# Center for the Exascale Simulation of Material Interfaces in Extreme Environments

## *Massachusetts Institute of Technology*

Cuong Nguyen: Performance and Accuracy of Machine Learning Potentials

Dionysios Sema: E(3)-Equivariant ML-DFT

Spencer Wyant: Developing Machine-Learning Interatomic Potentials with *julia*

# PERFORMANCE AND ACCURACY OF MACHINE LEARNING POTENTIALS

DIONYSIOS SEMA, YEONGSU CHO, NGOC NGUYEN, YOUSSEF MARZOUK, NICOLAS HADJICONSTANTINOU, HEATHER KULIK

MASSACHUSETTS INSTITUTE OF TECHNOLOGY

## INTRODUCTION

We introduce the proper orthogonal descriptors for efficient and accurate interatomic potentials of complex chemical systems [1, 2, 3]. We compose the proper orthogonal descriptors to develop two interatomic potentials by expressing the per-atom energies as a linear and then as a linear and quadratic combination of proper orthogonal descriptors. We demonstrate the weak and strong scaling of these potentials and perform MD simulations to calculate material properties. We also perform MD simulations using Allegro [4] on complex systems to model the oxidation process of Hf and map the vapor-liquid dome of Al. Finally, we use POD and Allegro for accurate prediction of melting points and mechanical properties of our target systems.
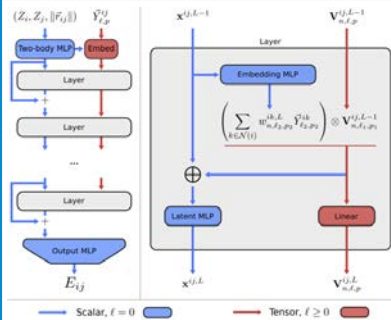
## E(3) EQUIVARIANT DEEP NNs



**Figure 1:** Allegro architecture (left) and the component of each layer (right). Features in the invariant and equivariant latent spaces interact at each layer with tensor products.

The Allegro Network contains separate latent spaces for invariant features (scalar, rotation order $l = 0$) and equivariant features of rotation order $l > 0$. The 2-body interactions are encoded with Bessel basis functions and a polynomial envelope and ACE-like many body interactions are constructed in each layer with: $N_{order} = N_{layers} + 2$. The final output are pairwise energies, $E_{ij}$

## ONGOING RESEARCH AND INTEGRATION

We plan to develop NN/GCN potentials based on the proper orthogonal descriptors. We plan to equip POD potentials with UQ and active learning methods developed by our CESMIX team. We are also working on a Kokkos implementation to gain significant performance gains using multiple GPUs. As we aim to make the POD potentials available to the LAMMPS community, we would like to collaborate the DOE labs to implement POD models in LAMMPS and FitSNAP.

## PROPER ORTHOGONAL DESCRIPTORS

We introduce the following set of snapshots on $(r_{min}, r_{max})$:

$$\xi_\ell(r_{ij}, \boldsymbol{\eta}) = V^{(2)}(r_{ij}, \boldsymbol{\eta}, \boldsymbol{\mu}_\ell), \quad \ell = 1, \dots, N_s \quad (1)$$

and compute the covariance matrix

$$C_{ij} = \frac{1}{N_s} \int_{r_{min}}^{r_{max}} \xi_i(x, \boldsymbol{\eta}) \xi_j(x, \boldsymbol{\eta}) dx, \quad 1 \le i, j \le N_s.$$

We then solve the eigenvalue problem $\boldsymbol{Ca} = \lambda \boldsymbol{a}$ to obtain the orthogonal basis functions

$$U_m^{(2)}(r_{ij}, \boldsymbol{\eta}) = \sum_{\ell=1}^{N_s} a_{\ell m}(\boldsymbol{\eta}) \, \xi_\ell(r_{ij}, \boldsymbol{\eta}),$$

$$U_{mn}^{(3)}(r_{ij}, r_{ik}, \theta_{ijk}, \boldsymbol{\eta}) = U_m^{(2)}(r_{ij}, \boldsymbol{\eta}) U_m^{(2)}(r_{ik}, \boldsymbol{\eta}) \cos(n\theta_{ijk})$$

and compute the proper orthogonal descriptors

$$D_{im}^{(2)}(\boldsymbol{\eta}) = \sum_j U_m^{(2)}(r_{ij}, \boldsymbol{\eta}),$$

$$D_{imn}^{(3)}(\boldsymbol{\eta}) = \sum_j \sum_k U_{mn}^{(3)}(r_{ij}, r_{ik}, \theta_{ijk}, \boldsymbol{\eta}). \quad (2)$$

