

Incorporating long-range physics in atomic-scale machine learning

Andrea Grisafi and Michele Ceriotti

*Laboratory of Computational Science and Modeling, IMX,
École Polytechnique Fédérale de Lausanne, 1015 Lausanne, Switzerland*

The most successful and popular machine learning models of atomic-scale properties derive their transferability from a locality ansatz. The properties of a large molecule or a bulk material are written as a sum over contributions that depend on the configurations within finite atom-centered environments. The obvious downside of this approach is that it cannot capture non-local, non-additive effects such as those arising due to long-range electrostatics or quantum interference. We propose a solution to this problem by introducing non-local representations of the system that are remapped as feature vectors that are defined locally and are equivariant in $O(3)$. We consider in particular one form that has the same asymptotic behavior as the electrostatic potential. We demonstrate that this framework can capture non-local, long-range physics by building a model for the electrostatic energy of randomly distributed point-charges, for the unrelaxed binding curves of charged organic molecular dimers, and for the electronic dielectric response of liquid water. By combining a representation of the system that is sensitive to long-range correlations with the transferability of an atom-centered additive model, this method outperforms current state-of-the-art machine-learning schemes, and provides a conceptual framework to incorporate non-local physics into atomistic machine learning.

INTRODUCTION

In recent years, atomistic machine learning models have become increasingly popular as a way to perform fast predictions of molecular and material properties with the accuracy of first-principle quantum mechanical calculations [1], but a much reduced cost. The success of these methods has gone hand-in-hand with the progress in constructing representations for molecular and materials configurations that are flexible enough to be transferred across a wide spectrum of different atomic arrangements, while satisfying, at the same time, stringent symmetry constraints [2–6].

At the core of the vast majority of transferable machine-learning model for physical properties lies the local nature of the underlying atomistic representation. This is usually constructed by considering the set of atomic coordinates that are included within spherical environments of a given radial cutoff around any arbitrary atomic center [7–9]. The prediction of a given physical property is therefore formally decomposed in the sum of atom-centered contributions that effectively incorporate information associated with many-body structural correlations between atoms in each local environment. This locality assumption is very convenient, as it keeps at bay the dimensionality of the regression problem that is modelled by ML, and is physically justified by the nearsightedness principle of electronic matter [10]. The major drawback is that it neglects long-range physical effects. Long-range electrostatic interactions, for example, are known to play a fundamental role in the description of ionic systems [11], macroscopically polarized interfaces [12], electrode surfaces [13] and nano-science in general [14]. In all these cases, the pathologically slow decay $\sim 1/r$ of the Coulomb interaction makes it virtually impossible to reach convergence while using a local machine-learning scheme, which

is reflected in an effective limit to the accuracy that can be reached by these models.

The problem of incorporating long-range effects in electronic energy predictions is usually tackled by explicitly separating the local many-body contribution to the total energy from a classical electrostatic term approximated via pairwise Coulomb interactions. This can be done either by direct subtraction of the Ewald-like electrostatic energy of the system [15, 16], or by machine learning, in turn, the partial charges and the atomic multipoles that determine the long-range electrostatics [17–22]. Other more sophisticated models, specifically designed for ionic systems, rely on a charge equilibration scheme [23, 24].

Beyond electrostatic energies, the breakdown of a local machine learning model is particularly pronounced when dealing with intrinsically non-local quantities like the dielectric response of a condensed-phase medium [6]. This non-locality has to do both with the effect of the far-field electrostatics [25], and to the topological quantum nature of the macroscopic polarization of an infinitely extended material [26, 27]. In this case, the problem can possibly be bypassed by adopting specific physical prescriptions. Examples of this can be found in Ref. [6], where the dielectric tensor ϵ_∞ of liquid water is learned indirectly by building a model for an effective molecular polarizability that is mapped to ϵ_∞ through the Clausius-Mossotti relationship [25]. In the context of reproducing the autocorrelation function of the macroscopic polarization of liquid water, another strategy has recently been adopted, where the selected learning targets are the positions of the Wannier centers that are used to recast the electron density of the system into a set of point-charges [28].

By and large, the learning models previously described tackle the problem of including long-range phenomena by making use of an *ad hoc* definition of the electrostatic energy, or dielectric response, in terms of local atomic

quantities. Although successful, these kind of approaches have the downside of being very system dependent and, as such, hardly transferable across systems that have a different nature, e.g., those related to charge transfer, or to charge polarizability in (near)-metallic systems [29]. Capturing long-range effects without any prior assumption on the nature of the learning target is a difficult task to accomplish with the methods currently available. Most of the approaches that have explicitly attempted to do so, such as Coulomb kernels [30], many-body tensor representations [31], or multi-scale invariants [32], are built upon a global representation of the system rather than on an additive atom-centred model.

