

# Derivative learning of tensorial quantities – Predicting finite temperature infrared spectra from first principles

Bernhard Schmiedmayer<sup>1</sup> and Georg Kresse<sup>1, a)</sup>  
(SFB TACO)

University of Vienna,  
Faculty of Physics and Center for Computational Materials Sciences,  
Kolingasse 14-16, 1090, Vienna, Austria

(\*Electronic mail: georg.kresse@univie.ac.at)

(\*Electronic mail: bernhard.schmiedmayer@univie.ac.at)

(Dated: 1 May 2024)

We develop a strategy that integrates machine learning and first-principles calculations to achieve technical accurate predictions of infrared spectra. Specifically, the methodology allows to predict infrared spectra for complex systems at finite temperatures. The method’s effectiveness is demonstrated in challenging scenarios, such as the analysis of water and the organic-inorganic halide perovskite MAPbI<sub>3</sub>, where our results consistently align with experimental data. A distinctive feature of the methodology is the incorporation of derivative learning, which proves indispensable for obtaining accurate polarization data in bulk materials and facilitates the training of a machine learning surrogate model of the polarization adapted to rotational and translational symmetries. We achieve polarisation prediction accuracies of about 1 % by training only on the predicted Born effective charges.

## I. INTRODUCTION

Infrared (IR) spectroscopy is an indispensable methodology for discerning local structures and the intricacies of functional groups within a material. It finds pervasive use across various scientific disciplines, firmly establishing itself as a key analytical instrument<sup>1</sup>. Notably, within the realm of catalysis research, IR spectroscopy is an important technique with the unique capability to elucidate surface structures and mechanisms at the molecular scale *in situ*<sup>2-4</sup>. Beyond catalysis, the application of IR spectroscopy extends its reach into the domains of clinical and biomedical analyses<sup>5-7</sup>. In the field of materials science, IR spectroscopy is also highly valuable, as it provides indirect insight on the bonding properties of the material<sup>8-10</sup>.

Although experimental IR spectroscopy is a mature technique, it is not always a simple matter to relate an experimental IR spectrum to a specific local structural feature. Prior knowledge is almost ubiquitously required to relate experimental data to structures (fingerprints). For instance, the frequencies of certain functional groups often relate to the oxidation state of the group or the electronegativity of the adsorption site. Computer simulations, particularly first-principles (FP) methods, have played a vital role in establishing these relationships. However, such FP calculations are very often limited to zero-temperature simulations<sup>11-13</sup>, whereas *in-situ* observations of catalytic reactions almost always involve finite temperature. In this study, we harness recent advancements in machine learning (ML) techniques in combination with FP calculations to develop a method that yields highly accurate computational IR spectra at finite temperatures.

In recent years, ML techniques have emerged as a potent new avenue in computational materials sciences<sup>14-18</sup>.

Our study uses an on-the-fly learning method<sup>19-22</sup> to generate transferable machine learned force fields (MLFFs). The first step is thus to harness MLFF to attain the required nano-second-long high-quality molecular dynamics (MD) trajectories. In the present work, we rely on a now fairly standard approach for the MLFF that uses rotational and translationally invariant descriptors and a kernel-based regression, although more refined approaches could be readily adopted<sup>23-25</sup>. Learning the polarization, though, is a more challenging task. The first approach to learn tensorial quantities dates back to  $\lambda$ -SOAP<sup>26-28</sup> that we have also decided to adapt in the present work. Equivariant message-passing networks would be equally suitable, but they are matter of fact closely related to the tensorial descriptors in  $\lambda$ -SOAP.

The second key issue that we address in the present work is that the polarization is not uniquely determined in bulk materials<sup>29-32</sup>. This is in stark contrast to molecules, where the polarization can be determined by appropriate integration of the density<sup>33</sup>. Direct learning of the polarization hence requires some curation of the polarization data *e.g.* deriving the polarization from Wannier-centres imposing some continuity condition<sup>34</sup> or manual “alignment” of the polarization data.

The solution presented in this study to overcome this challenge involves employing derivative learning<sup>35</sup>. Specifically, we utilize the derivative of the polarization with respect to the ionic positions, the Born effective charge tensors. The hypothesis is that by computing this derivative information across numerous structures (along any relevant “adiabatic” pathway), it becomes feasible to determine the anti-derivative, *i.e.*, the bulk polarization. This methodology offers a second crucial advantage: akin to the construction of MLFF where learning the forces proved pivotal, the Born effective charges are significantly more expressive yet nearly as computationally efficient to calculate in solid-state codes as a single polarization.

In the next section, we summarize our methodological approach, then demonstrate the feasibility of derivative-based

<sup>a)</sup>Also at VASP Software GmbH, Sensengasse 8, 1090 Vienna, Austria

learning for the water dimer, and discuss results for liquid water, where we find excellent agreement with the experiment only for a functional including van der Waals corrections. Lastly, we demonstrate very good agreement for the IR spectrum of an organic perovskite with experimental data. We finish with discussions and conclusions, and mention points worthwhile an improvement.

## II. METHOD

### A. General remarks:

In general, the energy of a molecule or material in the presence of an electric field is described by

$$E(\mathbf{x}, \mathcal{E}) = E_{\text{KS}}(\mathbf{x}) - \mathcal{E} \cdot \mathbf{P}(\mathbf{x}, \mathcal{E}) \quad (1)$$

where  $\mathcal{E}$  is the electric field,  $\mathbf{P}$  the polarization, and  $E_{\text{KS}}$  is the Kohn-Sham energy at zero field for the atoms at the position  $\mathbf{x}$ <sup>36</sup>. There is an implicit dependence of the Kohn-Sham energy on the electronic field, as the orbitals need to be determined by minimizing the energy in the presence of the field, but thanks to the variational properties and the Hellman-Feynman theorem, variations of the orbitals can be neglected for first derivatives. To calculate second derivatives, only the first derivatives of the orbitals are required. Obviously, the first derivative of the energy with respect to the electric field yields the polarization  $\mathbf{P}$ . The second derivative of the energy with respect to the field corresponds to the electronic polarizability, and the second mixed derivative with respect to the electric field and the positions yields the Born effective charge tensor:

$$\mathbf{Z}^* = \frac{\partial^2 E(\mathbf{x}, \mathcal{E})}{\partial \mathbf{x} \partial \mathcal{E}}. \quad (2)$$

This is a second-rank cartesian tensor. In the present work, we set out to learn the polarization as a function of the positions and neglect the dependence of the polarization on the electric field (implicitly assuming zero electric field). If we were to evaluate the energy using Eq. (1) this would only be correct to linear order in the field.