The linear POD potential is defined as

$$E(\boldsymbol{\eta}) = c^{(1)} + \sum_{m=1}^{N_d^{(2)}} c_m^{(2)} d_m^{(2)}(\boldsymbol{\eta}) + \sum_{n=1}^{N_d^{(3)}} c_n^{(3)} d_n^{(3)}(\boldsymbol{\eta}) \quad (3)$$

with $d_m^{(q)} = \sum_{i=1}^N D_{im}^{(q)}$. The quadratic POD potential is

$$E = c^{(1)} + \sum_{m=1}^{N_d^{(2)}} \left( c_m^{(2)} + b_m^{(2)} \right) d_m^{(2)} + \sum_{n=1}^{N_d^{(3)}} \left( c_n^{(3)} + b_n^{(3)} \right) d_n^{(3)} \quad (4)$$

with $b_m^{(2)} = \sum_{n=1}^{N_d^{(3)}} c_{mn}^{(23)} d_n^{(3)}$ and $b_n^{(3)} = \sum_{m=1}^{N_d^{(2)}} c_{mn}^{(23)} d_m^{(2)}$. The coefficients $c^{(1)}, c_m^{(2)}, c_n^{(3)}, c_{mn}^{(23)}$, and $\boldsymbol{\eta}$ are found by solving a nonlinear least-squares regression against DFT data. Extension to multi-element systems is carried out by computing the PODs in (2) for different atom types. The complexity of the resulting POD potentials is $O(NN_n^2 N_f)$, where $N_n$ is the number of neighbors and $N_f$ is the number of basis functions. The complexity of the multi-element SNAP potential is $O(NN_n N_f^2 N_e^2)$. Hence, the cost ratio between multi-element SNAP and POD potentials is $N_f N_e^2/N_n$.

## LITHIUM ION DIFFUSIVITY

| Lin. SNAP1 | NN SNAP2 | NN SNAP1 | Allegro | POD | POD+SNAP |
|---|---|---|---|---|---|
| 180.4 | 116.2 | 84.5 | 70.4 | 63.3 | 48.1 |

Table 1: Force training MAE (meV/Å) of various potentials for LGPS. SNAP1 refers to $j_{max} = 3$ (31 descriptors) and SNAP2 refers to $j_{max} = 4$ (56 descriptors). For POD+SNAP, we use $j_{max} = 2$ which includes 15 4-body SNAP descriptors on top of the 91 2-body and 3-body POD descriptors.
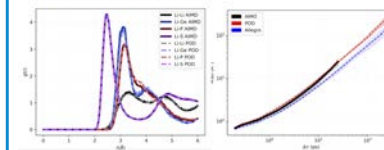


**Figure 2:** Simulation of Lithium ion diffusivity using AIMD, POD, and Allegro [5]. Our results show the ability of linear POD to accurately model diffusion in Li superionic conductors.
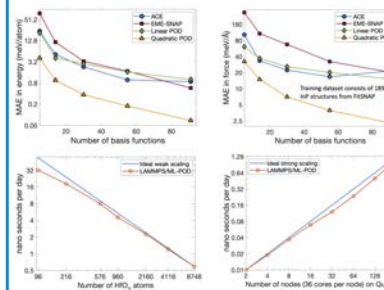
## PERFORMANCE SCALING



**Figure 4:** Accuracy comparison of different MLIP models for InP [3] (top) and performance scaling of POD using LAMMPS/ML-POD. For weak scaling (left), we performed MD on various HfO$_2$ sized systems on 36 cores. For strong scaling (right), we performed MD simulations of 1M HfO$_2$ atoms.
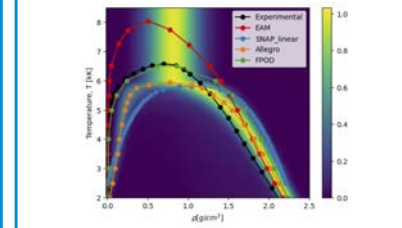
## VAPOR-LIQUID EQUILIBRIA