Here we propose a simple, yet elegant, solution to this problem, where the non-local character of the target property is incorporated in a symmetry-equivariant fashion into an atom-centered representation. In doing so, we construct a formalism that ensures that the resulting features exhibit the correct asymptotic dependence on the distribution of atoms in the far-field. This representation can be incorporated straightforwardly into conventional, additive machine learning models. While the idea is very general, we present as an example a model that has an asymptotic behavior consistent with electrostatic interactions. We show that it can be used successfully to build a local machine learning model that accurately reproduces Coulomb interactions between point particles, the binding curves of charged organic fragments, and the electronic dielectric response of bulk water.

LONG-DISTANCE EQUIVARIANT REPRESENTATION

Let us start from the same formal definition of a ML representation of a structure \mathcal{A} that was introduced in Ref. [33], written in the position basis as a decorated atom density

$$\langle \mathbf{r} | \mathcal{A} \rangle = \sum_i g(\mathbf{r} - \mathbf{r}_i) |\alpha_i\rangle, \quad (1)$$

where the index i runs over all the atoms in the structure, g is a Gaussian (or another localized function) peaked at each atom's position \mathbf{r}_i , and $|\alpha_i\rangle$ is an abstract vector that encodes the chemical nature of the atom. We now introduce an atom-density potential representation

$$\langle \mathbf{r} | \mathcal{V}^p \rangle = \sum_i |\alpha_i\rangle \int d\mathbf{r}' \frac{g(\mathbf{r}' - \mathbf{r}_i)}{|\mathbf{r}' - \mathbf{r}|^p}. \quad (2)$$

The rationale for performing this transformation (that can be seen as the action of a linear integral operator on $|\mathcal{A}\rangle$) is that, whereas $\langle \mathbf{r} | \mathcal{A} \rangle$ contains information only about the atoms in the vicinity of \mathbf{r} , $\langle \mathbf{r} | \mathcal{V}^p \rangle$ contains information about the position of *all* atoms in the structure, with a dependence on the position of the i -th atom that decays

asymptotically as $|\mathbf{r} - \mathbf{r}_i|^{-p}$ [34]. The physical significance of $|\mathcal{V}^p\rangle$ is obvious, if one considers typical forms of the interactions between atoms and molecules. For instance, if we had a single species and interpreted (1) as a charge density, $\langle \mathbf{r} | \mathcal{V}^1 \rangle$ would correspond to the electrostatic potential generated by such charge density. Analogously, the $p = 6$ case would provide the formally correct asymptotic limit of the energy per particle associated with dispersion interactions [35], which has inspired previous representations of local environments such as aSLATM [36].

Proceeding as in Ref. 33, we can symmetrize the representation over the continuous translation group, taking a tensor product with the density representation to preserve structural information. One obtains the symmetrized ket

$$\langle \mathbf{r} | \mathcal{A} \mathcal{V}^p \rangle_{\hat{t}} = \int d\hat{t} \langle \mathbf{0} | \hat{t} | \mathcal{A} \rangle \langle \mathbf{r} | \hat{t} | \mathcal{V}^p \rangle = \sum_j |\alpha_j\rangle \langle \mathbf{r} | \mathcal{V}_j^p \rangle, \quad (3)$$

where we introduced the shorthand notation (see the SI for a full derivation)

$$\langle \mathbf{r} | \mathcal{V}_j^p \rangle = \sum_i |\alpha_i\rangle \int d\mathbf{r}' \frac{(g \star g)(\mathbf{r}' - (\mathbf{r}_i - \mathbf{r}_j))}{|\mathbf{r}' - \mathbf{r}|^p}. \quad (4)$$

Modulo the re-definition of the atom density function as the auto-correlation of g , $\langle \mathbf{r} | \mathcal{V}_j^p \rangle$ is just the atom-density potential (2) computed using \mathbf{r}_j as the origin of the reference frame.

Symmetrization over the translation group leads naturally to a structural representation that amounts to a sum over atom-centred descriptors – foreshadowing an additive property model built on such feature vector. Particularly for low values of the potential exponent p , however, the integral in Eq. (4) introduces a substantially non-local behavior. The value of $\langle \mathbf{r} | \mathcal{V}_j^p \rangle$ in the vicinity of the central atom j can in principle depend on the position of atoms that are very far from it, *even if one introduces a cutoff function that restricts the range of $\langle \mathbf{r} | \mathcal{V}_j^p \rangle$ around the central atom, and hence its complexity*. One can then symmetrize further over the rotation group and over inversion symmetry. We will refer from now on to the resulting class of atomistic representations that capture long-range interactions based on the local value of an atom-density potential as the *long-distance equivariant* (LODE) framework. In the following we will focus on the case of $p = 1$, that corresponds to electrostatic interactions.