### B. Green-Kubo relation:

Using the Green-Kubo formalism, the ionic contribution to the polarizability, denoted as  $\chi(\omega)$ , is directly proportional to the Fourier transform of the autocorrelation function of the polarization  $\mathbf{P}$  and its time derivative and can be expressed as follows (SI units)<sup>37–39</sup>:

$$\chi_{\mu, \nu}(\omega) = \frac{\beta}{V \epsilon_0} \int_0^T \langle \mathbf{P}_\mu(0) \dot{\mathbf{P}}_\nu(t) \rangle e^{-i(\omega - i\delta)t} dt. \quad (3)$$

Here  $\mu$  and  $\nu$  represent Cartesian indices,  $V$  is the volume,  $\beta$  is the inverse temperature,  $\epsilon_0$  the vacuum permittivity,  $\omega$  the vibrational frequency, and  $\delta$  denotes a complex shift causing

a Lorentzian broadening. It is necessary that the product  $T\delta$  is small to avoid truncation artefacts. The time derivative of the polarization  $\dot{\mathbf{P}}$  can be written as

$$\frac{\partial \mathbf{P}}{\partial t} = \frac{\partial \mathbf{P}}{\partial \mathbf{x}} \frac{\partial \mathbf{x}}{\partial t} = \mathbf{Z}^* \dot{\mathbf{x}}. \quad (4)$$

To determine Eq. (3) accurately three prerequisites exist: first, accurate velocities  $\dot{\mathbf{x}}$  *i.e.* high-quality MD trajectories to describe the time evolution of the system, second, long simulation times  $T$ , and third, a reliable method to model the polarization and the Born effective charge tensors.

In infrared reflectivity or transmission experiments, an absorbance  $\alpha$  is measured and specified (the absorbance is proportional to  $\omega$  times the spectral function of the polarizability). We report the product of the absorbance  $\alpha(\omega)$  and the refractive index  $\eta(\omega)$ , isotropically averaged over the three diagonal components of the polarizability tensor:

$$\alpha(\omega) \eta(\omega) = \frac{\beta}{3V c \epsilon_0} \Re \int_0^T \langle \dot{\mathbf{P}}(0) \cdot \dot{\mathbf{P}}(t) \rangle e^{-i(\omega - i\delta)t} dt. \quad (5)$$

In this equation,  $c$  is the speed of light. We note that one can also autocorrelate  $\mathbf{P}(0)$  with  $\mathbf{P}(t)$ , or  $\mathbf{P}(0)$  with  $\dot{\mathbf{P}}(t)$ , adding simultaneously factors  $\omega^2$  and  $-i\omega$ , respectively. We have tested all three approaches and found identical results for the three versions.

### C. Molecular dynamics simulations:

For disordered materials, the required long simulation times are unattainable using first principles simulations, and hence surrogate models are required. To obtain FP data for constructing the MLFF, we employed the Vienna *Ab initio* Simulation Package (VASP)<sup>40–42</sup>. Both the training of the MLFF and its subsequent application were conducted using the ML framework integrated in the VASP code<sup>19–21</sup>. It is essential to emphasize that the quality of an MLFF is not only dependent on the underlying ML algorithm but also the quality and representativeness of the training dataset concerning the problem at hand<sup>43,44</sup>. Ideally, the training dataset should comprehensively cover all relevant regions of the potential energy surface, and this coverage should be compact to minimize the necessity for computationally expensive FP calculations. To accomplish this, we have employed the on-the-fly learning scheme integrated within VASP.

The training dataset was curated from a collection of multiple MD trajectories. These trajectories were obtained using Langevin thermostats<sup>45–47</sup> with a friction coefficient of 10 ps<sup>-1</sup>. To comprehensively cover all phases of interest, a series of heating and cooling runs were executed as discussed in the supplementary. Notably, the temperature range considered was slightly wider than the region of interest, a choice aimed at enhancing the stability of the MLFF.

To obtain an IR spectrum, we used multiple MD trajectories within a micro-canonical ensemble to avoid artefacts caused by thermostats. The initial configurations and velocities for these individual MD trajectories were drawn from an

isothermal-isobaric ensemble. The final IR spectrum was then computed as the average of the individual IR spectra obtained from these trajectories, providing a statistically accurate representation of the system’s vibrational modes. We note that IR spectra are only sensitive to the long-range vibrational modes, *i.e.* at zero temperature the calculations can be performed using the unit cell. However, at finite temperature, supercells are required to account for the finite temperature disorder. We checked carefully that the IR spectra are cell-size converged for all cases reported here.

#### D. Polarization model

As highlighted in the introduction, obtaining the polarization for bulk systems presents inherent challenges. This is due to an undetermined modulo resulting from the absence of a unique phase origin or reference point<sup>29–32</sup>. To address this problem, we use derivative learning. Specifically, the polarization  $\mathbf{P}$  of a system can be understood as the antiderivative of the Born effective charge tensor  $\mathbf{Z}^*$ <sup>48</sup>. The Born effective charge tensor for ion  $i$  is mathematically expressed as

$$\mathbf{Z}_{i,\alpha\beta}^* = Z_{J(i)} \delta_{\alpha,\beta} + V \left. \frac{\partial \mathbf{P}_\alpha^{\text{elec}}}{\partial \mathbf{x}_{i\beta}} \right|_{\mathcal{E}=0}. \quad (6)$$

Here,  $Z_{J(i)}$  represents the bare ionic charge of the  $i$ th ion and  $J(i)$  signifies the atomic species,  $V$  stands for the volume of the unit cell, and  $\mathbf{P}_\alpha^{\text{elec}}$  corresponds to the Cartesian component  $\alpha$  of the macroscopic electronic polarization. In the equation,  $\mathbf{x}_{i\beta}$  denotes the position of the  $i$ th ion along the Cartesian component  $\beta$ .

The central idea is that the surrogate ML model describes the polarization  $\mathbf{P}_\alpha^{\text{elec}}$ . However, we aim to avoid training this model on the polarization, but instead train the model on derivative data  $\mathbf{Z}^*$ . Clearly, the polarization must transform like a vectorial quantity under rotations. To this end, we employ ridge regression, with a linear kernel function  $K$ , constructed using the covariance of 3-dimensional descriptors  $D_n^\mu(\mathcal{X})$ . Here,  $\mathcal{X}$  represents an atomic environment,  $\mu$  corresponds to a Cartesian component, and  $n$  is the feature dimension. The linear kernel function is defined as follows:

$$K_{\mu\nu}(\mathcal{X}, \mathcal{X}') = \sum_n D_n^\mu(\mathcal{X}) D_n^\nu(\mathcal{X}'). \quad (7)$$

This equation captures the similarity between atomic environments  $\mathcal{X}$  and  $\mathcal{X}'$ . To describe the atomic environment, we utilized the  $\lambda$ -SOAP (Smooth Overlap of Atomic Potentials) descriptors developed by Grisafi *et al.*<sup>26</sup>. These descriptors conserve the rotational symmetry of tensorial quantities:

$$\hat{S} D_n^\mu(\mathcal{X}) = D_n^\mu(\hat{S} \mathcal{X}), \quad (8)$$

where  $\hat{S}$  is a generalized symmetry operator of the  $SO(3)$  group. The descriptor is tailored to describe the surroundings of an atom. In the present work we use two- and three-body descriptors and Bessel functions as radial basis sets, as in the original MLFF implementation of VASP<sup>19,20</sup>. To ensure

smoothness in derivatives and avoid abrupt discontinuities, we incorporated a Behler and Parrinello cutoff function<sup>14</sup>. The radial cutoffs are typically set to 5.5 Å.

