

## 428 Appendix

### 429 A Details of Datasets

Table A.1: Overview of total simulations performed for the HfO dataset.  $N_{\text{cond}}$  represents the total number of simulation conditions,  $N_{\text{str}}$  is the total number of MD snapshots, and  $N_{\text{f}}$  indicates the total number of atomic force data. Structures with index (idx.) indicate the same crystal family but have different lattice parameters.

Element	Structure (idx.)	Simulation Condition	$N_{\text{cond}}$	$N_{\text{str}}$	$N_{\text{f}}$
Hf & O	Monoclinic	m-MQA method	6	16,000	1,536,000
	Tetragonal	m-MQA method	6	16,000	1,536,000
	Cubic	m-MQA method	6	16,000	1,536,000
	Orthorhombic (1)	m-MQA method	6	16,000	1,536,000
	Orthorhombic (2)	m-MQA method	6	16,000	1,536,000
	Randomized Structures (ID)	m-MQA method	30	80,000	7,680,000
Hf & O	Randomized Structures (OOD)	m-MQA method	12	32,000	3,072,000

Table A.2: Overview of total simulations performed for the SiN dataset.  $N_{\text{cond}}$  represents the total number of simulation conditions,  $N_{\text{str}}$  is the total number of MD snapshots, and  $N_{\text{f}}$  indicates the total number of atomic force data.

Element	Structure	Simulation Condition	$N_{\text{cond}}$	$N_{\text{str}}$	$N_{\text{f}}$
Si & N	Amorphous	High T MD (2000, 4000 K; w/o strain)	11	13,650	1,006,500
		High T MD at 1500K (w/ strain $\pm 7\%$ )	4	4,000	306,000
		Quenching (4000 to 300 K; w/o strain)	8	26,150	1,305,600
		Quenching (4000 to 300 K; w/ strain $\pm 7\%$ )	4	2,000	153,000
		Structure Relaxation (w/o strain)	11	1,274	77,952
		Structure Relaxation (w/ strain $\pm 7\%$ )	4	570	38,760
	Crystals ( $\alpha, \beta, \gamma$ )	MD (300, 1200, 2100 K)	9	18,000	504,000
		Structure Relaxation	9	569	15,932
	Defect Structures	MD with Divacancy (2000 K)	3	300	87,600
		MD with Grain-Boundary (2000 K)	5	500	103,600
	Surfaces	MD (300, 1200 K)	6	6,000	480,000
		Structure Relaxation	6	1,400	105,600
	Isolated Clusters	MD (300, 1200 K)	6	1,200	142,400
		Structure Relaxation	6	600	70,800
Si	Crystals (Diamond, SC, BCC, FCC)	MD (2100 K; w/o strain)	4	2,000	69,500
		MD (2100 K; w/ strain $\pm 7\%$ )	8	4,000	139,000
	Defect Structures	MD with Divacancy (2100 K)	1	150	32,100
		MD with Two Vacancies (2100 K)	1	100	51,000
N	N <sub>2</sub> Molecules	MD (2100, 4000 K)	4	2,000	128,000
Si & N	Amorphous (OOD)	Melt, Quench, and Relaxation	3	3,700	388,500

#### 430 A.1 Guidance

431 Simulation data can be flexibly utilized by researchers to study interatomic potential models according  
432 to their specific preferences. For example, the HfO dataset was generated using the m-MQA method,  
433 which involved initial structures encompassing various crystals as well as randomized structures. In  
434 our HfO dataset, we included both types of structures to provide researchers with a comprehensive  
435 dataset. Researchers have the flexibility to train their models using a subset of the dataset consisting  
436 of randomized structures and evaluate their performance on the remaining portion of the dataset