**Figure 3:** Prediction of vapor-liquid Al phases using SNAP, POD, Allegro, and EAM. The critical point of the experiment is $(\rho_c, T_c) = (0.745 g/cm^3, 6500K)$. POD matches the experimental data better than Allegro, SNAP and the empirical EAM potentials.
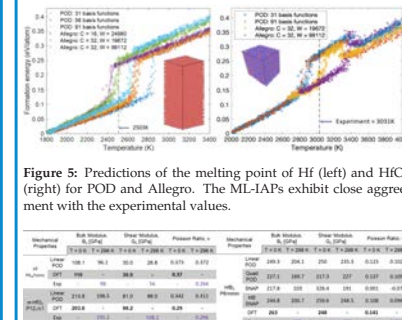
## MELTING AND MECHANICAL QOIs



**Figure 5:** Predictions of the melting point of Hf (left) and HfO$_2$ (right) for POD and Allegro. The ML-IAPs exhibit close agreement with the experimental values.



**Figure 6:** Predictions of mechanical properties for Hf, HfO$_2$ (left) and HfB$_2$ (right) for POD at 0K and room temperature in comparison with DFT and experiments.

## REFERENCES

[1] N. C. Nguyen and A. Rohskopf. *Journal of Computational Physics*, 480:112030, 2023.

[2] N. C. Nguyen. *Physical Review B*, 107(14):144103, apr 2023.

[3] N. C. Nguyen. *Journal of Computational Physics*, 2023. In submission.

[4] A. Musaelian, S. Batzner, A. Johansson, L. Sun, C. J. Owen, M. Kornbluth, and B. Kozinsky. *Nature Communications*, 14(1):579, 2023.

[5] A. Rohskopf et al. *Journal of Materials Research*, 2023. In submission.

# E(3)-EQUIVARIANT ML-DFT

## DIONYSIOS SEMA

### MASSACHUSETTS INSTITUTE OF TECHNOLOGY

**Sandia National Laboratories**

**Massachusetts Institute of Technology**

**MITCESMIX**

## INTRODUCTION

The development of an accurate and transferable machine-learned interatomic potentials (ML-IAPs) can take several months and requires computing resources and expertise across disciplines. In active learning approaches, we can use uncertainty metrics from an ensemble of trained potentials to drive an active learning approach, expand our datasets with additional DFT calculations and retrain the next generation of ML-IAPs. In this process, data generation is the most time-consuming part. Our goal is to accelerate this process and achieve potential-in-a-day development cycles. To this end, we propose a E(3)-equivariant ML-DFT that can be used as a proxy DFT model to accelerate DFT calculations or perform DFT calculations that predict the electron density, energy and atomic forces of large and complex systems ($\mathcal{O}(10^5)$ atoms) that are intractable for regular DFT codes.

## SELECTION SCHEME

We use a minimax sampling method to select $M$ configurations from a set of $K$ configurations $\{\mathcal{C}_k\}_{k=1}^K$ with $M \ll K$. Define the similarity matrix

$$S_{ij} = \frac{D(\mathcal{C}_i) \cdot D(\mathcal{C}_j)}{\|D(\mathcal{C}_i)\|\|D(\mathcal{C}_j)\|}, \quad i,j = 1,\dots,K, \quad (1)$$

where $D(\mathcal{C}_i)$ is a vector of descriptors for $\mathcal{C}_i$. Choose a set of $M$ indices $\{i_1,\dots,i_M\}$, where the first two indices are

$$(i_1, i_2) = \arg \min_{1 \le i,j,\le N} S_{ij} \quad (2)$$

and, for $m = 3,\dots,M$,

$$i_m = \arg \min_{1 \le i \le N} \max_{j \in \{i_1,\dots,i_{m-1}\}} S_{ij}. \quad (3)$$

This method is used in the following algorithm.

