium states of many-body systems with deep learning,” (2018), [arXiv:1812.01729](https://arxiv.org/abs/1812.01729).
- [46] A. P. Bartók, R. Kondor, and G. Csányi, *Phys. Rev. B* **87**, 184115 (2013).
- [47] G. Montavon, M. Rupp, V. Gobre, A. Vazquez-Mayagoitia, K. Hansen, A. Tkatchenko, K.-R. Müller, and O. A. von Lilienfeld, *New J. Phys.* **15**, 95003 (2013).
- [48] V. Botu and R. Ramprasad, *Phys. Rev. B* **92**, 094306 (2015).
- [49] F. Brockherde, L. Vogt, L. Li, M. E. Tuckerman, K. Burke, and K.-R. Müller, *Nat. Commun.* **8**, 872 (2017).
- [50] T. D. Huan, R. Batra, J. Chapman, S. Krishnan, L. Chen, and R. Ramprasad, *NPJ Comput. Mater.* **3**, 37 (2017).
- [51] T. Beraud, R. A. DiStasio Jr, A. Tkatchenko, and O. A. Von Lilienfeld, *J. Chem. Phys.* **148**, 241706 (2018).
- [52] N. Lubbers, J. S. Smith, and K. Barros, *J. Chem. Phys.* **148**, 241715 (2018).
- [53] K. Kanamori, K. Toyoura, J. Honda, K. Hattori, A. Seko, M. Karasuyama, K. Shitara, M. Shiga, A. Kuwabara, and I. Takeuchi, *Phys. Rev. B* **97**, 125124 (2018).
- [54] T. S. Hy, S. Trivedi, H. Pan, B. M. Anderson, and R. Kondor, *J. Chem. Phys.* **148**, 241745 (2018).
- [55] J. S. Smith, O. Isayev, and A. E. Roitberg, *Chem. Sci.* **8**, 3192 (2017).
- [56] J. Wang, C. Wehmeyer, F. Noé, and C. Clementi, “Machine learning of coarse-grained molecular dynamics force fields,” (2018), [arXiv:1812.01736](https://arxiv.org/abs/1812.01736).
- [57] R. Winter, F. Montanari, F. Noé, and D.-A. Clevert, *Chem. Sci.* (2019), [10.1039/C8SC04175J](https://doi.org/10.1039/C8SC04175J).
- [58] A. S. Christensen, F. A. Faber, and O. A. von Lilienfeld, “Operators in machine learning: Response properties in chemical space,” (2018), [arXiv:1807.08811](https://arxiv.org/abs/1807.08811).
- [59] S. Chmiela, A. Tkatchenko, H. E. Sauceda, I. Poltavsky, K. T. Schütt, and K.-R. Müller, *Sci. Adv.* **3**, e1603015 (2017).
- [60] S. Chmiela, H. E. Sauceda, K.-R. Müller, and A. Tkatchenko, *Nat. Commun.* **9**, 3887 (2018).
- [61] S. Chmiela, H. E. Sauceda, I. Poltavsky, K.-R. Müller, and A. Tkatchenko, “sgdml: Constructing accurate and data efficient molecular force fields using machine learning,” (2018), [arXiv:1812.04986](https://arxiv.org/abs/1812.04986).
- [62] K. Yao, J. E. Herr, D. W. Toth, R. Mckintyre, and J. Parkhill, *Chem. Sci.* **9**, 2261 (2018).
- [63] K. T. Schütt, P. Kessel, M. Gastegger, K. A. Nicoli, A. Tkatchenko, and K. R. Müller, *J. Chem. Theory Comput.* **15**, 448 (2019).
- [64] M. Alber, S. Lapuschkin, P. Seegerer, M. Hgele, K. T. Schütt, G. Montavon, W. Samek, K.-R. Müller, S. Dähne, and P.-J. Kindermans, “iNNvestigate neural networks!” (2018), [arXiv:1808.04260](https://arxiv.org/abs/1808.04260).
- [65] M. Meila, S. Koelle, and H. Zhang, “A regression approach for explaining manifold embedding coordinates,” (2018), [arXiv:1811.11891](https://arxiv.org/abs/1811.11891).
- [66] J. Friedman, T. Hastie, and R. Tibshirani, *The elements of statistical learning*, Vol. 1 (Springer series in statistics New York, NY, USA:, 2001).
- [67] V. N. Vapnik, *The Nature of Statistical Learning Theory* (Springer-Verlag, Berlin, Heidelberg, 1995).
- [68] J. Behler, S. Lorenz, and K. Reuter, *J. Chem. Phys.* **127**, 014705 (2007).
- [69] J. Behler, *J. Chem. Phys.* **134**, 074106 (2011).
- [70] J. Behler, *Phys. Chem. Chem. Phys.* **13**, 17930 (2011).

- [71] E. Solak, R. Murray-smith, W. E. Leithead, D. J. Leith, and C. E. Rasmussen, in *Advances in Neural Information Processing Systems 15* (MIT Press, 2003) pp. 1057–1064.
- [72] A. Choudhary, K. J. Kamer, and R. T. Raines, *J. Org. Chem.* **76**, 7933 (2011).
- [73] S. Blanco and J. C. López, *J. Phys. Chem. Lett.* **9**, 4632 (2018).
- [74] J. P. Perdew, K. Burke, and M. Ernzerhof, *Phys. Rev. Lett.* **77**, 3865 (1996).