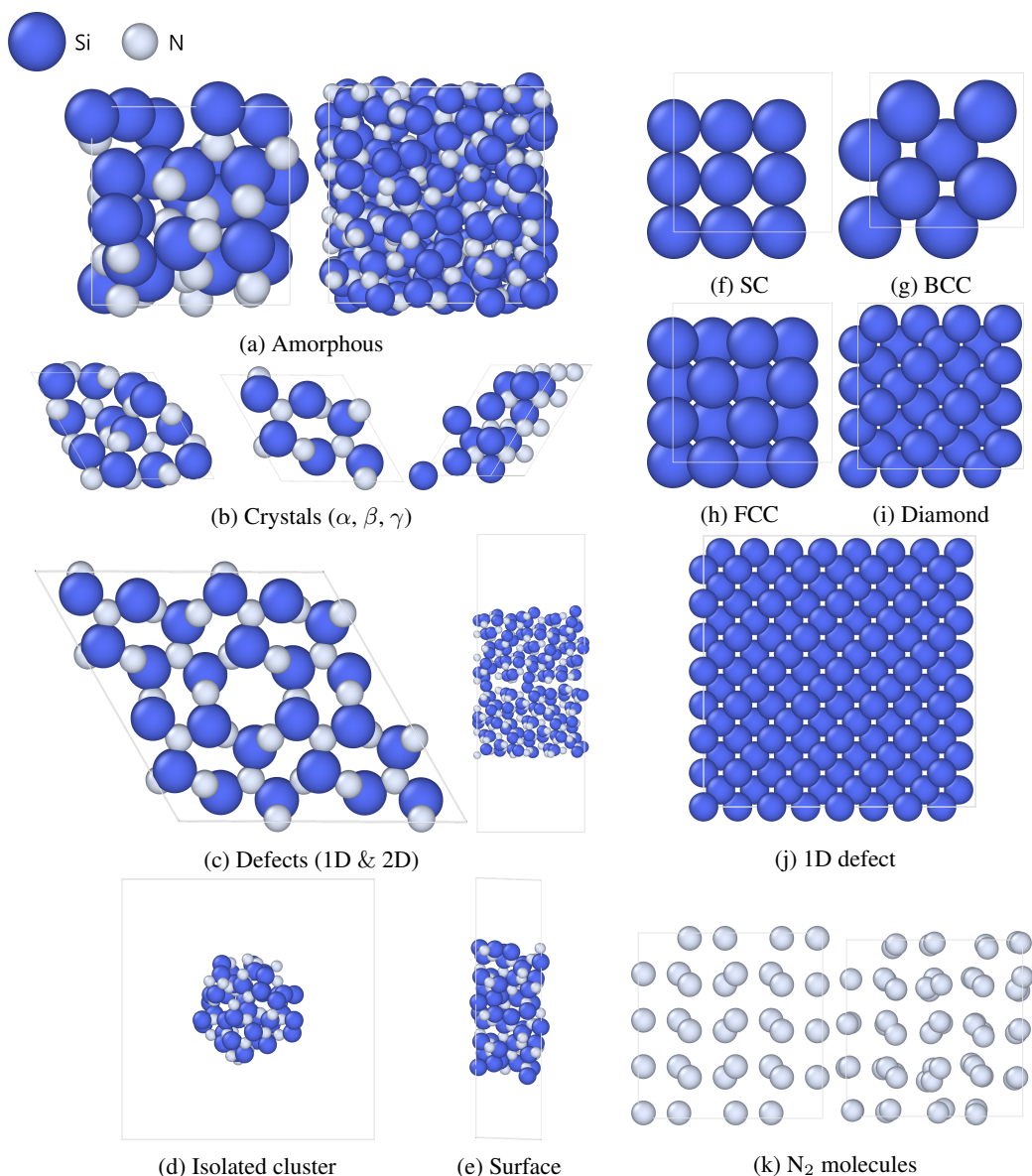


Figure A.1: Reference atomic structures employed in the SiN dataset.

comprising crystal structures, and vice versa. This approach enables thorough testing and assessment of model performance across different structural configurations within the HfO dataset.

Researchers can modify the SiN simulation data to investigate the generalization performance of MLFF models. Despite the integrated SiN dataset containing not only SiN compounds but also Si-only structures and N-only structures, researchers have the choice to exclusively use SiN compounds for model training. Subsequently, they can evaluate the generalization performance of these models on Si- and N-only structures. However, it is crucial to recognize the significant challenge arising from the substantial disparities in sample distributions between SiN compounds and the Si- and N-only structures.

The generation of specialized datasets for silicon nitride necessitates a high degree of domain expertise. However, our proposed m-MQA scheme represents a notable advancement by presenting a streamlined approach for generating the HfO dataset. In contrast to conventional methods reliant on specialized knowledge, this scheme empowers researchers familiar with DFT to effortlessly generate the dataset. By simplifying the dataset generation process and eliminating the requirement for

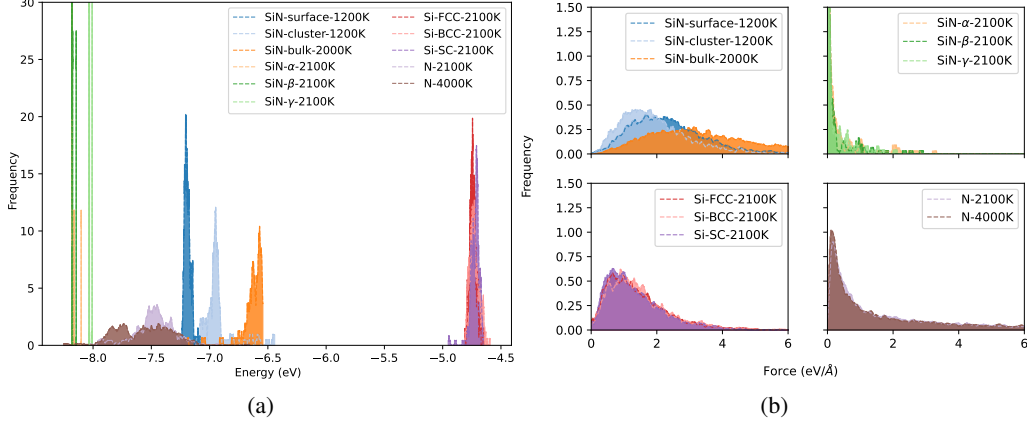


Figure A.2: Distribution of energy per atom and forces for SiN datasets: (a) energy per atom and (b) force per atom.

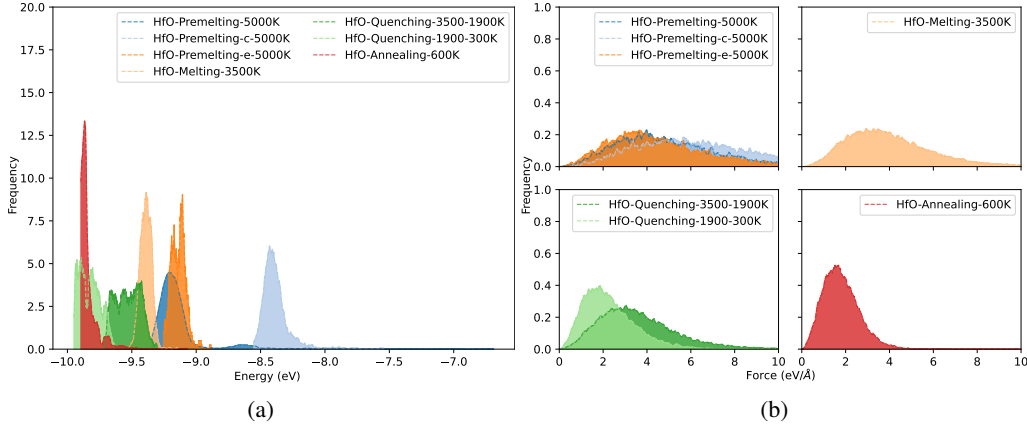


Figure A.3: Distribution of energy per atom and forces for selected MD simulations for HfO: (a) energy per atom and (b) force per atom.

extensive domain expertise, our m-MQA scheme enhances accessibility for researchers, facilitating advancements in MLFF research.

