# Learning Atomistic Force Fields On-the-Fly with Bayesian Inference

Jonathan Vandermause,[1,2] Steven B. Torrisi,[1] Simon Batzner,[3] and Boris Kozinsky[2]

[1]*Department of Physics, Harvard University, Cambridge, MA 02138, USA*
[2]*John A. Paulson School of Engineering and Applied Sciences,*
*Harvard University, Cambridge, MA 02138, USA*
[3]*Center for Computational Engineering, Massachusetts Institute of Technology, Cambridge, MA 02139, USA*
(Dated: March 25, 2019)

Machine learning provides a path toward fast, accurate, and large-scale materials simulation, promising to combine the accuracy of *ab initio* methods with the computational efficiency of classical interatomic potentials. However, training current state-of-the-art models often requires hundreds of CPU hours and databases containing thousands of structures. We present an on-the-fly Bayesian inference scheme for automating and accelerating the construction of interatomic force fields. Gaussian process regression is coupled to a first principles DFT code to learn two- and three-body force fields on-the-fly with minimal trianing data. The resulting force field is easily extended to structures outside the training set and compares favorably to state-of-the-art classical and machine learned potentials.

*Ab initio* molecular dynamics is a powerful tool for accurately probing the dynamics of molecules and solids, but it is fundamentally limited by the cubic scaling of the most commonly used density functional theory (DFT) codes [1]. A common solution to this problem involves bypassing a quantum mechanical treatment of the electrons and instead directly modelling the Born-Oppenheimer potential energy surface of the ions. This is the approach taken when constructing classical interatomic potentials, which trade the accuracy of DFT and other first principles approaches for the speed and scalability of a local and analytic model, making possible the fully atomistic simulation of many thousands of atoms over nanosecond timescales. Classical potentials, however, have limited accuracy, flexibility, and transferability, and are inadequate in many settings.

A number of machine learning (ML) approaches for fitting interatomic potentials have recently been shown to approach the accuracy of first principles methods. However, most of these approaches provide only point estimates of the quantities of interest (typically energies, forces, and stesses) rather than a predictive distribution reflecting model uncertainty. Without knowledge of the highest uncertainty training points, a laborious fitting procedure is required, in which thousands of reference structures are selected *ad hoc* from a database of first principles calculations. At test time, lack of predictive uncertainty makes it difficult to determine when the fitted model is out-of-sample, leading to unreliable results and making the model difficult to update in the presence of new data.

Here, we show that on-the-fly Bayesian inference can be used to both accelerate the training of a high-quality machine learned force field and flexibly adapt the model to out-of-sample atomic environments. By coupling Gaussian process regression and density functional theory in a single molecular dynamics trajectory, it is shown that the number of DFT runs needed to train a high quality potential can be dramatically reduced from several thousand to a few dozen. By dramatically reducing the computational cost of both training and updating a high quality potential, our technique promises to extend fast and accurate ML modelling to a much wider class of systems than would otherwise be feasible.

## SUPPLEMENTARY INFORMATION

### Covariant kernels for direct force prediction

The total energy $E$ of a system of atoms in a periodic cell is modelled as a sum over two- and three-body contributions,

$$E = \sum_{ij} \varepsilon_{ij} + \sum_{ijk} \varepsilon_{ijk}, \tag{1}$$

where the sums range over all unique pairs and triplets of atoms containing at least one atom from the unfolded primary cell. In practice, the sums are truncated by considering local atom-centered environments surrounding each atom in the primary cell and neglecting contributions from atoms beyond a chosen cutoff distance from the central atom. The energy may then be expressed as

$$E = \sum_i \left( \frac{1}{2} \sum_{j \in \rho_i} \varepsilon_{ij} + \frac{1}{3} \sum_{j,k \in \rho_i} \varepsilon_{ijk} \right), \tag{2}$$

where $\rho_i$ denotes the local environment of atom $i$ containing all atoms within the cutoff sphere and the fractional factors take care of multiple counting due to the repeated appearance of bonds and triplets in neighboring environments. This may be written more compactly as

$$E = \sum_i \varepsilon_i, \tag{3}$$

where $\varepsilon_i \equiv \frac{1}{2}\sum_{j\in\rho_i}\varepsilon_{ij} + \frac{1}{3}\sum_{j,k\in\rho_i}\varepsilon_{ijk}$ may be viewed as the local energy of atom $i$.