The polarization  $\mathbf{P}$  of a configuration, characterized by atomic environments  $\mathcal{X}_j$ , is determined using descriptors  $D_n^v(\mathcal{X}_{I_{\text{ref}}})$  of reference atomic environments  $\mathcal{X}_{I_{\text{ref}}}$  following the equation:

$$\mathbf{P}_\alpha = \sum_{v_{I_{\text{ref}}}} \omega_{I_{\text{ref}}}^v \sum_{jn} D_n^\alpha(\mathcal{X}_j) D_n^v(\mathcal{X}_{I_{\text{ref}}}). \quad (9)$$

Here  $\omega_{I_{\text{ref}}}^v$  represents weights that are paired with each reference descriptor  $D_n^v(\mathcal{X}_{I_{\text{ref}}})$ . These weights are determined through derivative learning of the Born effective charge tensor, as given by the equation:

$$\mathbf{Z}_{\alpha\beta}^*(i) = \sum_{v_{I_{\text{ref}}}} \omega_{I_{\text{ref}}}^v \sum_{jn} \frac{\partial D_n^\alpha}{\partial \mathbf{x}_{i\beta}}(\mathcal{X}_j) D_n^v(\mathcal{X}_{I_{\text{ref}}}). \quad (10)$$

The weights are obtained by utilizing a linear regression model using the least squares method, allowing for simultaneous calculation of  $\omega_{I_{\text{ref}}}^v$  across all training configurations. We use sparse regression, *i.e.* reduce the number of kernel-basis functions  $\mathcal{X}_{I_{\text{ref}}}$  using farthest point sampling.

Our investigation revealed a substantial improvement in the quality of the fit by specially treating the diagonal elements of the Born effective charge tensor. To achieve this improvement, we applied a preprocessing step that involved subtracting the mean value of the diagonal elements of the Born effective charge tensor for each atomic species  $\bar{Z}_J^*$  before the training process. Here,  $J$  represents the atomic species. Subsequently, we add back to the polarization  $\mathbf{P}$  the following term:

$$\bar{\mathbf{P}}_\alpha = \sum_{i=1}^{N_{\text{atom}}} \bar{Z}_{J(i)}^* \mathbf{x}_{i\alpha}. \quad (11)$$

The computation of Born effective charges was carried out using density functional perturbation theory (DFPT)<sup>49–51</sup> by computing the static ion-clamped dielectric matrix, as discussed by Baroni and Resta<sup>52</sup> and Gajdoš *et al.* for the projector-augmented wave (PAW) method<sup>50</sup>. It is noteworthy that the calculation of the Born effective charge tensor is computationally three times more intensive than a comparable Density Functional Theory (DFT) groundstate calculation. Essentially, it requires calculating the orbitals’ first derivative with respect to the three directions of the external fields  $\mathcal{E}$ , which suffices to obtain the mixed derivatives as defined in Eq. (2). Alternatively, one could also determine the first derivative of the orbitals with respect to the positions, and then predict the mixed second derivatives, however, this scales linearly with system size and is thus computationally more involved (the results for the Born effective charges are independent of the order of differentiation). To the best of our knowledge, any electronic structure code can determine the Born-effective charges via the orbital derivative with respect to the field. The present approach is therefore applicable to most first-principles codes. The inclusion of derivatives provides considerably more information, in particular, an addi-

tional  $3N_{\text{atom}}$  of data, so that only a small number of FP calculations are required to achieve a high degree of accuracy as demonstrated below.

## E. Results and Discussion

To demonstrate the versatility of the developed methodology, we have applied it to three distinct systems. First, we consider the water dimer  $2(\text{H}_2\text{O})$ . In the case of molecules, the polarization is a well-defined property. Second, we examine water. Lastly, we extend our analysis to a complex solid state system, focusing on an organic perovskite  $\text{MAPbI}_3$ , since it exhibits sizable anharmonicities.

## F. Water dimer

We start our analysis with a water dimer  $2(\text{H}_2\text{O})$ . Given that molecules possess a well-defined polarization, we can make a comparative assessment between the derivative learning approach and the conventional method of directly learning the polarization. This provides a valuable benchmark for assessing the reliability of the methodology.

The training and test configurations were extracted from an MD trajectory. We commenced with a water dimer placed in a simulation box with ample vacuum space. One of the oxygen atoms was constrained using selective dynamics to anchor the system. To simulate thermalization, we executed a heating run, spanning temperatures from 10 K to 320 K, utilizing a Langevin thermostat. Importantly, the MD run was enhanced with on-the-fly machine learning to speed up the computational process. For the FP calculations, we opted for a revised Perdew-Burke-Ernzerhof (RPBE) functional<sup>53</sup> with van der Waals (vdW) dispersion energy corrections of Grimme *et al.* with zero-damping function<sup>54</sup> (RPBE-D3). This choice ensured that our calculations account for vdW interactions.

During the MD, the mass of the hydrogen atom was increased to 8 u, to allow for larger time steps of 1.5 fs. Out of a total of 200 000 MD time steps, we selected 1000 uniformly distributed configurations for the computation of the polarization. For calculating the Born effective charge tensor 150 configurations were chosen. To determine polarization, we integrated the total (electronic and ionic) charge times the position operator by switching on dipole corrections in VASP (see *e.g.* Refs.<sup>55,56</sup>). To improve the accuracy of the Born effective charge tensor, we used a strict convergence criterion of  $1 \times 10^{-7}$  eV for the electronic self-consistency loop and large cells. To confirm the internal consistency between polarization and Born effective charges, we also calculated numerical derivatives of the polarization. These derivatives showed an excellent agreement with the Born charges with a root mean square error (RMSE) less than  $5 \times 10^{-6} |e|$ , confirming the reliability of the resulting database.

To construct the learning curves, we used a dataset consisting of 1000 polarization calculations. The dataset was split in half by alternately selecting configurations for training and validation. Additionally, we selected 50 Born effective charge

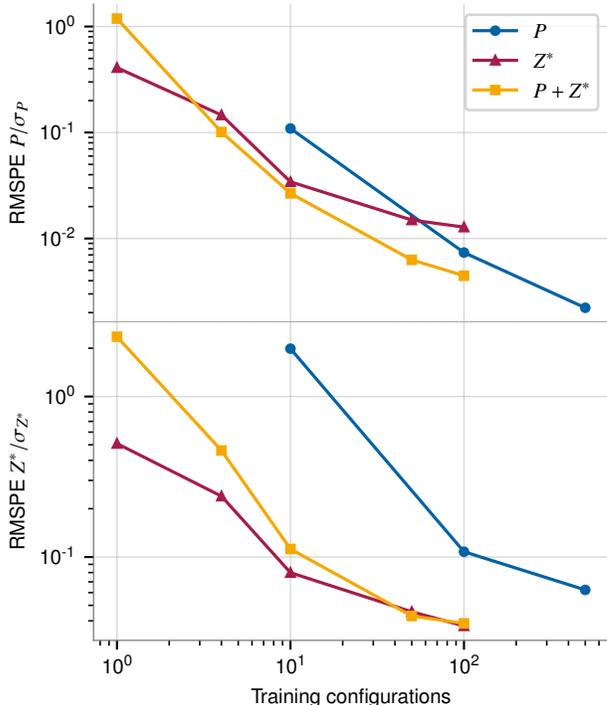


FIG. 1. Learning curves for the polarization ( $P$ ) and Born effective charges ( $Z^*$ ) of a water dimer. Errors on the test set are expressed in root mean square percentage errors (RMSPE), where the RMSE is normalized by the standard deviation ( $\sigma$ ) of the test set to yield a dimensionless quantity (note that  $10^{-1}$  and  $10^{-2}$  at the y-axis correspond to 10% and 1%, respectively). The blue line (circle) denotes training solely on polarization ( $P$ ), and the red line (triangle) represents training solely on Born effective charges ( $Z^*$ ). The yellow line (square) demonstrates combined learning of polarization and Born effective charges ( $P + Z^*$ ).

calculations as validation data, chosen to be uniformly distributed across the dataset of 150 Born effective charge tensor calculations. The remaining 100 Born effective charge calculations served as the training dataset. During the combined training process, the polarization data was weighted ten times higher than the Born effective charge data. To determine the optimal number of fitting parameters, we selected the number of kernel functions that minimized the error in polarization prediction. It is worth noting that, as proposed by Cortes *et al.* for regression<sup>57</sup>, there should be a relationship between the RMSE and the number of training data, characterized by a power-law decay. The learning curves, which provide insights into the model’s performance, are presented in Fig. 1.