## A.2 Dataset Characteristics

To illustrate the diversity of datasets of our DFT calculations, we analyze the energy and force distribution for each DFT calculation. Figure A.2 shows the energy and force distributions of the representative crystal structures in the SiN datasets. For brevity, we only plot the energy and force distribution of single MD scenario for each structure. The SiN dataset contains numerous phases, showing well-separated variations for both energy and forces. The force distribution is particularly interesting since it reflects the structural properties of each MD simulation. The force distribution shows a considerable variation ranging from 0 to 6 eV/Å for amorphous phases. However, the highly symmetric crystalline phases ( $\alpha$ ,  $\beta$ , and  $\gamma$ -SiN) show that most of the forces on atoms are zero. MD simulations on face-centered cubic (FCC), body-centered cubic (BCC), and simple cubic (SC) also show relatively large variations of forces, but their distributions are narrower than amorphous cases. The N-alone case shows similar distribution as the SiN-crystalline case.

Unlike SiN, the HfO datasets consist of 96 atoms with a fixed Hf:O stoichiometry of 1:2. Each MD scenario shows variations in energy and force, as demonstrated in Figure A.3 (a). The high-temperature pre-melting stage at 5000K shows the largest energy distribution. Depending on the compressive (HfO-premelting-c-5000K) or expansive (HfO-premelting-e-5000K) strain, they exhibit

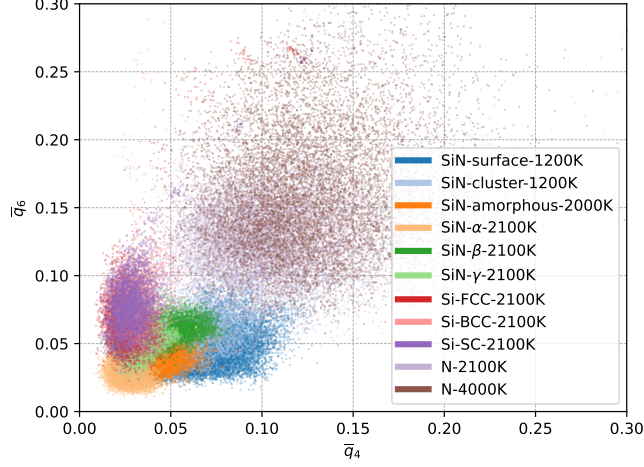


Figure A.4: Distribution of the averaged bond order parameter  $\bar{q}_4$  and  $\bar{q}_6$  for SiN.

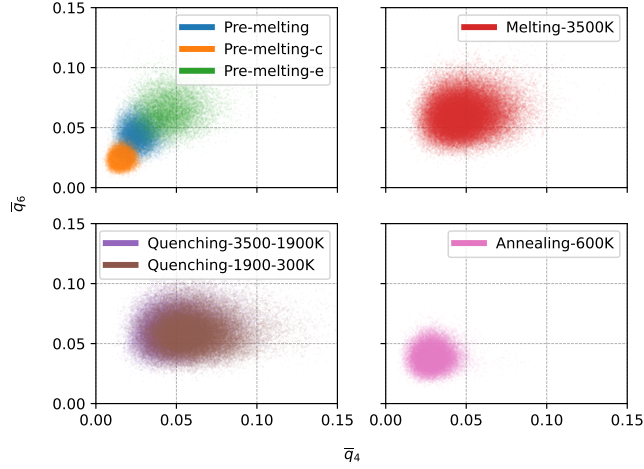


Figure A.5: Distribution of the averaged bond order parameter  $\bar{q}_4$  and  $\bar{q}_6$  for HfO.

well-separated energy distribution. The average energy also lowers as the temperature decreases through the melting and quenching process. For the annealing process, they show the narrowest energy distribution. The force distributions of these MD scenarios also follow the typical patterns shown in MD simulations. For liquid-like pre-melting and melting stages (Figure A.3 (b)), we see relatively broader force distributions. As the temperature decreases through the quenching and annealing process, the average force on each atom decreases.

Another essential aspect of the MLFF dataset is the inclusion of diverse atomic environments. Since the local atomic environment determines the atomic energy and forces, having a diverse atomic environment in the dataset is crucial. One way of quantifying the local environment is the bond order parameter (BOA) [51]. We employ the averaged BOA, usually denoted by two-dimension vectors  $\bar{q}_4$  and  $\bar{q}_6$  [52]. This parameter can be a simple but good indicator for quantifying the diversity of the local atomic environment.