It is instructive to first consider the case of the first order spherical invariant, and to take the limit in which the atom density is represented by Dirac- δ distributions. It is easy to see that in this limit

$$\langle \alpha \mathbf{r} | \mathcal{V}_j^1 \rangle = \sum_{i \in \alpha} \frac{1}{|\mathbf{r} - \mathbf{r}_{ij}|}, \quad (5)$$

where $\mathbf{r}_{ij} = \mathbf{r}_i - \mathbf{r}_j$. Integrating over the SO(3) group

yields the first invariant

$$\langle \alpha \mathbf{r} | \mathcal{V}_j^{1(1)} \rangle = \int d\hat{R} \langle \alpha \mathbf{r} \hat{\mathbf{r}} | \hat{R} | \mathcal{V}_j^1 \rangle = \sum_{i \in \alpha} \min \left[\frac{1}{r}, \frac{1}{r_{ij}} \right], \quad (6)$$

that simply sums up $1/r_{ij}$ terms for all atoms *outside* the region over which the LODE representation is computed. Ignoring the contribution from the atoms within the cutoff, that can be better characterized by other atomic structure representations, a linear model built on these features is equivalent to a fixed point-charge electrostatic model. In other words, in this limit the radial dependence of the regression weights is integrated out, and the weights associated with each pair of central atom type α' and neighbor type α corresponds to the product of the atomic charges $q_{\alpha'}$ and q_{α} .

While this construction is very revealing, it is clear that its descriptive power is limited. Non-linear kernel models can provide a more flexible functional form, but higher-order invariants provide a systematic way of incorporating more information on structural features. As in the SOAP framework for the atom density [3, 33, 37], the most convenient way to compute such invariants involves writing the scalar field associated with the species α on a basis of radial functions $R_n(r)$ and spherical harmonics $Y_m^l(\hat{\mathbf{r}})$,

$$\langle \alpha n l m | \mathcal{V}_j^p \rangle = \int d\mathbf{r} R_n(r) Y_m^l(\hat{\mathbf{r}})^* \langle \alpha \mathbf{r} | \mathcal{V}_j^p \rangle \quad (7)$$

and then computing the appropriate spherically-covariant combinations. For example, for rotationally invariant representations of order $\nu = 2$ (the form that is equivalent to the SOAP power spectrum and that we will use in applications)

$$\langle \alpha n \alpha' n' l | \mathcal{V}_j^{p(2)} \rangle = \sum_{|m| \leq l} \frac{\langle \alpha n l m | \mathcal{V}_j^p \rangle^* \langle \alpha' n' l m | \mathcal{V}_j^p \rangle}{\sqrt{2l+1}}. \quad (8)$$

The extension to higher orders in spatial correlations $\nu > 2$ and/or to rotationally covariant representations of a given spherical-tensor order $\lambda > 0$ is straightforward based on the analogous density-based counterparts [6, 33, 38]. Note that it is also possible to compute representations that combine different values of p , and even $p = 0$, corresponding to the atom-density field. A systematic investigations of the various combinations, and their physical meaning, is left for future work.

Efficient evaluation of the LODE representation

As discussed in the SI, for molecules and clusters the expansion (7) can be computed conveniently in real space, by numerical integration on appropriate atom-centred grids. For a bulk system, described by a periodically-repeated supercell, the long-range nature of the integral

kernel that appears in (4) would make computing the expansion prohibitive. This is exactly the same problem one faces when evaluating electrostatic interactions in the condensed phase, and fortunately it has long been solved, e.g., with the many techniques based on the use of a plane-waves auxiliary basis [39, 40]. Consider the plane-wave definition as $\langle \mathbf{r} | \mathbf{k} \rangle = e^{i\mathbf{k} \cdot \mathbf{r}}$, with $\{\mathbf{k}\}$ representing a set of wave-vectors that are compatible with the simulation box. The fact we start from a smooth, Gaussian atom density, means that in practice one needs only a manageable number of plane waves. In particular, the width σ of the Gaussian density determines the minimum wavelength that should be introduced in the the plane-wave expansion, so that \mathbf{k} -vectors only need to be generated within a sphere of radius k_{\max} of the order of $2\pi/\sigma$. In order to evaluate the local potential projections, it is then enough to include the identity resolution $\sum_{\mathbf{k}} |\mathbf{k}\rangle \langle \mathbf{k}|$ within the bracket of Eq. (7), i.e.

$$\langle \alpha n l m | \mathcal{V}_j^p \rangle = \sum_{\mathbf{k}} \langle n l m | \mathbf{k} \rangle \langle \alpha \mathbf{k} | \mathcal{V}_j^p \rangle. \quad (9)$$

As detailed in the SI, $\langle n l m | \mathbf{k} \rangle$ corresponds to the expansion in plane waves of the basis of the local environment representation, and can be computed analytically once and for all if the radial functions are taken to be Gaussian type orbitals [41]. Conversely, $\langle \alpha \mathbf{k} | \mathcal{V}_j^p \rangle$ represents the Fourier components of the potential generated by the Gaussian density of element α for the entire system, and can be readily computed analytically [42]. As a result, the geometric local nature of the representation of Eq. (9) is formally factorized from its system-dependent global character. It has not escaped our attention that Eqn. (9) could also be used to compute efficiently the coefficients of the density expansion that enter, for instance, the SOAP framework. In the context of electrostatic interactions, one should note that although the fictitious charge density distribution of Eq. (1) does not satisfy charge neutrality, one can avoid a divergence of the potential by ignoring the $\mathbf{k}=\mathbf{0}$ component from the sum of Eq. (9). Similarly, divergences in the potential for $p > 1$ can be eliminated by appropriately regularizing the $1/r^p$ divergence in reciprocal space.