In Gaussian process models, the covariance between targets is set equal to a kernel or similarity measure between inputs. The covariance between total energy observations $E_l, E_m$ of two distinct structures $\sigma_l, \sigma_m$ may be written as

$$\langle E_l E_m \rangle = \sum_{i\in\sigma_l}\sum_{j\in\sigma_m} \langle \varepsilon_i \varepsilon_j \rangle, \tag{4}$$

where the covariance between local energies is

$$\langle \varepsilon_i \varepsilon_j \rangle = \frac{1}{4}\sum_{n\in\rho_i}\sum_{p\in\rho_j}\langle \varepsilon_{in}\varepsilon_{jp}\rangle + \frac{1}{9}\sum_{n,q\in\rho_i}\sum_{p,r\in\rho_j}\langle \varepsilon_{inq}\varepsilon_{jpr}\rangle. \tag{5}$$

Letting $F_{i\xi} = -\frac{dE}{d\xi_i}$ denote the force on atom $i$ along Cartesian component $\xi$, the covariance between force observations may be written as

$$\langle F_{i\xi} F_{j\chi}\rangle = \sum_{n\in\rho_i}\sum_{p\in\rho_j}\frac{\partial^2}{\partial\xi_i\partial\chi_j}\langle \varepsilon_{in}\varepsilon_{jp}\rangle + \sum_{n,q\in\rho_i}\sum_{p,r\in\rho_j}\frac{\partial^2}{\partial\xi_i\partial\chi_j}\langle \varepsilon_{inq}\varepsilon_{jpr}\rangle, \tag{6}$$

where here the fractional factors do not appear as the sums are restricted to the local environments of atoms $i$ and $j$. This is convenient in practice, as it allows local information about individual atoms to be used at test time without having to store the atomic environments of all the other atoms in the structure.

In order to infer energies from force observations, it is necessary to consider the covariance between energies and forces.
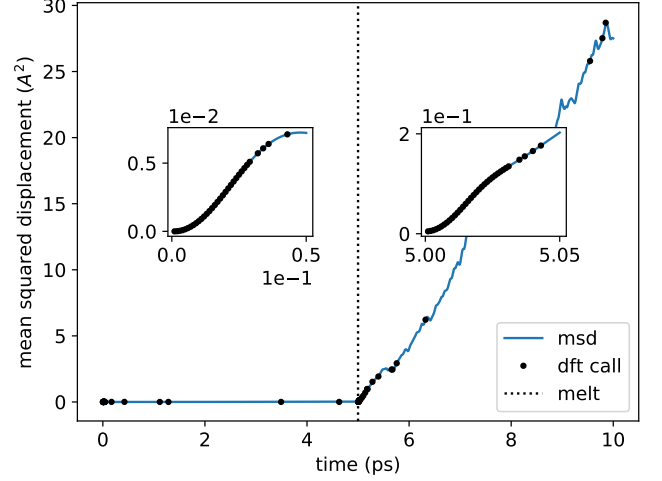


FIG. 1. On-the-fly learning of an aluminum force field at multiple temperatures.
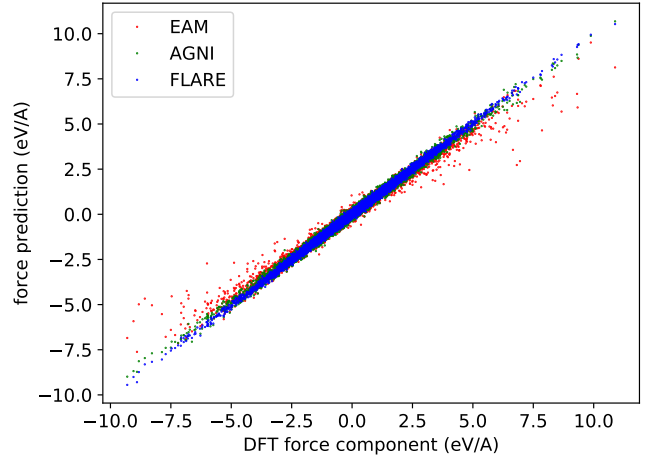


FIG. 2. Mean square displacement of aluminum melt.

| | | Solid | Liquid | Slab | Vacancy |
|---|---|---|---|---|---|
| OTF | | | | | |
| EAM | | | | | |
| AGNI | | | | | |