Figures A.4 and A.5 shows the averaged BOA of SiN and HfO, respectively. SiN shows distinct phase separation in the  $\bar{q}_4 - \bar{q}_6$  plane, with most BOA values falling between  $0 \leq \bar{q}_4 \leq 0.3$  and  $0 \leq \bar{q}_6 \leq 0.3$ . HfO also exhibits similar phase separation, with unique patterns for each MD scenario. For the high-temperature pre-melting step, the phase boundary is separated by a change in volume. For the melting and quenching steps, their coverage overlaps. However, the energy and force distribution among these MD scenarios are different as shown in Figure A.3. This means that the



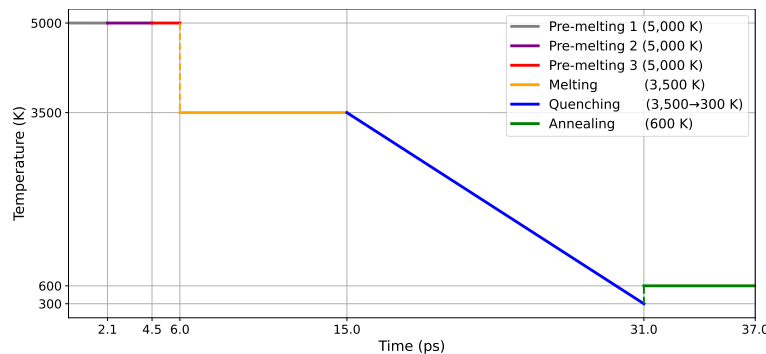


Figure A.6: A modified Melting-Quenching-Annealing scheme used for generating HfO dataset.

local atomic environment also differs even though  $\bar{q}_4$  and  $\bar{q}_6$  shows similar distribution. It is important to acknowledge that the BOA does not take into account atomic species and distance information, which can limit its reliability as a comprehensive descriptor. However, this analysis can provide an initial insight into the dataset's quality, serving as a starting point for further investigations.

### A.3 Utilization of Large-Scale Structural Databases

The Materials Project [42] and AFLOW [43] are comprehensive databases that encompass a wide range of atomic structures, spanning various elements and structural motifs. These structures are predominantly derived from precise quantum mechanical calculations and represent materials in their energetically stable states, where the forces on most atoms approach zero.

For those diving into MD simulations via MLFF models, understanding the nature of these databases is crucial. They primarily present configurations where atoms are at their stable states with minimum-energy. Molecular dynamics simulations, on the other hand, often explore high-energy conditions and a variety of force environments. Such scenarios are underrepresented in these databases, making them less suitable for training MLFF models directly. However, with their extensive collection of stable crystal structures, these databases are invaluable for selecting initial configurations when curating an MLFF training dataset.

### A.4 Details of the Modified Melt-Quench-Annealing Method

The generation of the HfO dataset involved the use of a modified melt-quench-annealing (m-MQA) method, where the temperature was scheduled as depicted in Figure A.6. In the m-MQA method, aimed at obtaining high entropy structures, we incorporated three pre-melting stages at an exceedingly high temperature of 5000 K. These stages involved pre-melting at crystal's original volume (pre-melting 1), followed by 10% isotropic compression (pre-melting 2) and 10% expansion (pre-melting 3) relative to the original volume. These steps were taken to ensure the presence of both dense and sparse atomic environments in the resulting dataset. Initial structures underwent three pre-melting steps at 5000 K for 6 ps, melting at 3500 K for 9 ps, quenching from 3500 to 300 K for 16 ps. Finally, the structures were restored to their original volume and annealed at 600 K for 6 ps. All simulations were performed under the NVT ensemble with the Nose-Hoover thermostat.

### A.5 Dataset Instruction

The extended-xyz format begins with a line declaring the total number of atoms. The subsequent line typically encompasses lattice information, energy, and other relevant details associated with the atomic structure for that specific simulation step. From the third line, each line represents an individual atom: the first column designates the atomic type, followed by its Cartesian coordinates and force components. This structured pattern is interspersed throughout the file, capturing snapshots

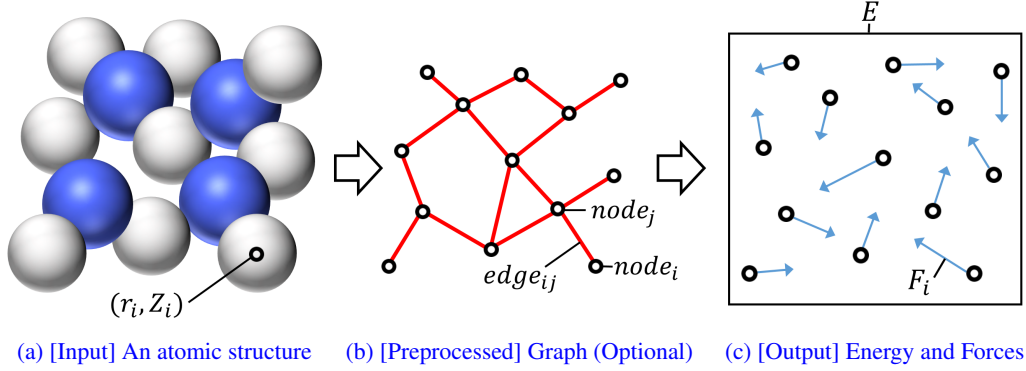


Figure B.1: General flow of MLFF models.

from various points in the simulation, providing a comprehensive record that offers insights into the dynamic behavior of atoms at different stages.

While the extended-xyz format can be handled with conventional text-processing utilities, employing specialized computational toolkits considerably streamlines the process. The Atomic Simulation Environment (ASE) [53], renowned for its versatility, effortlessly accommodates a wide range of atomistic formats, including the extended-xyz. From intuitive data visualization to comprehensive atomic analyses, its functionality significantly enhances user efficiency. Similarly, the Python Materials Genomics (pymatgen) [54] presents itself as an invaluable tool for researchers focused on atomic structural analysis. It offers deep insights into atomic structural data, enabling seamless structural analysis and transformations. Both ASE and pymatgen serve as commendable aids, optimizing the user experience in managing atomic structural databases.