RESULTS

We now proceed to test the performance of the LODE representation in the context of predicting scalar electrostatic properties. In all cases we use Gaussian process regression using simple polynomial kernels, to emphasize the role of the features - as opposed to the regression scheme - on the performance of the model. Details of the parameters used in each example are reported in the SI. We use the SOAP framework as the baseline for a comparison, which is appropriate given the close relation between the two approaches, and the excellent

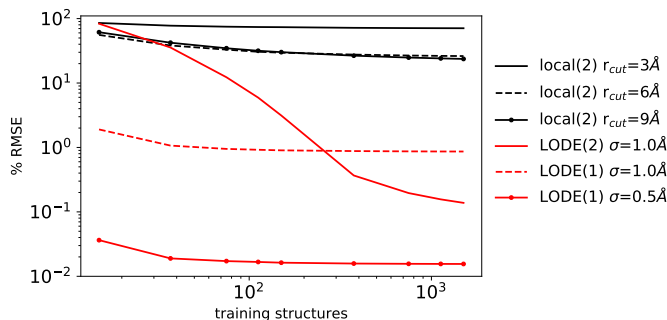


FIG. 1. Learning curves for the electrostatic energy of an idealized random gas of point charges. The model is trained on 1500 randomly selected configurations and tested on other 500 independent configurations. (*black full and dashed lines*) Local ML (SOAP) results at environment cutoffs of 3, 6 and 9 Å. (*red lines*) LODE($\nu = 1$) results at an environment cutoff of 2 Å and Gaussian smearing of 0.5 and 1.0 Å, and LODE($\nu = 2$) results with a cutoff of 3Å.

performances demonstrated by SOAP-based models. It is important however to stress that *any* local model with a finite cutoff will exhibit similar behavior as what we observe with SOAP. We also benchmark the combination of SOAP and LODE, that incorporates the advantages of both short-range and long-range models, realizing a kind of range-separated machine learning framework.

A gas of point charges

We begin by considering a toy system made of randomly distributed point-charges in a cubic box that is infinitely repeated in the three dimensions using periodic boundary conditions. The number of positive charges is equal to the number of negative charges, so that the system is overall neutral. To limit the amplitude of energy fluctuations, we discard configurations in which two charges are closer together than 2.5 Å. Following these prescriptions, we generate a total of 2000 configurations, each of which contains 64 atoms in cubic boxes spanning a broad range of densities, with side lengths between 12 and 20 Å. For each of these configurations, we compute the electrostatic energy using the Ewald method, as implemented in LAMMPS [43]. Fig. 1 compares the learning performance obtained using a local SOAP representation with different cutoffs, to the one obtained by direct application of the LODE representation. In both cases, a Gaussian width of $\sigma=1.0$ Å has been used to construct the density distribution of Eq. (1).

The figure clearly demonstrate the inefficiency of a local model when attempting to learn a property that is dominated by long-range effects. Given that the training set contains few configurations with atoms closer than 3 Å, the model with $r_{\text{cut}}=3$ Å is almost completely ineffective. Even increasing the cutoff up to 9 Å, a SOAP model barely

reaches an accuracy of about 20% RMSE when using the maximum number of training structures. A linear model built using the LODE($\nu = 1$) representation, on the other hand yields an error below 1% by using a handful of training points. As discussed above, this model represents exactly Coulomb interactions between fixed point charges, and the only reason the error does not converge to zero is the fact we use a Gaussian smearing in the definition of LODE, rather than δ distributions. This is apparent in the dramatic reduction of the error when halving the value of σ . A LODE($\nu = 2$) model, although initially less effective, possesses sufficient descriptive power to reach, and then overcome, the accuracy of the linear $\nu = 1, \sigma = 1$ Å model. This simple example highlights how difficult it is to incorporate long-range physics with a conventional local structure representation, and demonstrates that the LODE features can, on their own, be used as a very efficient description to predict the electrostatic energy of a system of fixed point charges.