## E(3)-GCN ELECTRON DENSITY

The electron density is a scalar value over all 3D space. We typically represent it using a "basis set". The functions of the basis set have the mathematical form:

$$\Phi_{l,m} = Y_l^m \exp(-\alpha_{l,m}\|r - R_i\|^2), \quad (4)$$

where the first term are the spherical harmonics, and the second is a Gaussian radial basis. The density on a given atom, $i$, is represented by a linear combination of the basis functions projected onto a delta Dirac function (the origin of DFT formulation). Each basis function has a coefficient that is the weight of that function's contribution:

$$\rho_i = \sum_\lambda \delta(r - r_\lambda)\|\psi_\lambda(r)\|^2 = \sum_l C_i^{l,m}\Phi_{l,m}. \quad (5)$$
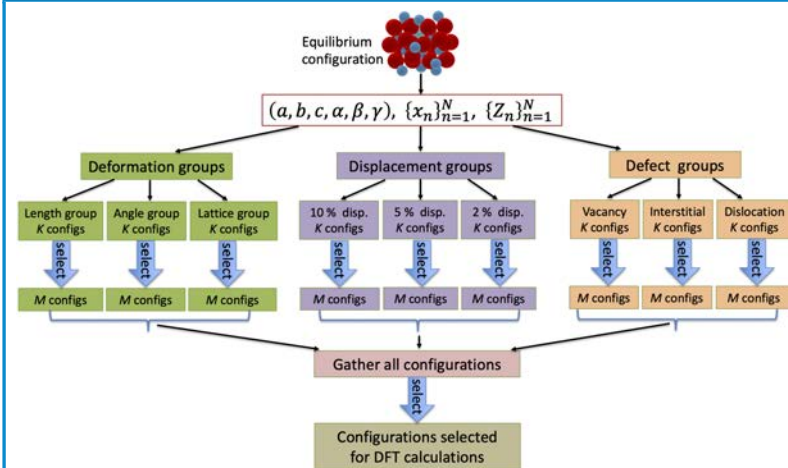
We perform DFT runs with pyscf with pbc and project the electron density onto a density fitting basis. We calculate the 3-center 2e tensor integral, $P|ij$, and the 2c2e integral, $P|Q$:

```
mol = gto.M(atom='C 0 0 0', a=np.eye(3)*5, basis='aug-cc-pvtz')
kpts = mol.make_kpts([1,1,1])
mol.spin = 2
mf = dft.KUKS(mol, kpts, xc="pbe").density_fit()
mf.kernel()
auxbasis = 'def2-universal-jfit'
auxmol = df.addons.make_auxmol(mol, auxbasis)
ints_3c2e = df.incore.aux_e1(mol, auxmol, intor='int3c2e').transpose(1, 2, 0)
ints_2c2e = auxmol.intor('int2c2e')
evals, evecs = np.linalg.eigh(ints_2c2e)
evals = np.where(evals < 1e-10, 0.0, 1.0/evals)
ints_2c2e_inv = np.einsum('ik,k,jk->ij', evecs, evals, evecs)
```

We can then extract the final irreps that are determined for each atom type, following the approach of Dunlap et al. [1]. The coefficients and exponents of the basis functions is the data we will train the model with. The raw data is subtracted from the density of the isolated atoms and converted to a molecular graph.

The GCN contains 3 layers of fully connected tensor products with gated block non-linearities. The input are the coordinates and atoms types concatenated to the radial basis as a one-hot vector of length, $N$, with irreps $Nxoe$. The hidden features have 16 copies with $l_{max} = 4$ with even and odd parity, $p = \{-1, 1\}$. The cutoff radius was set to $r_{cut} = 4 \text{Å}$.

## AUTOMATIC CONFIGURATION GENERATION

# Developing Machine-Learning Interatomic Potentials with Julia

Spencer Wyant, Youssef Marzouk

*(with package contributions from Emmanuel Lujan, Dallas Foster, Joanna Zou, and other CESMIX contributors)*
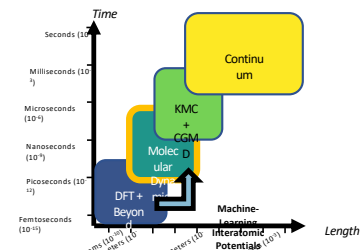
MIT CESMIX

## Background

When performing **molecular dynamics (MD)**, the motion of a given particle $i$ in a system of $N$ atoms is governed by Newton's equation of motion:

$$-\nabla_i E(\boldsymbol{r}, \boldsymbol{z}) = \boldsymbol{F}_i(\boldsymbol{r}, \boldsymbol{z}) = m_i \frac{d^2 \boldsymbol{r}_i}{dt^2}$$