## B Details of Benchmark Models

In this Section, we briefly review the MLFF models used in this benchmark. Figure B.1 describes general flow of MLFF models. Figure B.1a illustrates a structure of input data, which is an 3D point cloud data consist of coordinates  $r_i = (x_i, y_i, z_i) \in \mathbb{R}^3$ , and atom numbers  $Z_i \in \mathbb{Z}$  for  $1 \leq i \leq n$  where  $n$  represents the number of atoms. Figure B.1b describes the preprocessed data with 3D graph form, by constructing edges between points within a cutoff radius. As shown in Figure B.1c, model output consists of total system energy  $E \in \mathbb{R}$  and force  $F_i = (F_i^x, F_i^y, F_i^z) \in \mathbb{R}^3$ .

**BPNN [13]** As a pioneer study for MLFF, BPNN introduced a base idea for the MLFF field: the total energy is represented by a sum of atomic features. For a regression model, a neural network consisting of fully-connected layers is employed. To train the neural network, the coordinates of atoms are converted into hand-crafted features by using symmetry functions to describe local environment of atoms; the hand-crafted features are generally referred to as descriptors.

**DPA-1 [49]** DPA-1, which was specifically designed to be operated on DeePMD-kit [48], is an end-to-end deep potential energy model that includes trainable descriptors. These descriptors are invariant under translation, rotation, and permutation. It also employs a self-attention mechanism to effectively incorporate angular information for improved performance.

**SchNet [16]** SchNet is a deep neural network to adaptively learn atomic features representing local atomic environment from atom-centered symmetry functions without relying on pre-defined descriptors. It combines atomic intermediate features using continuous-filter convolutions to mimic atomic interactions.

**DimeNet++ [18]** DimeNet [17] is an MLFF model based on graph neural networks (GNNs), where directional information of atomic environment is embedded into directional messages. It utilizes distance features and dihedral angle features based on radial basis function and spherical basis function, respectively. As an improved version of DimeNet, DimeNet++ introduce an efficient implementation of the directional message passing layer of DimeNet, and its neural architecture is modified from that of DimeNet, enhancing predictive capabilities.

**GemNet (GemNet-T and -dT) [19]** GemNet is a GNN-based model that incorporates directed edge embeddings and two-hop message passing, allowing the model to capture complex directional information. Quadruplet, a tuple with four atoms, can be used for these scheme; GemNet using the quadruplet is called GemNet-Q. In GemNet-T, the message and intermediate features are described by using atom pairs and triplets, which is a tuple with three atoms. As an alternative to calculate forces, direction force prediction method is proposed and applied into GemNet-T arhitecture, which is GemNet-dT. GemNet achieves universal approximation capabilities for translation-invariant and permutation/rotation-equivariant predictions.

**NequIP [22]** NequIP utilizes E(3)-equivariant convolutions to capture interactions of geometric tensors, including vectors and higher-order tensors, of atom embeddings. NequIP can learn a representation of atomic environments more comprehensively, compared to models that use invariant convolutions and operate solely on scalars. Thus, it achieves prominent accuracy on a diverse range of molecules and materials while requiring significantly fewer training data.

**Allegro [23]** In GNN-based MLFF models, message passing is considered as a necessary technique to learn many-body interactions. However, the message passing needs information exchange among atoms, leading a large communication overhead on distributed computing system, which is unsuitable for large-scale simulation. To address this issue, Allegro combines equivariant features and tensor products to represent a strictly local equivariant atomic features without relying on message passing, achieving both accuracy and scalability.

**MACE [24]** MACE is theoretically based on multi-ACE framework [55] that was extended from atomic cluster expansion (ACE) framework by introducing GNNs equipped with high body-order messages. Due to the usage of high body-order messages, MACE can reduce message passing layers to learn the atomic potentials, enabling fast and parallelizable computations and enhancing the training efficiency.

**SCN [25]** SCN is a GNN-based model where atom features are a set of spherical channels represented by spherical harmonics, enabling it satisfy rotation equivariance. By relaxing the constraint of rotational equivariance in message passing and aggregation, SCN demonstrated that the performance of energy and force prediction can be improved. SCN is specifically designed for OC20 and requires heavy computation and memory usage. To deal with this issue, we reduced its model size by using 4 interaction (message passing) layers and 64 spherical channels; the small SCN was used for the all of our experiments.

## **C Details of Numerical Benchmark**

### **C.1 Error Metrics and Losses**

Energy and force errors are commonly employed as evaluation metrics to assess the performance of MLFF models, primarily due to their ease of use and straightforward interpretability. The errors and loss functions can be computed in several ways, which are listed in Table C.1. As in discussed in Section 4.1, we chose **the RMSE of per-atom energy**, which can be considered invariant among various snapshot sizes, and **the axis-wise RMSE of forces** (briefly called the RMSE of forces). In addition, we also list the **MSE-based loss** and **MAE-based loss**.

Table C.1: Error metrics and losses used in MLFF research and this benchmark.  $N$  is the number of snapshots,  $n_i$  is the number of atoms in snapshot with index  $i$ . The normal characters such as  $E_i$ ,  $f_x^{i,k}$  indicate reference data (*i.e.*, ground truth). The characters with a hat such as  $\hat{E}_i$ ,  $\hat{f}_x^{i,k}$  indicate predicted values by MLFF models.  $\lambda_f$  is a coefficient of force loss.