Binding curves of charged dimers

We now consider a more realistic scenario, namely the problem of predicting the binding curves of a dataset of organic molecular dimers that carry an electric charge. We extract 661 different dimers containing H, C, N and O atoms from the BioFragment Database (BFDdb) [44], where at least one of the two monomers in each dimer configuration has a net charge. This choice ensures that we focus the exercise on a problem for which permanent electrostatic interactions play a prominent role. Contrary to the NaCl toy system, however, one cannot expect that a fixed point-charge model would suffice to predict the binding curves. The dataset contains a multitude of chemical moieties, including neutral polar fragments, highly polarizable groups, and provides a realistic assessment of how well a LODE model can perform in practice. For each of the 661 dimers, we consider 13 configurations where the reciprocal distance between the two monomers, defined as the distance between their geometric centers, spans an interval that can go from a minimum of ~ 3 Å to a maximum of ~ 8 Å. For each of these configurations, unrelaxed binding curves are computed at the DFT/B3LYP level using the FHI-aims quantum-chemistry package [45]. The training dataset is defined by considering the binding curves of the first 600 dimers out of the total of 661, while predictions are tested on the remaining 61. We also include the isolated monomers in the training set, so that the ML model has knowledge of the dissociation limit, and compute a few additional reference energies at larger separations, which are however not used for training. SOAP and LODE representations are defined within spherical environments of $r_{\text{cut}} = 3.0$ Å, while the Gaussian width of the density field is chosen to be $\sigma=0.3$ and 1.0 Å respectively.

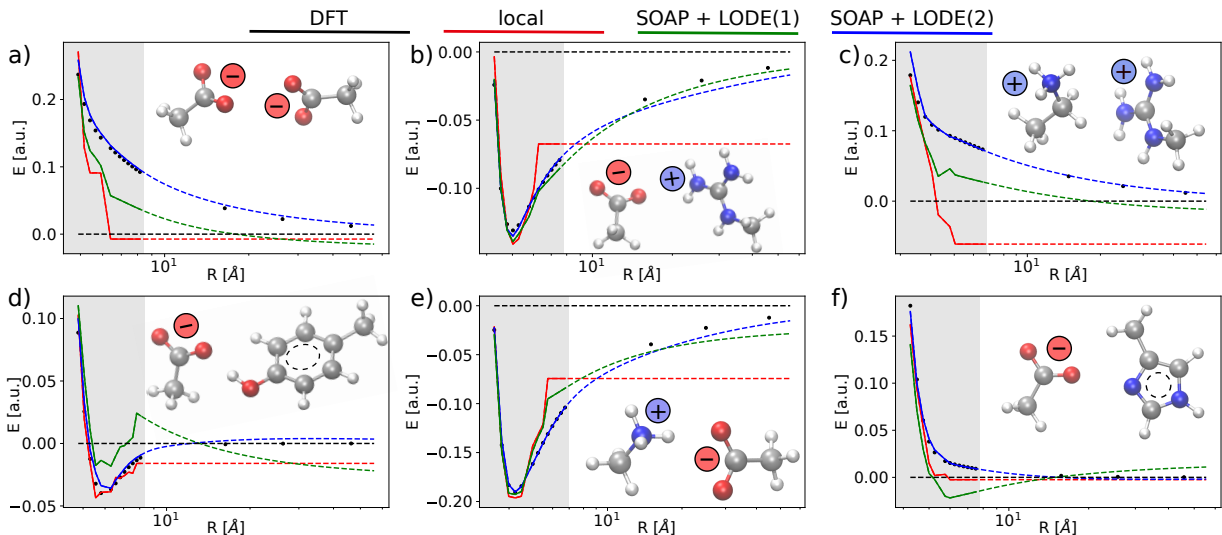


FIG. 2. Comparison of reference and predicted binding curves of six molecular dimers. (*black dots*) DFT reference calculations, (*red lines*) local SOAP predictions, (*green lines*) combined SOAP and LODE(1) predictions, (*blue lines*) combined SOAP and LODE(2) predictions. Full lines and shaded background represent the range of distances that is comparable to the geometries included in the training set. Dashed lines refer to predictions carried out in an extrapolative (long-range) regime.

Before carrying out the learning exercise, the reference DFT energies are baselined with respect to the monomer energies, so that the model only has to reproduce the interaction energies between the two fragments. Upon this baselining, we find that optimal SOAP performances correspond to a RMSE $\sim 20\%$, whereas a suitable combination between SOAP and LODE($\nu = 2$) allows us to bring the error down to $\sim 4\%$. This substantial improvement can be justified by the large discrepancy between the SOAP and SOAP+LODE accuracy in representing the interaction between the monomers at intermediate and large distance. To clarify the issue further, we plot in Fig. 2 the predicted binding curves of 6 test dimers, against the reference DFT calculations. We observe that a SOAP-based local description is overall able to capture the short-range interactions with good accuracy. However, it becomes less and less effective as the distance between the monomers increases, to the point of being completely blind to changes in interatomic distances when the environment's cutoff distance is overcome. Note that the performance of the local model at small separations is degraded substantially by the inclusion of fully dissociated dimers in the training set, because the representation cannot distinguish these configurations from those barely beyond the cutoff distance, that correspond to a non-zero value of the binding curve. The SOAP+LODE multi-scale description, in contrast, can recognize the changes in separation between the monomers, leading to a smooth asymptotic behavior of the predicted binding curve. Although a linear model incorporating LODE($\nu = 1$) allows us to halve the error made by SOAP down to $\sim 10\%$, it is not sufficiently expressive to achieve predictive accuracy - particularly for binding curves that involve neutral

monomers that do not have a $1/r$ asymptotic behavior.