The learning curves demonstrate the power-law relation between the RMSE and the number of training configurations. However, it is important to note that with a large number of training data points, the improvement in the mean squared errors seems to plateau somewhat, indicating that our linear regression with two- and three-body descriptors will likely yield some residual (but acceptable) model errors. Overall, the most favourable results are achieved through combined learning, with percentage errors below 0.5% and 3% for the polariza-

tion and the Born effective charges, respectively. Crucially, learning only Born charges from 100 training configurations also attains a high relative accuracy of approximately 1 % for the polarization, and is as accurate as combined learning for the Born effective charges. We also note that there is no noticeable offset in the ML polarization compared to the FP calculations. Likely this is so since the molecule is free to rotate (and does rotate during the MD) and the surrogate model reliably determines the offset. On the other hand, attempting to train on the polarization data only requires more training data (blue line), but even with 500 training data points, the Born effective charges still show twice the errors as combined training does using 100 training data points.

### G. Liquid water

For our second proof-of-concept, we selected liquid  $\text{H}_2\text{O}$  at room temperature. While the IR spectrum of water has been extensively studied and is well understood<sup>58–63</sup>, theoretical interpretations of these spectra remain challenging<sup>64–71</sup>. Even with the use of modern DFT methods, accurately reproducing the experimental properties of water is a complex task and often yields results that are far from accurate<sup>72</sup>. As a result, this system serves as an ideal test case for validating the reliability and effectiveness of the methodology as well as gleaning some insight on the underlying dynamics of water.

We chose the training configurations from an MD trajectory using an on-the-fly learning scheme. The simulation was conducted within a cubic box with periodic boundary conditions, containing a total of 64  $\text{H}_2\text{O}$  molecules. The lattice constant of the cubic box was adjusted to attain a density of approximately  $997 \text{ kg m}^{-3}$ , closely resembling the density of water at room temperature<sup>73</sup>. Data was collected during multiple heating and cooling runs conducted within a canonical MD ensemble. We employed various temperatures, spanning from 270 K to 420 K. The MLFF was trained using FP data obtained from DFT calculations. Specifically, we again employed the RPBE-D3 functional but found it necessary to use hard PAW potentials and a cutoff of 800 eV to obtain accurate stretch frequencies<sup>74</sup>. Additionally, we trained an MLFF using SCAN<sup>75</sup> to determine whether SCAN offers a reasonable description of the water dynamics.

The training configurations for the tensorial machine learning framework were once more chosen from a canonical MD ensemble. This ensemble was maintained at a constant temperature of 298.2 K and controlled by a Langevin thermostat. A total of 10000 MD steps, accelerated by the MLFF, were executed. From the MD trajectory, 100 configurations were uniformly selected as the training dataset for the Born effective charges. A scatter plot of the trained and predicted Born effective charges is shown in Fig. 2.

The computational IR spectrum presented alongside experimental data in Fig. 3 was computed by averaging the results over 20 individual IR spectra. Each of these individual spectra was calculated from a micro-canonical MD trajectory. For each MD run, we initiated the simulation from an uncorrelated starting configuration, ensuring the appropriate aver-

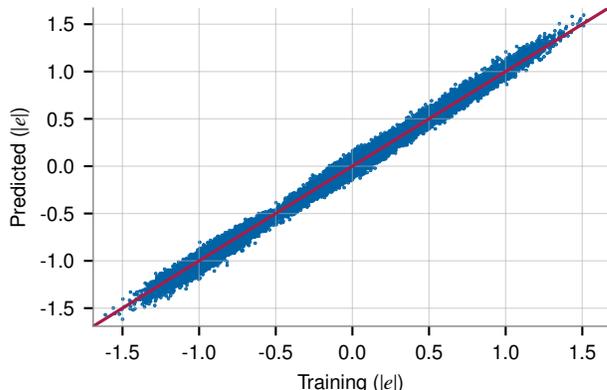


FIG. 2. Scatter plot of the trained and predicted Born effective charges for water.

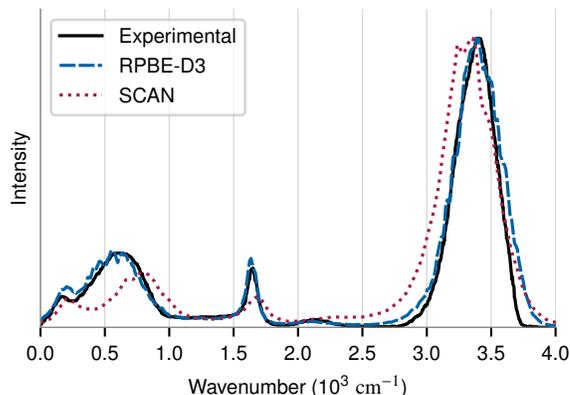


FIG. 3. Experimental and computational IR spectra for liquid water. The experimental reference data is from Ref.<sup>62</sup>

age corresponding to room temperature. In each run, we performed a total of 100 000 MD steps, with a time step of 0.25 fs. This strategy of conducting multiple calculations from uncorrelated starting configurations was employed to enhance the statistical accuracy of our results. Instead of Lorentzian broadening, a Gaussian filter was applied before performing the Fourier transform. This improves the feature sharpness in the computed IR spectra.

As illustrated in Fig. 3, the present methodology allows the computation of the IR spectrum of liquid water with a remarkable level of agreement with experimental data. There are two important conclusions to draw from this result. That the intensity follows so closely the experimental data, is related to an accurate description of the Born effective charges. Since we train on the Born effective charges, the good agreement is likely not astonishing. Second, RPBE+D3 yields an excellent description of the dynamics of liquid water, both for the high-frequency stretch, but also for the medium-frequency bending motions. Most important are the lower frequency modes that are related to intra-molecular motions. It is important to

note that the results for the SCAN functional are significantly worse for these modes. Specifically, the low-frequency mode around  $500\text{ cm}^{-1}$  is wrong by almost 40% indicating serious deficiencies in the description of the molecular motion of water molecules en-caged by the four surrounding water molecules.

Previous works, such as those by Sommers *et al.*<sup>28</sup>, Zhang *et al.*<sup>34</sup>, and Gastegger *et al.*<sup>33</sup>, have undertaken similar approaches utilizing MD trajectories in conjunction with symmetry-preserving machine learning frameworks to obtain Raman spectra for water and IR spectra for molecules. In these prior studies, the focus has been on learning the polarization directly. This required significantly more training data and a careful calculation of the polarization in order to avoid any discontinuities.