Type	Name	Formulation
Energy error metric	MAE of (total) energy	$\sum_{i=0}^N  E_i - \hat{E}_i  / N$
	RMSE of (total) energy	$\left[ \sum_{i=0}^N  E_i - \hat{E}_i ^2 / N \right]^{1/2}$
	MAE of per-atom energy	$\sum_{i=0}^N \left  \frac{E_i - \hat{E}_i}{n_i} \right  / N$
	RMSE of per-atom energy	$\left[ \sum_{i=0}^N \left  \frac{E_i - \hat{E}_i}{n_i} \right ^2 / N \right]^{1/2}$
Energy loss	MAE of per-atom energy	$\left  \frac{E_i - \hat{E}_i}{n_i} \right $
	MSE of per-atom energy	$\left  \frac{E_i - \hat{E}_i}{n_i} \right ^2$
Force error metric	Axis-wise MAE of force	$\sum_{i=0}^N \sum_{k=0}^{n_i} ( f_x^{i,k} - \hat{f}_x^{i,k}  +  f_y^{i,k} - \hat{f}_y^{i,k}  +  f_z^{i,k} - \hat{f}_z^{i,k} ) / \sum_{i=0}^N 3n_i$
	Axis-wise RMSE of force	$\left[ \sum_{i=0}^N \sum_{k=0}^{n_i} ( f_x^{i,k} - \hat{f}_x^{i,k} ^2 +  f_y^{i,k} - \hat{f}_y^{i,k} ^2 +  f_z^{i,k} - \hat{f}_z^{i,k} ^2) / \sum_{i=0}^N 3n_i \right]^{1/2}$
Force loss	Axis-wise MAE of force	$\sum_{k=0}^{n_i} ( f_x^{i,k} - \hat{f}_x^{i,k}  +  f_y^{i,k} - \hat{f}_y^{i,k}  +  f_z^{i,k} - \hat{f}_z^{i,k} ) / 3n_i$
	Axis-wise MSE of force	$\sum_{k=0}^{n_i} ( f_x^{i,k} - \hat{f}_x^{i,k} ^2 +  f_y^{i,k} - \hat{f}_y^{i,k} ^2 +  f_z^{i,k} - \hat{f}_z^{i,k} ^2) / 3n_i$
	L2-MAE of force	$\sum_{i=0}^{n_i} ( f_x^{i,k} - \hat{f}_x^{i,k} ^2 +  f_y^{i,k} - \hat{f}_y^{i,k} ^2 +  f_z^{i,k} - \hat{f}_z^{i,k} ^2)^{1/2} / n_i$
	L2-MSE of force	$\sum_{i=0}^{n_i} ( f_x^{i,k} - \hat{f}_x^{i,k} ^2 +  f_y^{i,k} - \hat{f}_y^{i,k} ^2 +  f_z^{i,k} - \hat{f}_z^{i,k} ^2) / n_i$
<b>EF metric</b> = RMSE of per-atom energy + Axis-wise RMSE of force		
<b>MAE-based loss</b> = MAE of per-atom energy + $\lambda_f \times$ L2-MAE of force		
<b>MSE-based loss</b> = MSE of per-atom energy + $\lambda_f \times$ Axis-wise MSE of force		

## 596 C.2 Training Hyperparameters

597 In the ML community, where the reproducibility of experiments is highly valued, sharing training  
598 recipes that include hyperparameters for training MLFF models is considered a valuable endeavor. In  
599 alignment with this attitude, we present the hyperparameters of the models, which we trained for this  
600 benchmark study. **The configuration files are included in codes of the supplementary materials.**  
601 By doing so, we aim to promote transparency and enable others to replicate and build upon our work  
602 in a meaningful way.

603 The all of configurations including neural architecture hyperparameters and optimization hyperpa-  
604 rameters are prepared into our codes in Supplementary materials: codes.zip (in configs/train/SiN/ and  
605 configs/train/HfO/). Also, the training factors are listed in Table C.2.

606 As a variant of Adam [56], AMSGrad [57] maintains the maximum of the past squared gradients  
607 instead of using exponential moving average, leading to more stable convergence. However, we  
608 empirically observed that training losses of MLFF models are often skyrocketed, resulting in relatively  
609 slow convergence. Thus, for models that used AMSGrad in the original papers, we tried to train  
610 models with or without AMSGrad, and then the better results are reported in this paper. Recently-  
611 proposed models tend to employ exponential moving average (EMA), which may also help to obtain  
612 stable training results.

613 While most of the hyperparameters of each model including the neural architecture introduced in the  
614 original paper (except SCN) were preserved, we established a common rule of some hyperparameters  
615 which can affect the training cost: epochs, batch size, and learning rate schedule. For both HfO and  
616 SiN datasets, MLFF models were trained during **200 epochs**. The training batch size for each model  
617 was selected from the options of 3, 4, 8, 16 to ensure compatibility with the memory capacity of the  
618 V100 GPU; the batch size selected by each model can be seen in Table C.2.

619 We trained models using a **linear decaying schedule** (LinearLR) that an learning rate is linearly  
620 decayed every step until the training is completed, *i.e.*, 200 epochs. The most existing MLFF models

Table C.2: Details of training factors for 10 MLFF models.