This limitation can be addressed using a non-linear kernel based on SOAP+LODE($\nu = 2$). The resulting model is able to accurately predict the binding curves in the entire domain of distances, demonstrating its transferability across a vast spectrum of different chemical species and intermolecular configurations. This is particularly remarkable, as the SOAP+LODE($\nu = 2$) model does not only predict accurately systems that are dominated by monopole electrostatics (Fig.3-(a,b,c,e)), but also systems in which only one of the molecules is charged, and so interactions involve polarization as well as charge-dipole electrostatics (Fig.3-(d,f)). It should be noted, however, that the current scheme cannot transparently describe the physics of polarization or charge transfer. While the use of a composite SOAP+LODE kernel can describe how the environment of an atom affects its response to an external field, there is no explicit provision to represent how the field generated by far-away atoms depends on their neighboring structure.

Dielectric response of liquid water

As a final example, we revisit the problem of constructing a model of the infinite-frequency dielectric response tensor ϵ_∞ of liquid water. Details about the dataset generation and the computation of the dielectric tensors are reported in Ref. [6]. In that work, we argued that a local model was inefficient in learning dielectric response because of its collective nature, and showed that using the Clausius-Mossotti relationship to map ϵ_∞ to more local quantities was greatly improving the model. Here, LODE

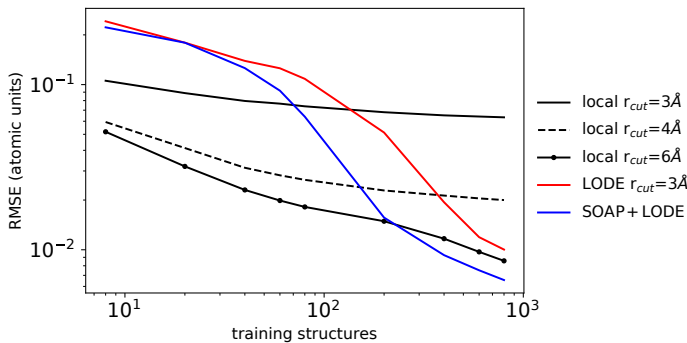


FIG. 3. Learning curves for the isotropic component of the dielectric response tensor ϵ_∞ of liquid water. The model is trained on up to 800 randomly selected configurations and tested on other 200 independent configurations. (black full and dashed lines) SOAP results with $r_{\text{cut}} = 3, 4$ and 6 \AA . (red line) LODE results with $r_{\text{cut}} = 3 \text{ \AA}$. (blue line) combined results of SOAP and LODE, both using $r_{\text{cut}} = 3 \text{ \AA}$.

learning performances are only tested for the isotropic component of the tensor $\epsilon_0 = \text{Tr}[\epsilon_\infty]$, which was shown to be most sensitive to the collective nature of the physics of dielectrics. Similarly to the case of the BFDB, we use a non-linear kernel that combines a SOAP representations computed using an optimal Gaussian width of $\sigma = 0.3 \text{ \AA}$, and LODE ($\nu = 2$) features constructed starting from a Gaussian density of $\sigma = 1.0 \text{ \AA}$. Figure 3 reports results obtained when learning on 800 randomly selected structures and predicting on other 200 independent configurations.

Similarly to what has been observed in the previous example, LODE performs much better than SOAP when relying upon a local description of $r_{\text{cut}} = 3 \text{ \AA}$. In this case, however, we observe a substantial improvement of the performance of SOAP when increasing the size of the local environments, eventually overcoming the LODE accuracy with a radial cutoff of $r_{\text{cut}} = 6 \text{ \AA}$. This might be a consequence of a less pronounced contribution of long-range tails, or - likely - of the fact that a cutoff of 6 \AA encompasses the entirety of the supercell, and therefore effectively provides a complete description of the input space of this specific dataset. Optimal ML predictions can be obtained when combining the fine-grained local description of SOAP at $r_{\text{cut}} = 3 \text{ \AA}$ with the coarse-grained and non-local description of LODE at the same cutoff. This behaviour highlights the multiscale character of ϵ_0 , meaning that both the local many-body information and the long-range electrostatic effects need to be considered to get accurate predictions. It is also important to stress that a combination of SOAP and LODE is not only beneficial in terms of learning performance, but can also reduce the computational effort in evaluating the feature vector - much like efficient methods for evaluating empirical potentials often treat separately short-range and long-range interactions.