- [75] A. Tkatchenko and M. Scheffler, *Phys. Rev. Lett.* **102**, 073005 (2009).
- [76] V. Blum, R. Gehrke, F. Hanke, P. Havu, V. Havu, X. Ren, K. Reuter, and M. Scheffler, *Comput. Phys. Commun.* **180**, 2175 (2009).
- [77] J. M. Turney, A. C. Simmonett, R. M. Parrish, E. G. Hohenstein, F. A. Evangelista, J. T. Fermann, B. J. Mintz, L. A. Burns, J. J. Wilke, M. L. Abrams, N. J. Russ, M. L. Leininger, C. L. Janssen, E. T. Seidl, W. D. Allen, H. F. Schaefer, R. A. King, E. F. Valeev, C. D. Sherrill, and T. D. Crawford, *WIREs Comput. Mol. Sci.* **2**, 556 (2012).
- [78] R. M. Parrish, L. A. Burns, D. G. A. Smith, A. C. Simmonett, A. E. DePrince, E. G. Hohenstein, U. Bozkaya, A. Y. Sokolov, R. Di Remigio, R. M. Richard, J. F. Gonthier, A. M. James, H. R. McAlexander, A. Kumar, M. Saitow, X. Wang, B. P. Pritchard, P. Verma, H. F. Schaefer, K. Patkowski, R. A. King, E. F. Valeev, F. A. Evangelista, J. M. Turney, T. D. Crawford, and C. D. Sherrill, *J. Chem. Theory Comput.* **13**, 3185 (2017).
- [79] D. G. A. Smith, L. A. Burns, D. A. Sirianni, D. R. Nascimento, A. Kumar, A. M. James, J. B. Schriber, T. Zhang, B. Zhang, A. S. Abbott, E. J. Berquist, M. H. Lechner, L. A. Cunha, A. G. Heide, J. M. Waldrop, T. Y. Takeshita, A. Alenaizan, D. Neuhauser, R. A. King, A. C. Simmonett, J. M. Turney, H. F. Schaefer, F. A. Evangelista, A. E. DePrince, T. D. Crawford, K. Patkowski, and C. D. Sherrill, *J. Chem. Theory Comput.* **14**, 3504 (2018).
- [80] K. Hansen, F. Biegler, R. Ramakrishnan, W. Pronobis, O. A. von Lilienfeld, K.-R. Müller, and A. Tkatchenko, *J. Phys. Chem. Lett.* **6**, 2326 (2015).
- [81] R. Ramakrishnan, P. O. Dral, M. Rupp, and O. A. von Lilienfeld, *J. Chem. Theory Comput.* **11**, 2087 (2015).
- [82] B. Huang and O. A. von Lilienfeld, “The ”dna” of chemistry: Scalable quantum machine learning with ”amons”,” (2017), [arXiv:1707.04146](https://arxiv.org/abs/1707.04146).
- [83] B. Guillot, *J. Mol. Liq.* **101**, 219 (2002).
- [84] T. Oroguchi and M. Nakasako, *Sci. Rep.* **7**, 15859 (2017).
- [85] S. Scheiner, *Molecules* **22**, 1521 (2017).
- [86] B. Kuhn, P. Mohr, and M. Stahl, *J. Med. Chem.* **53**, 2601 (2010).
- [87] P. Hobza, *Int. J. Quantum Chem.* **90**, 1071 (2002).
- [88] A. Karpfen and E. S. Kryachko, *J. Phys. Chem. A* **113**, 5217 (2009).
- [89] C. Wang, D. Danovich, S. Shaik, and Y. Mo, *J. Chem. Theory Comput.* **13**, 1626 (2017).
- [90] K. Goldsby and R. Chang, *Chemistry* (McGraw-Hill Higher Education, 2015).
- [91] A. Martyniak, I. Majerz, and A. Filarowski, *RSC Adv.* **2**, 8135 (2012).
- [92] R. W. Newberry and R. T. Raines, *Acc. Chem. Res.* **50**, 1838 (2017).
- [93] R. Deepak and R. Sankararamakrishnan, *Biophys. J.* **110**, 1967 (2016).
- [94] R. Sarkar, S. R. Reddy, S. Mahapatra, and H. Köppel, *Chem. Phys.* **482**, 39 (2017).
- [95] P. Gruene, D. M. Rayner, B. Redlich, A. F. G. van der Meer, J. T. Lyon, G. Meijer, and A. Fielicke, *Science* **321**, 674 (2008).
- [96] C. Romanescu, D. J. Harding, A. Fielicke, and L.-S. Wang, *J. Chem. Phys.* **137**, 014317 (2012).
- [97] R. M. Balabin, *Phys. Chem. Chem. Phys.* **12**, 5980 (2010).
- [98] J. A. Ruiz-Santoyo, J. Wilke, M. Wilke, J. T. Yi, D. W. Pratt, M. Schmitt, and L. Ivarez Valtierra, *J. Chem. Phys.* **144**, 044303 (2016).
- [99] J. A. Davies, L. E. Whalley, and K. L. Reid, *Phys. Chem. Chem. Phys.* **19**, 5051 (2017).
- [100] F. Gmerek, B. Stuhlmann, E. Pehlivanovic, and M. Schmitt, *J. Mol. Struct* **1143**, 265 (2017).