Model	Optimizer	Initial learning rate	Batch size		etc.
			SiN	HfO	
BPNN	Adam	5.e-3	16	16	weight decay: 1.e-6
DPA-1	Adam	1.e-3	16	16	
SchNet	Adam	1.e-4	16	16	
DimeNet++	Adam	1.e-4	3	8	AMSGrad, ema decay: 0.999, gradient clipping: 10
GemNet-T	AdamW	5.e-4	4	8	
GemNet-dT	AdamW	5.e-4	4	8	
NequIP	Adam	5.e-3	16	16	
Allegro	Adam	5.e-3	16	16	
MACE	Adam	1.e-2	16	16	ema decay: 0.99
SCN	AdamW	4.e-4	6	8	weight decay: 5.e-7, ema decay: 0.99
					AMSGrad, ema decay: 0.999

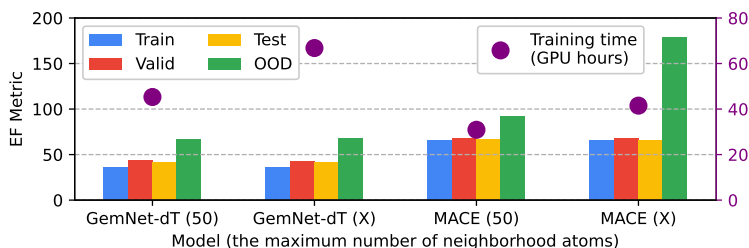


Figure C.1: EF metrics and training times of GemNet-dT and MACE trained with or without restricting the number of neighborhood atoms.

621 primarily employed one of the two types of learning rate schedules: ReduceLROnPlateau [19, 22,  
622 23, 24] and StepLR [17, 18, 25] implemented in PyTorch. The former stably gives decent training  
623 results but makes it difficult to estimate when the training will be completed. The latter allows for  
624 predicting the training completion point by specifying the training cost. However, it requires careful  
625 consideration of factors such as the decaying points, the decay factor, and the number of decay  
626 steps. Therefore, to control the training cost, we opted for a linear decaying schedule with solely one  
627 required hyperparameter, an initial learning rate, and thus we can control the training budget for fair  
628 benchmark. During training, the learning rate linearly decreases from the initial learning rate to 0 at  
629 each iteration.

### 630 C.3 Data Preprocessing

#### 631 C.3.1 Graph Generation

632 For the fairness comparison among MLFF models, we train MLFF models with a **radius cutoff of 6.0**  
633 and a maximum limit of **50 neighborhood atoms**. As illustrated in Figure C.1, despite the 1.5x and  
634 1.3x increase in training costs for GemNet-dT and MACE, respectively, the test EF metrics of the two  
635 models are hardly improved. This implies that the features of an atom, derived by message-passing  
636 in GNNs, may be more strongly influenced by the closer neighborhood atoms. In MACE, utilizing  
637 more neighborhood atoms does not lead to improvements, but instead results in higher EF metric on  
638 the OOD set. Thus, analyzing the effect of graph generation conditions is also important to obtain  
639 models that have lower OOD errors and are more suitable for simulations.

#### 640 C.3.2 Normalization for Energy and Forces

641 In this benchmark, when training MLFF models, we adopt a common normalization strategy for  
642 energy and forces. Here, we introduce normalization for energy and forces, and show our empirical  
643 results that present little difference regardless of a normalization strategy.

644 Normalizing energy and forces are commonly utilized in MLFF research, either through explicit  
645 normalization of energy and forces [29] or by incorporating scale and shift factors within the model

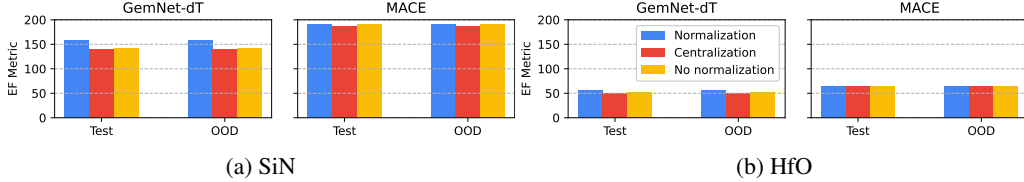


Figure C.2: EF metrics of GemNet-dT and MACE trained with different normalization strategies.

to fulfill the normalization role [22, 23, 24]. In this Section, after defining the normalization of energy and forces, we analyze the effect of the normalization to training results.

First, to introduce the normalization technique for energy and forces, we have to be aware that the force of an atom is equal to the derivative of total energy with respect to the atom position in various MLFF models. Calculating forces via taking the derivatives naturally induces two properties: sum of all forces in a single snapshot is equal to  $\mathbf{0}$ , and multiplying a constant to an energy affects exactly same to the forces. Fortunately, the first property implies that the forces does not require centralization (*i.e.*, shifting the forces by their mean,  $\mathbf{0}$ ). Thus, it suffices to scale the forces by dividing their standard deviation  $\sigma_f$  [22, 23, 24, 29]. Moreover, by the second property, the scaling factor of the energy should be identical to the scaling factor of the forces, *i.e.*,  $\sigma_f$ .

The only remaining part in defining the normalization is how to centralize the energy, and several methods for the centralization have been proposed. As discussed about error metrics and losses, the energy types to be targeted for obtaining the shift factor are two: total energy and per-atom energy.

The total energy can be centralized by  $\tilde{E}_i = E_i - \bar{E}$ , where  $E_i$  is the total energy of snapshot  $i$  and  $\bar{E}$  is the mean of total energies [29]. However, for MLFF models that calculate the total energy by predicting and summing up individual atomic energies, the centralization of total energy only works properly for the datasets, which consist of snapshots with the identical size.