CONCLUSIONS

Machine-learning of atomic-scale properties that are dominated by short-range interactions has reached a stage of maturity, with a substantial consensus about the ingredients of a successful model. The most commonly used frameworks incorporate symmetries and physical principles into the representation of atomic configurations, and achieve transferability by building additive property models. Furthermore, there is a growing understanding of the deep connections that exist between many of these methods, which is reflected in the fact that in most applications they reach similar levels of accuracy. In this paper we show how to extend these schemes in a way that makes it possible to incorporate long-range physics, without sacrificing the transferability of additive property models and the general applicability of rather abstract measures of atomic structure correlations. The crux lies in the definition of an atom-density potential that folds global information on the structure and composition of a system into a local representation, that (1) has a physically-motivated asymptotic behavior with interatomic separation and (2) can be efficiently computed in a symmetry-consistent fashion using similar ideas as those that underlie the SOAP framework and related approaches.

We apply this long-distance equivariant (LODE) representation focusing on the version that is based on a Coulomb-like atom-density potential. We demonstrate that, alone or in combination with SOAP, it outperforms local machine-learning methods in capturing long-range physics, for tasks that involve learning the electrostatic energy of a point-charge model, the binding curve of dimers of electrically charged organic fragments, and the dielectric constant of bulk water.

These examples are little more than an assay that proves that this scheme can incorporate efficiently long-range information in atomistic machine learning. More work is needed to draw a systematic, formal connection between a ML model built on LODE features and long-range interatomic potentials, much like a connection has been shown between linear models built on density-based features and short-range many-body potentials [33, 46, 47]; whether choosing other exponents in $\langle \mathbf{r} | \mathcal{V}^p \rangle$ can improve models of dispersion and of long-range effects that do not imply a characteristic asymptotic behavior; whether equivariant local features can be obtained by combining the expansion of the density and that of $\langle \mathbf{r} | \mathcal{V}^p \rangle$; whether the combination of SOAP and LODE can be used to improve the accuracy and the computational efficiency of existing ML forcefields; whether it is possible to incorporate polarizable atoms physics into the LODE framework. Future investigation will address these and many other questions, and unearth the full potential of this physics-inspired approach to atomistic machine learning.

ACKNOWLEDGMENTS

The Authors would like to thank Clemence Corminboeuf and Gábor Csányi for insightful comments on an early version of the manuscript. M.C and A.G. were supported by the European Research Council under the European Union's Horizon 2020 research and innovation programme (grant agreement no. 677013-HBMAP), and by the NCCR MARVEL, funded by the Swiss National Science Foundation. A.G. acknowledges funding by the MPG-EPFL Center for Molecular Nanoscience and Technology. We thank CSCS for providing CPU time under project id s843.