## H. Anharmonic solids

Our final test system is the organic-inorganic halide perovskite, methylammonium lead iodide (MAPbI<sub>3</sub>). This material has been the subject of numerous experimental and theoretical studies, including state-of-the-art vibrational studies<sup>20,73,76–82</sup>. Notably, MAPbI<sub>3</sub> undergoes three entropy-driven phase transitions: from an orthorhombic phase to a tetragonal phase at 160 K, and from the tetragonal phase to a cubic phase at 330 K. The thermodynamic nature of this material makes it challenging to model the IR spectra of the tetragonal phase using traditional 0 K methods like DFPT<sup>49–51</sup>. Therefore, MAPbI<sub>3</sub> serves as a valuable test case for validating the reliability of the scheme for strongly anharmonic solids.

The training process for the MLFF applied to MAPbI<sub>3</sub> closely mirrored the methodology employed for water and previous studies of MAPbI<sub>3</sub><sup>20</sup>. Multiple heating runs, encompassing the two phases— orthorhombic, and tetragonal—were conducted using Langevin thermostat-driven MD simulations. The temperature range spanned from 80 K to 430 K. Initially, fixed cell volumes were used, followed by simulating an isothermal-isobaric ensemble using the Parrinello-Rahman method<sup>83,84</sup>. The training was performed using a strongly-constrained and appropriately normed (SCAN) meta-gradient corrected functional<sup>75</sup>. We note that we found in previous studies that SCAN is better suited for the simulation of MAPbI<sub>3</sub><sup>85</sup> than say RPBE-D3, as RPBE-D3 does not account for screening of the vdW interactions by the strongly polarizable cage atoms.

In Fig. 4, we present the IR spectra for the orthorhombic phase at 107 K and the tetragonal phase at 228 K, alongside experimental results for comparison. The spectrum of the well-ordered orthorhombic phase and the tetragonal phase was calculated using a  $4 \times 4 \times 4$  supercell for better statistics and to allow for some reordering of the MA molecules. The starting configurations for the individual MD runs were chosen from an isothermal-isobaric ensemble. Aside from using starting configurations with varying cell vectors, the procedure for calculating the IR spectrum closely mirrors the one described for water.

The analysis of the computed IR spectra reveals excellent agreement with the experimental data but also some small discrepancies, particularly, in the intensities of the individual peaks. Notably, the peaks around  $900\text{ cm}^{-1}$  appear with higher intensity. The experimental reference spectra, obtained from a single crystal<sup>82</sup>, may also be influenced by surface effects and crystal structure orientation, contributing to variations in intensity between computational and experimental results.

We start with a comparison for the orthorhombic low-temperature phase (107 K). The agreement between the DFPT and the finite temperature simulation is very good, but there are some marked improvements in the finite temperature data. The first peak around  $900\text{ cm}^{-1}$  shows a double peak in both the DFPT and FT simulation, in agreement with the experiment. The peak at  $970\text{ cm}^{-1}$  is completely missing in the DFPT data but visible in the finite temperature data. It is a result of anharmonic interactions. The shoulder at  $1420\text{ cm}^{-1}$  is pronounced using DFPT but washed out at finite temperature. This peak corresponds to the CH<sub>3</sub>-bending motions<sup>51</sup>. In the supplementary, we show that using the ionic charges only, this peak is visible in the spectrum, but the intensity of the peak is overestimated. Calculating the electronic contribution (el) only shows an almost identical peak, however, the phase of the electronic polarization is opposite to the ionic contribution, so when autocorrelating the sum of the electronic and ionic dipoles, the peak is strongly suppressed. This brings better overall agreement with the experiment, where the peak is also weak, but likely our Born effective charges are not quite sufficiently accurate (linear regression). The physics behind the vanishing of the peak is fairly simple: in this particular mode the H atoms move orthogonal to their bond direction, a direction in which the total Born effective charges  $\mathbf{Z}^*$  are small to start with (electronic contribution cancels ionic one). Furthermore, the movement of three hydrogen atoms is concerted being close to a helicopter motion which reduces the IR intensity further.

In the experiment, we see quite some intensity between  $1450\text{ cm}^{-1}$  and  $1600\text{ cm}^{-1}$ , which is also nicely reproduced by the FT simulations. We note that this frequency range becomes even more populated in the tetragonal phase in the FT simulations, and this population is a result of the strongly anharmonic rattling motion of the molecules in the cage in turn affecting the bending motion of the hydrogens.

The tetragonal phase spectrum (228 K) is also in excellent agreement with the experimental spectrum. We note that the DFPT spectrum now shows many deficiencies, with a complete lack of peaks at  $950\text{ cm}^{-1}$  and a tiny peak around  $1250\text{ cm}^{-1}$ , as well as sharp features around  $1580\text{ cm}^{-1}$ . The FT data resolves these issues, albeit the two main peaks around  $900\text{ cm}^{-1}$  and  $1480\text{ cm}^{-1}$  are somewhat too broad and washed out, and again the peak around  $1480\text{ cm}^{-1}$  is suppressed by the electronic contribution.

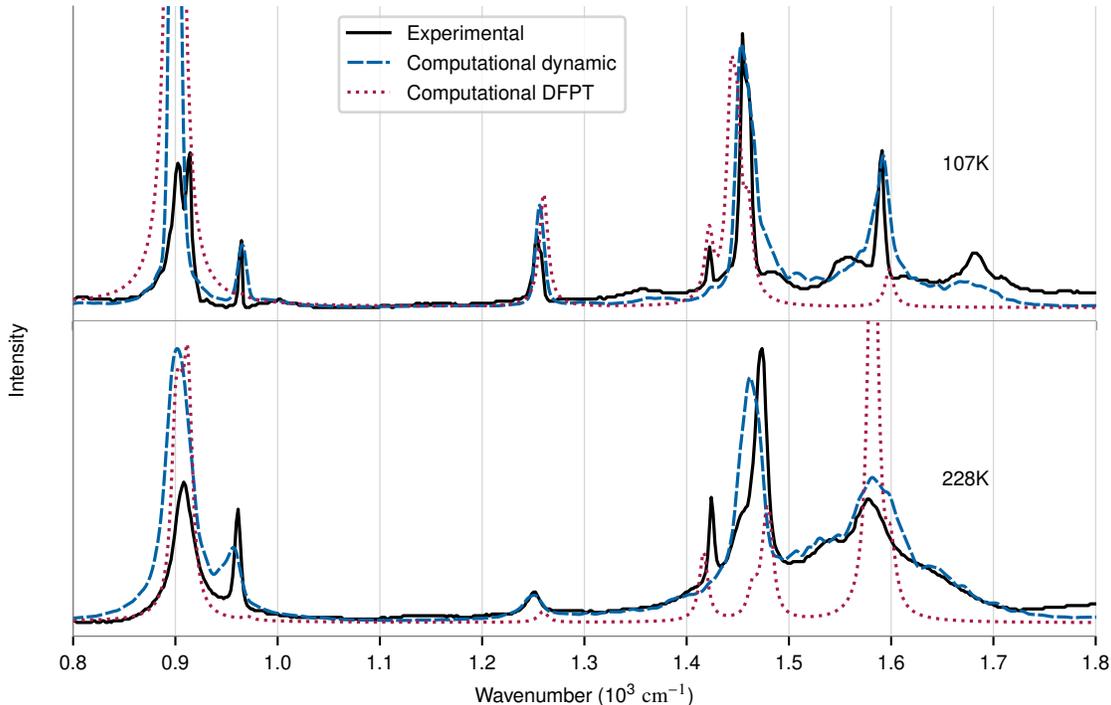


FIG. 4. Experimental and computational IR spectra for the orthorhombic and tetragonal phases of  $\text{MAPbI}_3$ . The vibrational frequencies of the computational IR spectra have been red-shifted by 1.5 % for alignment with experimental data. The experimental reference data is from Ref.<sup>82</sup>.