### Potential Energy Surface (PES)

$$E(\boldsymbol{r}, z): \mathcal{R} \times \mathcal{Z} \to \mathbb{R} \qquad \mathcal{R} := \bigcup_{N=1}^{\infty} \mathcal{R}_N, \ \mathcal{R}_N := \mathbb{R}^{3N}$$

$$\mathcal{Z} := \{1, \dots, N_e\}$$

The **PES** maps the positions and chemical identities of each atom in the system to the total potential energy (in the absence of external forces). A model of the PES should be capable of treating systems with arbitrary numbers of atoms.
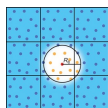


While **density functional theory (DFT)** calculations can provide an accurate PES and downstream properties, their computational cost and O($N^3$) scaling prevents accessing systems larger than a few hundred atoms and timescales beyond tens of picoseconds. **Machine-learning interatomic potentials (MLIPs)**, which leverage flexible models trained on DFT-computed forces and energies, act as a bridge to the larger spatial and temporal scales accessed with traditional MD methods, bringing the *ab initio* accuracy of more expensive methods to these scales.
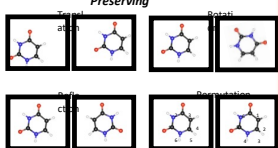
### Key Features of MLIPs

**Additivity and Locality**

$$E(\boldsymbol{r}) = \sum_{i=1}^{N} E_i\left(\{\boldsymbol{r}_j\}_{j \in N(i)}\right)$$

**Symmetry Preserving**

MLIPs can model systems with arbitrary number of atoms by decomposing the total energy into a sum of atom-wise or pairwise terms. Likewise, they exploit the typically local nature of interparticle forces by employing a cutoff function, which only considers atoms within some neighborhood $\mathcal{N}$. Finally, interatomic potentials need to respect the Euclidean group E(3) symmetries (energy invariance and force equivariance) and, from a computational perspective, should be invariant to the ordering of atoms.

## Goals

### Problem Input

**Material System(s)** (Examples):
- AlN
- HfO2
- LiEuGe

**Types of Atomic Environments** (Target QoIs) (Examples):
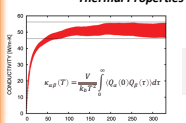- mono, b, s, amorphous, vacancy

**Constraints**
- **Minimum Inference Speed** — Need fast force routines, especially if simulating long time scales and/or many atoms
- **Data Budget** — Since DFT calculations are expensive, want to minimize # of fitting data needed (even more so if using more expensive methods like CCSD); *Available compute resources*; *Soft constraints depend on CCSD*
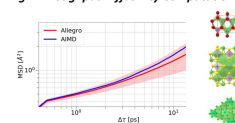
MLIPs are typically developed to model particular material system(s) and specific quantities of interest (QoIs), which in turn require modeling different kinds of atomistic environments. Some QoIs require diverse atomistic environments, effectively sampling a larger "volume" of configuration space. Additionally, it may not be known which environment are relevant *a priori*, necessitating a generalized MLIP. Thus, the choice of material systems and the scope of atomic environments are the primary inputs to the MLIP development process, subject to the constraints listed on the right.
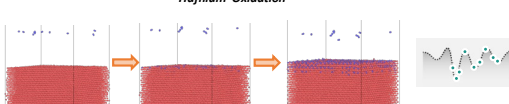
### Example Problems

**Thermal Properties**



$$\kappa_{\alpha\beta}(T) = \frac{V}{k_B T^2} \int_0^{\infty} \langle Q_\alpha(0) Q_\beta(\tau) \rangle d\tau$$

- Only need to sample a small portion of configuration space (i.e., a single minima of the PES).
- However, need high accuracy to capture curvature of the minima, which controls thermal properties

**High-Throughput Diffusivity Computation**

- Allegro
- AIMD

- Lower accuracies are acceptable, as the goal is down-selection.
- However, need highly sample-efficient strategies (i.e., small number of fitting data for each material system)
- One possible strategy could be to develop a universal MLIP that is subsequently fine-tuned on specific material systems

**Hafnium Oxidation**

- Requires sampling a wide configuration space, including hafnium metal, hafnia and related sub-oxides, surfaces and surface adsorption, and oxygen gas.
- A major research focus of the CESMIX project