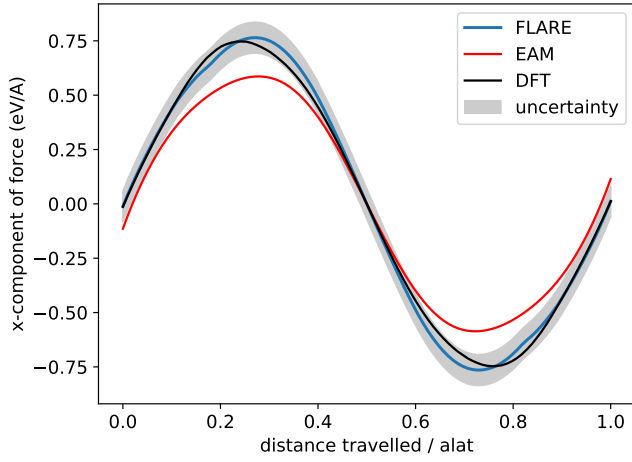TABLE I. On-the-fly force field error compared to a recent EAM potential.
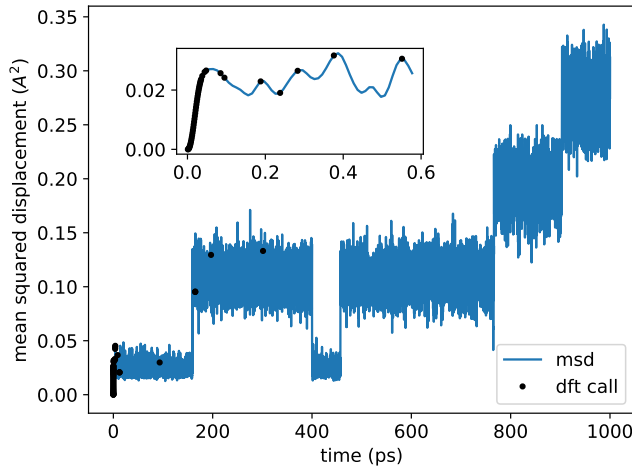
FIG. 3. Mean square displacement of aluminum melt.
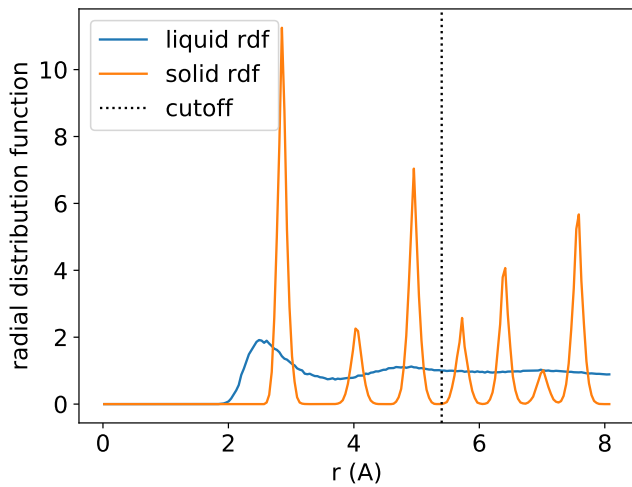


FIG. 4. Mean square displacement of aluminum melt.



FIG. 5. RDF of Al melt.

---

**Algorithm 1** Active Learning of Atomistic Force Fields

---

**Require:** initial structure (positions, velocities, periodic cell)
**Require:** initial GP model (kernel and hyperparameters)
**Require:** $\Delta t$: molecular dynamics time step
**Require:** $T$: total simulation time
**Require:** $\mathcal{U}$: initial uncertainty threshold

1: Initialize time: t = 0
2: **while** $t < T$ **do**
3:     predict forces and uncertainties with GP model
4:     **if** uncertainty above threshold **then**
5:         compute forces with DFT
6:         add highest uncertainty atom to training set
7:         update GP hyperparameters
8:         update structure with DFT forces
9:     **else**
10:         update structure with GP forces
11:     **end if**
12:     update time: $t = t + \Delta t$
13: **end while**

---

[1] W. Kohn, Reviews of Modern Physics **71**, 1253 (1999).

| Energy Kernel | $k_{\mathrm{inv}}$ | $\sigma^2 \sum_{c,p} k f_{\mathrm{cut}}(\vec{d_c}) f_{\mathrm{cut}}(\vec{d_p})$ |
|---|---|---|
| - | $k$ | $\exp\left(-\frac{\|\vec{d_c}-\vec{d_p}\|^2}{2\ell^2}\right)$ |
| - | $\vec{d}^{(2)}$ | $(r_{i_1})$ |
| - | $\vec{d}^{(3)}$ | $(r_{i_1}, r_{i_2}, r_{i_1,i_2})$ |
| Force Kernel | $\frac{\partial^2 k_{\mathrm{inv}}}{\partial \xi_i \partial \chi_j}$ | $\sigma^2 \sum_{c,p}(k_0 + k_1 + k_2 + k_3)$ |
| - | $k_0$ | $k \frac{\partial f_{\mathrm{cut}}(\vec{d_c})}{\partial \xi_i} \frac{\partial f_{\mathrm{cut}}(\vec{d_p})}{\partial \chi_j}$ |
| - | $k_1$ | $\frac{\partial k}{\partial \xi_i} f_{\mathrm{cut}}(\vec{d_c}) \frac{\partial f_{\mathrm{cut}}(\vec{d_p})}{\partial \chi_j}$ |