### III. CONCLUSION

The present study showcases the effectiveness of a computational framework combining first-principles simulations with machine-learning methodologies. Machine-learned force fields can be used to access timescales that were previously very difficult to attain. This allows us to obtain vibrational spectra with excellent statistical accuracy, which would be very expensive to calculate without force fields. The main advance of the present work is, however, that we learn the polarization from its derivative the Born effective charges. As mentioned in the main text, in VASP— but likely so in any plane wave-based code—the calculation of the Born-effective charges via the first derivative of the orbitals with respect to external fields is only roughly three times more costly than a groundstate Kohn-Sham calculation. Learning of force fields using first-principle nuclear derivatives is now ubiquitous, and learning the polarization from its nuclear derivative is a natural extension to this idea. Crucially, we have shown that the inclusion of polarization data, which is difficult to determine without some arbitrary modulus, is not required. As the polarization is the anti-derivative of the Born-effective charges, it can be directly calculated from the machine-learning model as long as derivative data is supplied along all relevant adiabatic pathways. We successfully applied this approach to challenging scenarios, including water and the organic-inorganic halide perovskite  $\text{MAPbI}_3$ . In both cases, our method demonstrates excellent agreement with experimental results, high-

lighting its capacity to capture the vibrational properties of diverse materials.

We finish with a few comments on further developments. In the present work, we have only used linear regression with two-body and three-body descriptors. Although the prediction accuracies are good (condensed matter) to excellent (molecules), we feel that the inclusion of higher body-order terms, or non-linear kernel-based regression might improve the predicted Born-effective charges further. Furthermore, the subtraction of a diagonal component from the Born-effective charges, albeit not particularly cumbersome, seems somewhat unsatisfactory, and one would prefer to avoid it.

Overall, the present methodology is already very robust and can be readily applied to many relevant problems. Further research could be directed towards infrared simulations of more complex adsorbates on surfaces or of water interacting with electrodes.

### IV. ACKNOWLEDGEMENTS

This research was funded in whole by the Austrian Science Fund (FWF) 10.55776/F81. For open access purposes, the author has applied a CC BY public copyright license to any author accepted manuscript version arising from this submission. The computational results presented have been achieved in part using the Vienna Scientific Cluster (VSC).