-
- [1] J. Behler, *Angewandte Chemie International Edition* **56**, 12828 (2017).
- [2] J. Behler and M. Parrinello, *Phys. Rev. Lett.* **98**, 146401 (2007).
- [3] A. P. Bartók, R. Kondor, and G. Csányi, *Phys. Rev. B* **87**, 184115 (2013).
- [4] A. Shapeev, *Multiscale Model. Sim.* **14**, 1153 (2016).
- [5] A. Glielmo, P. Sollich, and A. De Vita, *Phys. Rev. B* **95**, 214302 (2017).
- [6] A. Grisafi, D. M. Wilkins, G. Csányi, and M. Ceriotti, *Phys. Rev. Lett.* **120**, 036002 (2018).
- [7] A. P. Bartók, S. De, C. Poelking, N. Bernstein, J. R. Kermode, G. Csányi, and M. Ceriotti, *Sci. Adv.* **3** (2017).
- [8] S. Chmiela, A. Tkatchenko, H. E. Sauceda, I. Poltavsky, K. T. Schütt, and K.-R. Müller, *Sci. Adv.* **3** (2017).
- [9] L. Zhang, J. Han, H. Wang, R. Car, and W. E, *Phys. Rev. Lett.* **120**, 143001 (2018).
- [10] E. Prodan and W. Kohn, *Proceedings of the National Academy of Sciences* **102**, 11635 (2005).
- [11] R. Kjellander, *The Journal of Chemical Physics* **148**, 193701 (2018).
- [12] Z. Guo, F. Ambrosio, W. Chen, P. Gono, and A. Pasquarello, *Chemistry of Materials* **30**, 94 (2018).
- [13] R. Jorn, R. Kumar, D. P. Abraham, and G. A. Voth, *The Journal of Physical Chemistry C* **117**, 3747 (2013).
- [14] R. H. French, V. A. Parsegian, R. Podgornik, R. F. Rajter, A. Jagota, J. Luo, D. Asthagiri, M. K. Chaudhury, Y.-m. Chiang, S. Granick, S. Kalinin, M. Kardar, R. Kjellander, D. C. Langreth, J. Lewis, S. Lustig, D. Wesolowski, J. S. Wettlaufer, W.-Y. Ching, M. Finnis, F. Houlihan, O. A. von Lilienfeld, C. J. van Oss, and T. Zemb, *Rev. Mod. Phys.* **82**, 1887 (2010).
- [15] A. P. Bartók, M. C. Payne, R. Kondor, and G. Csányi, *Phys. Rev. Lett.* **104**, 136403 (2010).
- [16] Z. Deng, C. Chen, X.-G. Li, and S. P. Ong, *npj Computational Materials* **5**, 75 (2019).
- [17] N. Artrith, T. Morawietz, and J. Behler, *Phys. Rev. B* **83**, 153101 (2011).
- [18] T. Beraud, D. Andrienko, and O. A. von Lilienfeld, *J. Chem. Theory Comput.* **11**, 3225 (2015).
- [19] T. Beraud, R. A. DiStasio, A. Tkatchenko, and O. A. von Lilienfeld, *J. Chem. Phys.* **148**, 241706 (2018).
- [20] P. Bleiziffer, K. Schaller, and S. Riniker, *Journal of Chemical Information and Modeling* **58**, 579 (2018).
- [21] B. Nebgen, N. Lubbers, J. S. Smith, A. E. Sifain, A. Lokhov, O. Isayev, A. E. Roitberg, K. Barros, and S. Tretiak, *J. Chem. Theory Comput.* **14**, 4687 (2018).
- [22] K. Yao, J. E. Herr, D. Toth, R. Mckintyre, and J. Parkhill, *Chem. Sci.* **9**, 2261 (2018).
- [23] S. A. Ghasemi, A. Hofstetter, S. Saha, and S. Goedecker, *Phys. Rev. B* **92**, 045131 (2015).
- [24] S. Faraji, S. A. Ghasemi, S. Rostami, R. Rasoulkhani, B. Schaefer, S. Goedecker, and M. Amsler, *Phys. Rev. B* **95**, 104105 (2017).
- [25] C. Böttcher, O. van Belle, P. Bordewijk, and A. Rip, *Theory of electric polarization* (Elsevier Scientific Pub. Co., 1978).
- [26] R. Resta, *Rev. Mod. Phys.* **66**, 899 (1994).
- [27] R. Resta, *Journal of Physics: Condensed Matter* **22**, 123201 (2010).
- [28] L. Zhang, M. Chen, X. Wu, H. Wang, W. E, and R. Car, arXiv:1906.11434 (2019).
- [29] D. M. Wilkins, A. Grisafi, Y. Yang, K. U. Lao, R. A. DiStasio, and M. Ceriotti, *Proc. Natl. Acad. Sci.* **116**, 3401 (2019).
- [30] M. Rupp, A. Tkatchenko, K.-R. Müller, and O. A. von Lilienfeld, *Phys. Rev. Lett.* **108**, 058301 (2012).
- [31] H. Huo and M. Rupp, arXiv:1704.06439 (2017).
- [32] M. Hirn, S. Mallat, and N. Poilvert, *Multiscale Modeling & Simulation* **15**, 827 (2017).
- [33] M. J. Willatt, F. Musil, and M. Ceriotti, *The Journal of Chemical Physics* **150**, 154110 (2019).
- [34] Evaluation of the integral for $p > 1$ require some form of regularization or short-distance cutoff to remove the singularity for $\mathbf{r} \rightarrow \mathbf{r}_i$.
- [35] R. Dreizler and E. Gross, *Density Functional Theory: An Approach to the Quantum Many-Body Problem* (Springer Berlin Heidelberg, 2012).
- [36] B. Huang and O. A. von Lilienfeld, arXiv:1707.04146 (2017).
- [37] S. De, A. A. P. Bartók, G. Csányi, and M. Ceriotti, *Phys. Chem. Chem. Phys.* **18**, 13754 (2016).
- [38] A. Grisafi, D. M. Wilkins, M. J. Willatt, and M. Ceriotti, arXiv:1904.01623 (2019).
- [39] P. P. Ewald, *Annalen der Physik* **369**, 253 (1921).
- [40] U. Essmann, L. Perera, M. L. Berkowitz, T. Darden, H. Lee, and L. G. Pedersen, *The Journal of Chemical Physics* **103**, 8577 (1995).
- [41] K. Cahill, *Physical Mathematics* (Cambridge University Press, 2013).
- [42] M. P. Allen and D. J. Tildesley, *Computer Simulation of Liquids* (Clarendon Press, 1989).
- [43] S. Plimpton, *J. Comp. Phys.* **117**, 1 (1995).
- [44] L. A. Burns, J. C. Faver, Z. Zheng, M. S. Marshall, D. G. Smith, K. Vanommeslaeghe, A. D. MacKerell, K. M. Merz, and C. D. Sherrill, *J. Chem. Phys.* **147**, 161727 (2017).
- [45] V. Blum, R. Gehrke, F. Hanke, P. Havu, V. Havu, X. Ren, K. Reuter, and M. Scheffler, *Computer Physics Communications* **180**, 2175 (2009).
- [46] A. Glielmo, C. Zeni, and A. De Vita, *Phys. Rev. B* **97**, 184307 (2018).
- [47] R. Drautz, *Phys. Rev. B* **99**, 014104 (2019).