### Primary Tasks

**Need to Select and/or Optimize**
1. *Model Class* (e.g., linear, neural network, gaussian process, descriptor type, etc.)
2. *Model Hyperparameters* (# of basis vectors, NN architecture, etc.)
3. *Model Parameters* (standard optimization)
4. *Fitting Data Generation & Selection* (Experimental design, active learning)
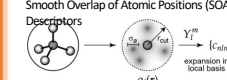
**Two Perspectives**

*Scientist Perspective:*
- Pre-implemented workflows to generate well-validated MLIPs that satisfy inputs and constraints
- Workflows are automated/semi-automated

*Methods Developer Perspective:*
- Facilitate the development of new MLIP models and/or data selection and fitting strategies
- Provide an ability to comparatively evaluate different MLIPs and different
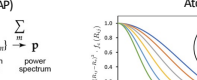
## Packages Overview

### InteratomicPotentials.jl

**Atomistic Descriptors** (examples below)

Smooth Overlap of Atomic Positions (SOAP) Descriptors
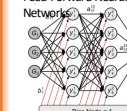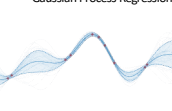


Atom-centered Symmetry Functions

**Regression Models** (examples below)

Feed-Forward Neural Networks

Gaussian Process Regression

Linear Models
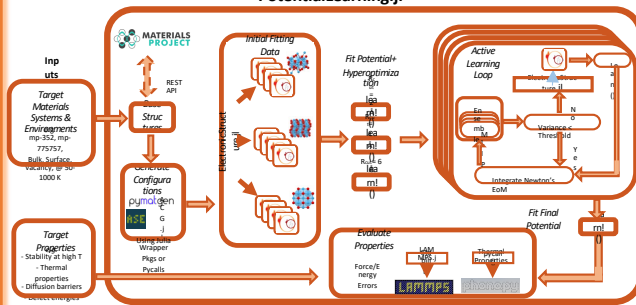
$$E_i(\boldsymbol{r}) = \sum_k \theta_k \phi_k(\boldsymbol{r})$$

InteratomicPotentials.jl is intended to provide a library of existing MLIPs and classical potentials implemented in Julia, along with abstractions that makes it easy to develop and compose new MLIPs. In part, this involves providing implementations of different atomistic descriptors (i.e., ways of representing local atomistic environments) that can be fed into different kinds of regression models, with the flexibility for users to create their own descriptors and/or regressors. More recent MLIP methods that *learn* descriptors (e.g., graph neural network methods) are an important development target, though they likely require a concurrent expansion/improvement of Julia's machine learning ecosystem. InteratomicPotentials.jl can integrate with Molly.jl, a Julia-native MD code, as well as with LAMMPS using the LAMMPS.jl wrapper package, also developed as part of CESMIX.

### PotentialLearning.jl



The PotentialLearning.jl package facilitates the development of MLIPs, including data generation and selection, potential fitting, and hyperparameter optimization. It acts as a high-level orchestration script, providing the necessary abstractions to key actions like model fitting and data selection that can be stitched together into larger workflows like that presented above. As a development goal, this package should provide sufficient flexibility to recreate most of the learning strategies that exist in the MLIP literature (including a number of active learning variants), while making it easy for users to develop and test new strategies. One important challenge is to ensure that workflows like the one above are HPC-compatible, i.e., can be easily run in an HPC environment in a nearly automated fashion. Two strategies are being considered: one approach leverages the Flux resource manager via FluxRM.jl; the other uses the AiiDA workflow manager, a package written in python that would require some Julia integration, but with key benefits including full provenance tracking and an error management & recovery framework.

**Figure Credits:**
- Behler, Jörg. *Chemical Reviews* 121.16 (2021): 10037-10072.
- Deringer, Volker L., et al. *Chemical Reviews* 121.16 (2021): 10073-10141.
- Musaelian, Albert, et al. *Nature Communications* 14.1 (2023): 579.
- Jones, Reese E., and Kranthi K. Mandadapu. s. *The Journal of chemical physics* 136.15 (2012): 154102.
- Behler, Jörg. *International Journal of Quantum Chemistry* 115.16 (2015): 1032-1050.
- https://docs.jaxgaussianprocesses.com/

Visit us at: github.com/orgs/cesmix-