- <sup>1</sup>B. H. Stuart, *Infrared spectroscopy: fundamentals and applications* (John Wiley & Sons, 2004).
- <sup>2</sup>A. Vimont, F. Thibault-Starzyk, and M. Daturi, “Analysing and understanding the active site by ir spectroscopy,” *Chemical Society Reviews* **39**, 4928–4950 (2010).
- <sup>3</sup>Y. J. Chabal, “Surface infrared spectroscopy,” *Surface Science Reports* **8**, 211–357 (1988).
- <sup>4</sup>J. Ryczkowski, “Ir spectroscopy in catalysis,” *Catalysis Today* **68**, 263–381 (2001).
- <sup>5</sup>A. Barth, “Infrared spectroscopy of proteins,” *Biochimica et Biophysica Acta (BBA)-Bioenergetics* **1767**, 1073–1101 (2007).
- <sup>6</sup>L. M. Ng and R. Simmons, “Infrared spectroscopy,” *Analytical chemistry* **71**, 343–350 (1999).
- <sup>7</sup>J. Luypaert, D. Massart, and Y. Vander Heyden, “Near-infrared spectroscopy applications in pharmaceutical analysis,” *Talanta* **72**, 865–883 (2007).
- <sup>8</sup>T. Theophile, *Infrared spectroscopy: Materials science, engineering and technology* (BoD–Books on Demand, 2012).
- <sup>9</sup>L. Fernández-Carrasco, D. Torrens-Martín, L. Morales, and S. Martínez-Ramírez, “Infrared spectroscopy in the analysis of building and construction materials,” *Infrared spectroscopy—Materials science, engineering and technology* **510** (2012).
- <sup>10</sup>S. M. Silva, C. R. Braga, M. V. Fook, C. M. Raposo, L. H. Carvalho, and E. L. Canedo, “Application of infrared spectroscopy to analysis of chitosan/clay nanocomposites,” *Infrared spectroscopy—materials science, engineering and technology*, 43–62 (2012).
- <sup>11</sup>M. Biczysko, J. Bloino, and C. Puzzarini, “Computational challenges in astrochemistry,” *Wiley Interdisciplinary Reviews: Computational Molecular Science* **8**, e1349 (2018).
- <sup>12</sup>T. L. Jansen, “Computational spectroscopy of complex systems,” *The Journal of Chemical Physics* **155** (2021).
- <sup>13</sup>K. B. Beć, J. Grabska, and C. W. Huck, “Current and future research directions in computer-aided near-infrared spectroscopy: A perspective,” *Spectrochimica Acta Part A: Molecular and Biomolecular Spectroscopy* **254**, 119625 (2021).
- <sup>14</sup>J. Behler and M. Parrinello, “Generalized neural-network representation of high-dimensional potential-energy surfaces,” *Physical review letters* **98**, 146401 (2007).
- <sup>15</sup>A. P. Bartók, M. C. Payne, R. Kondor, and G. Csányi, “Gaussian approximation potentials: The accuracy of quantum mechanics, without the electrons,” *Physical review letters* **104**, 136403 (2010).
- <sup>16</sup>A. P. Bartók, J. Kermode, N. Bernstein, and G. Csányi, “Machine learning a general-purpose interatomic potential for silicon,” *Physical Review X* **8**, 041048 (2018).
- <sup>17</sup>L. Bonati and M. Parrinello, “Silicon liquid structure and crystal nucleation from ab initio deep metadynamics,” *Physical review letters* **121**, 265701 (2018).
- <sup>18</sup>T. Morawietz, A. Singraber, C. Dellago, and J. Behler, “How van der waals interactions determine the unique properties of water,” *Proceedings of the National Academy of Sciences* **113**, 8368–8373 (2016).
- <sup>19</sup>R. Jinnouchi, F. Karsai, and G. Kresse, “On-the-fly machine learning force field generation: Application to melting points,” *Physical Review B* **100**, 014105 (2019).
- <sup>20</sup>R. Jinnouchi, J. Lahnsteiner, F. Karsai, G. Kresse, and M. Bokdam, “Phase transitions of hybrid perovskites simulated by machine-learning force fields trained on the fly with bayesian inference,” *Physical review letters* **122**, 225701 (2019).
- <sup>21</sup>R. Jinnouchi, F. Karsai, C. Verdi, R. Asahi, and G. Kresse, “Descriptors representing two- and three-body atomic distributions and their effects on the accuracy of machine-learned inter-atomic potentials,” *The Journal of Chemical Physics* **152** (2020).
- <sup>22</sup>R. Jinnouchi, K. Miwa, F. Karsai, G. Kresse, and R. Asahi, “On-the-fly active learning of interatomic potentials for large-scale atomistic simulations,” *The Journal of Physical Chemistry Letters* **11**, 6946–6955 (2020).
- <sup>23</sup>R. Drautz, “Atomic cluster expansion for accurate and transferable interatomic potentials,” *Physical Review B* **99**, 014104 (2019).
- <sup>24</sup>S. Batzner, A. Musaelian, L. Sun, M. Geiger, J. P. Mailoa, M. Kornbluth, N. Molinari, T. E. Smidt, and B. Kozinsky, “E (3)-equivariant graph neural networks for data-efficient and accurate interatomic potentials,” *Nature communications* **13**, 2453 (2022).
- <sup>25</sup>I. Batatia, D. P. Kovacs, G. Simm, C. Ortner, and G. Csányi, “Mace: Higher order equivariant message passing neural networks for fast and accurate force fields,” *Advances in Neural Information Processing Systems* **35**, 11423–11436 (2022).
- <sup>26</sup>A. Grisafi, D. M. Wilkins, G. Csányi, and M. Ceriotti, “Symmetry-adapted machine learning for tensorial properties of atomistic systems,” *Physical review letters* **120**, 036002 (2018).
- <sup>27</sup>D. M. Wilkins, A. Grisafi, Y. Yang, K. U. Lao, R. A. DiStasio Jr, and M. Ceriotti, “Accurate molecular polarizabilities with coupled cluster theory and machine learning,” *Proceedings of the National Academy of Sciences* **116**, 3401–3406 (2019).
- <sup>28</sup>G. M. Sommers, M. F. C. Andrade, L. Zhang, H. Wang, and R. Car, “Raman spectrum and polarizability of liquid water from deep neural networks,” *Physical Chemistry Chemical Physics* **22**, 10592–10602 (2020).
- <sup>29</sup>R. King-Smith and D. Vanderbilt, “Theory of polarization of crystalline solids,” *Physical Review B* **47**, 1651 (1993).
- <sup>30</sup>D. Vanderbilt and R. King-Smith, “Electric polarization as a bulk quantity and its relation to surface charge,” *Physical Review B* **48**, 4442 (1993).
- <sup>31</sup>R. Resta, “Macroscopic polarization in crystalline dielectrics: the geometric phase approach,” *Reviews of modern physics* **66**, 899 (1994).
- <sup>32</sup>R. Resta, “Quantum-mechanical position operator in extended systems,” *Physical Review Letters* **80**, 1800 (1998).
- <sup>33</sup>M. Gastegger, J. Behler, and P. Marquetand, “Machine learning molecular dynamics for the simulation of infrared spectra,” *Chemical science* **8**, 6924–6935 (2017).
- <sup>34</sup>Y. Zhang, S. Ye, J. Zhang, C. Hu, J. Jiang, and B. Jiang, “Efficient and accurate simulations of vibrational and electronic spectra with symmetry-preserving neural network models for tensorial properties,” *The Journal of Physical Chemistry B* **124**, 7284–7290 (2020).
- <sup>35</sup>S. Chmiela, A. Tkatchenko, H. E. Sauceda, I. Poltavsky, K. T. Schütt, and K.-R. Müller, “Machine learning of accurate energy-conserving molecular force fields,” *Science advances* **3**, e1603015 (2017).
- <sup>36</sup>P. Umari and A. Pasquarello, “Ab initio molecular dynamics in a finite homogeneous electric field,” *Physical Review Letters* **89**, 157602 (2002).
- <sup>37</sup>R. Kubo, “Statistical-mechanical theory of irreversible processes. i. general theory and simple applications to magnetic and conduction problems,” *Journal of the physical society of Japan* **12**, 570–586 (1957).
- <sup>38</sup>R. Zwanzig, “Time-correlation functions and transport coefficients in statistical mechanics,” *Annual Review of Physical Chemistry* **16**, 67–102 (1965).
- <sup>39</sup>D. Sangalli, A. Marini, and A. Debernardi, “Pseudopotential-based first-principles approach to the magneto-optical kerr effect: From metals to the inclusion of local fields and excitonic effects,” *Physical Review B* **86**, 125139 (2012).
- <sup>40</sup>G. Kresse and J. Furthmüller, “Efficiency of ab-initio total energy calculations for metals and semiconductors using a plane-wave basis set,” *Computational materials science* **6**, 15–50 (1996).
- <sup>41</sup>G. Kresse and J. Furthmüller, “Efficient iterative schemes for ab initio total-energy calculations using a plane-wave basis set,” *Physical review B* **54**, 11169 (1996).
- <sup>42</sup>G. Kresse and D. Joubert, “From ultrasoft pseudopotentials to the projector augmented-wave method,” *Physical review b* **59**, 1758 (1999).
- <sup>43</sup>J. Li, B. Jiang, and H. Guo, “Permutation invariant polynomial neural network approach to fitting potential energy surfaces. ii. four-atom systems,” *The Journal of chemical physics* **139** (2013).
- <sup>44</sup>V. Botu and R. Ramprasad, “Adaptive machine learning framework to accelerate ab initio molecular dynamics,” *International journal of quantum chemistry* **115**, 1074–1083 (2015).
- <sup>45</sup>M. P. Allen and D. J. Tildesley, *Computer simulation of liquids* (Oxford university press, 2017).
- <sup>46</sup>W. G. Hoover, A. J. Ladd, and B. Moran, “High-strain-rate plastic flow studied via nonequilibrium molecular dynamics,” *Physical Review Letters* **48**, 1818 (1982).
- <sup>47</sup>D. J. Evans, “Computer “experiment” for nonlinear thermodynamics of couette flow,” *The Journal of Chemical Physics* **78**, 3297–3302 (1983).
- <sup>48</sup>C.-Z. Wang, R. Yu, and H. Krakauer, “Polarization dependence of born effective charge and dielectric constant in kno<sub>3</sub>,” *Physical Review B* **54**, 11161 (1996).
- <sup>49</sup>X. Wu, D. Vanderbilt, and D. Hamann, “Systematic treatment of displacements, strains, and electric fields in density-functional perturbation theory,” *Physical Review B* **72**, 035105 (2005).

- <sup>50</sup>M. Gajdoš, K. Hummer, G. Kresse, J. Furthmüller, and F. Bechstedt, “Linear optical properties in the projector-augmented wave methodology,” *Physical Review B* **73**, 045112 (2006).
- <sup>51</sup>M. A. Pérez-Osorio, R. L. Milot, M. R. Filip, J. B. Patel, L. M. Herz, M. B. Johnston, and F. Giustino, “Vibrational properties of the organic–inorganic halide perovskite  $\text{CH}_3\text{NH}_3\text{PbI}_3$  from theory and experiment: factor group analysis, first-principles calculations, and low-temperature infrared spectra,” *The Journal of Physical Chemistry C* **119**, 25703–25718 (2015).
- <sup>52</sup>S. Baroni and R. Resta, “Ab initio calculation of the macroscopic dielectric constant in silicon,” *Physical Review B* **33**, 7017 (1986).
- <sup>53</sup>B. Hammer, L. B. Hansen, and J. K. Nørskov, “Improved adsorption energetics within density-functional theory using revised Perdew-Burke-Ernzerhof functionals,” *Physical Review B* **59**, 7413 (1999).
- <sup>54</sup>S. Grimme, J. Antony, S. Ehrlich, and H. Krieg, “A consistent and accurate ab initio parametrization of density functional dispersion correction (dft-d) for the 94 elements h-pu,” *The Journal of chemical physics* **132**, 154104 (2010).
- <sup>55</sup>G. Makov and M. C. Payne, “Periodic boundary conditions in ab initio calculations,” *Physical Review B* **51**, 4014 (1995).
- <sup>56</sup>J. Neugebauer and M. Scheffler, “Adsorbate-substrate and adsorbate-adsorbate interactions of na and k adlayers on al (111),” *Physical Review B* **46**, 16067 (1992).
- <sup>57</sup>C. Cortes, L. D. Jackel, S. Solla, V. Vapnik, and J. Denker, “Learning curves: Asymptotic values and rate of convergence,” *Advances in neural information processing systems* **6** (1993).
- <sup>58</sup>M. Falk and T. Ford, “Infrared spectrum and structure of liquid water,” *Canadian Journal of Chemistry* **44**, 1699–1707 (1966).
- <sup>59</sup>D. A. Draegert, N. Stone, B. Curnutte, and D. Williams, “Far-infrared spectrum of liquid water,” *JOSA* **56**, 64–69 (1966).
- <sup>60</sup>G. Walrafen, “Water: a comprehensive treatise,” by F. Franks, Plenum Press, New York **1**, 151 (1972).
- <sup>61</sup>J. Hasted, S. Husain, F. Frescura, and J. Birch, “Far-infrared absorption in liquid water,” *Chemical physics letters* **118**, 622–625 (1985).
- <sup>62</sup>J. E. Bertie and Z. Lan, “Infrared intensities of liquids xx: The intensity of the oh stretching band of liquid water revisited, and the best current values of the optical constants of  $\text{H}_2\text{O}(\text{l})$  at 25 °c between 15,000 and 1  $\text{cm}^{-1}$ ,” *Applied Spectroscopy* **50**, 1047–1057 (1996).
- <sup>63</sup>D. A. Schmidt and K. Miki, “Structural correlations in liquid water: A new interpretation of ir spectroscopy,” *The journal of physical chemistry A* **111**, 10119–10122 (2007).
- <sup>64</sup>P. Madden and R. Impey, “On the infrared and raman spectra of water in the region 5–250  $\text{cm}^{-1}$ ,” *Chemical physics letters* **123**, 502–506 (1986).
- <sup>65</sup>R. Bansil, T. Berger, K. Toukan, M. Ricci, and S. Chen, “A molecular dynamics study of the oh stretching vibrational spectrum of liquid water,” *Chemical physics letters* **132**, 165–172 (1986).
- <sup>66</sup>B. Guillot, “A molecular dynamics study of the far infrared spectrum of liquid water,” *The Journal of chemical physics* **95**, 1543–1551 (1991).
- <sup>67</sup>G. Corongiu, “Molecular dynamics simulation for liquid water using a polarizable and flexible potential,” *International journal of quantum chemistry* **42**, 1209–1235 (1992).
- <sup>68</sup>P. L. Silvestrelli, M. Bernasconi, and M. Parrinello, “Ab initio infrared spectrum of liquid water,” *Chemical Physics Letters* **277**, 478–482 (1997).
- <sup>69</sup>B. Auer and J. Skinner, “Ir and raman spectra of liquid water: Theory and interpretation,” *The Journal of Chemical Physics* **128** (2008).
- <sup>70</sup>J. Xu, M. Chen, C. Zhang, and X. Wu, “First-principles study of the infrared spectrum in liquid water from a systematically improved description of h-bond network,” *Physical Review B* **99**, 205123 (2019).
- <sup>71</sup>G. R. Medders and F. Paesani, “Infrared and raman spectroscopy of liquid water through “first-principles” many-body molecular dynamics,” *Journal of Chemical Theory and Computation* **11**, 1145–1154 (2015).
- <sup>72</sup>M. J. Gillan, D. Alfe, and A. Michaelides, “Perspective: How good is dft for water?” *The Journal of chemical physics* **144** (2016).
- <sup>73</sup>M. Tanaka, G. Girard, R. Davis, A. Peuto, and N. Bignell, “Recommended table for the density of water between 0 °c and 40 °c based on recent experimental reports,” *Metrologia* **38**, 301 (2001).
- <sup>74</sup>P. M. de Hijes, C. Dellago, R. Jinnouchi, B. Schmiedmayer, and G. Kresse, “Comparing machine learning potentials for water: Kernel-based regression and behler-parrinello neural networks,” *arXiv preprint arXiv:2312.15213* (2023).
- <sup>75</sup>J. Sun, A. Ruzsinszky, and J. P. Perdew, “Strongly constrained and appropriately normed semilocal density functional,” *Physical review letters* **115**, 036402 (2015).
- <sup>76</sup>D. Weber, “ $\text{CH}_3\text{NH}_3\text{PbX}_3$ , ein pb (ii)-system mit kubischer perowskitstruktur/ $\text{CH}_3\text{NH}_3\text{PbX}_3$ , a pb (ii)-system with cubic perovskite structure,” *Zeitschrift für Naturforschung B* **33**, 1443–1445 (1978).
- <sup>77</sup>N. Onoda-Yamamuro, T. Matsuo, and H. Suga, “Calorimetric and ir spectroscopic studies of phase transitions in methylammonium trihalogenoplumbates (ii),” *Journal of Physics and Chemistry of Solids* **51**, 1383–1395 (1990).
- <sup>78</sup>Y. Kawamura, H. Mashiyama, and K. Hasebe, “Structural study on cubic-tetragonal transition of  $\text{CH}_3\text{NH}_3\text{PbI}_3$ ,” *Journal of the Physical Society of Japan* **71**, 1694–1697 (2002).
- <sup>79</sup>C. C. Stoumpos, C. D. Malliakas, and M. G. Kanatzidis, “Semiconducting tin and lead iodide perovskites with organic cations: phase transitions, high mobilities, and near-infrared photoluminescent properties,” *Inorganic chemistry* **52**, 9019–9038 (2013).
- <sup>80</sup>T. Baikie, Y. Fang, J. M. Kadro, M. Schreyer, F. Wei, S. G. Mhaisalkar, M. Graetzel, and T. J. White, “Synthesis and crystal chemistry of the hybrid perovskite ( $\text{CH}_3\text{NH}_3$ ) $\text{PbI}_3$  for solid-state sensitised solar cell applications,” *Journal of Materials Chemistry A* **1**, 5628–5641 (2013).
- <sup>81</sup>M. Bokdam, T. Sander, A. Stroppa, S. Picozzi, D. Sarma, C. Franchini, and G. Kresse, “Role of polar phonons in the photo excited state of metal halide perovskites,” *Scientific reports* **6**, 28618 (2016).
- <sup>82</sup>G. Schuck, D. M. Többsen, M. Koch-Müller, I. Efthimiopoulos, and S. Schorr, “Infrared spectroscopic study of vibrational modes across the orthorhombic–tetragonal phase transition in methylammonium lead halide single crystals,” *The Journal of Physical Chemistry C* **122**, 5227–5237 (2018).
- <sup>83</sup>M. Parrinello and A. Rahman, “Crystal structure and pair potentials: A molecular-dynamics study,” *Physical review letters* **45**, 1196 (1980).
- <sup>84</sup>M. Parrinello and A. Rahman, “Polymorphic transitions in single crystals: A new molecular dynamics method,” *Journal of Applied physics* **52**, 7182–7190 (1981).
- <sup>85</sup>M. Bokdam, J. Lahnsteiner, B. Ramberger, T. Schäfer, and G. Kresse, “Assessing density functionals using many body theory for hybrid perovskites,” *Physical review letters* **119**, 145501 (2017).