As an alternative, some researches [22, 23, 24] centralize per-atom energy instead of the total energy, which leads to

$$\tilde{E}_i = E_i - \frac{n_i}{N} \sum_i \frac{E_i}{n_i}, \quad (2)$$

where  $n_i$  is the size of snapshot  $i$  and thus  $\frac{1}{N} \sum_i \frac{E_i}{n_i}$  is the mean of per-atom energies. When shifting the total energy, the mean of per-atom energies is compensated by multiplying the snapshot size. Since the SiN dataset consists of various snapshots with different number of atoms, Eq. 2 is employed for the centralization. Finally, the normalization of energy and forces for this benchmark is formulated as

$$\tilde{E}_i = \left( E_i - \frac{n_i}{N} \sum_i \frac{E_i}{n_i} \right) / \sigma_f, \quad \tilde{\mathbf{F}}_i = \mathbf{F}_i / \sigma_f \quad (3)$$

MACE [24] can use per-species normalization, which is an extension of the per-atom normalization, by computing per-species energy for each chemical species by solving linear system set by the datasets. However, such method is not appropriate to the datasets which only contain snapshots with uniform stoichiometric ratio, such as the HfO dataset; the linear system becomes singular, and thus the per-species energy is not determined.

To explore the influence of the normalization strategies on the training results, we trained models using three different strategies: per-atom normalization, per-atom centralization, no normalization. As presented in Figure C.2, the differences in errors obtained from models trained using these three methodologies are marginal, even in OOD errors. The energy prediction method employed by the MLFF models in this benchmark entails obtaining per-atom values within the model, which are then summed up to calculate the total energy of a snapshot. Thus, even not employing normalization, internal representations in MLFF models can inherently aligns with the per-atom values and be used in simulations for structures whose sizes are different from those in the training dataset. This is empirically demonstrated by the small error differences between test and OOD data.



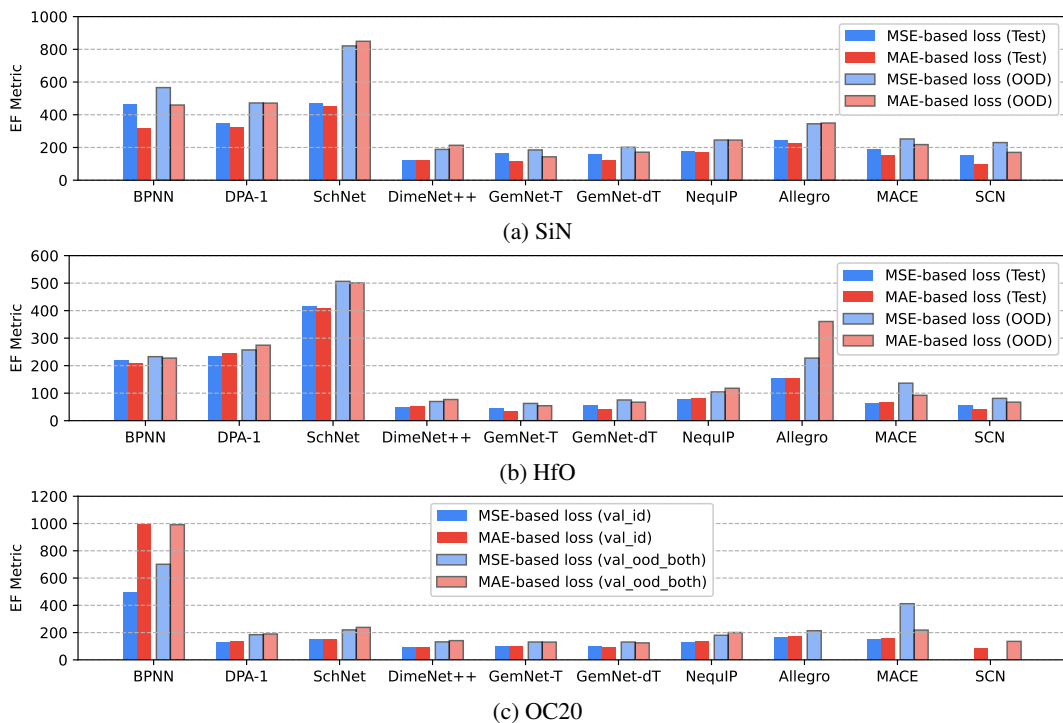


Figure C.3: EF metrics of 10 MLFF models trained with MSE- and MAE-based loss functions.

## C.4 Supplementary Results

### C.4.1 Results of Benchmark for Energy and Force

Table C.3 presents numerical results, which are visualized in Figure 1. As in discussed in Section 4.1, we trained MLFF models with two loss functions. For SiN and HfO datasets, we illustrate the training results of these two loss functions in Figure C.3 (a) and (b), respectively; blue and red bars represent the training results of MSE- and MAE-based loss functions, respectively. Section 4.2 covers the analysis about these results.

In addition, maintaining the same training recipes except setting the number of epochs as 50, we trained the models using OC20 dataset [29]; the training results are visualized in Figure C.3 (c). Training recipes of OC20 benchmark are different from ours. For instance, the target energy values of snapshots are adjusted by subtracting so-called reference energy, which corresponds to the energy of an initial snapshot with the same trajectory as the snapshots, and then apply the normalization technique for the shifted energy values; in contrast, we use the target energy values and apply per-atom energy normalization. Furthermore, DimeNet++ and GemNet-dT used for OC20 benchmark have larger architectures than those introduced in their original papers, while we utilize the original architectures; nevertheless, they show prominent results in energy and force prediction.

Overall, GNN-based models experience high-ranked performance. GNNs incorporating equivariance, such as Allegro and MACE, demonstrate inferior performance in energy and force prediction compared to SchNet, which is already considered the least favorable model in SiN and HfO datasets. For models except BPNN and Allegro, differences between the training results of two loss functions are less than those in SiN and HfO datasets, and EF metrics of using MAE-based loss are very slightly higher than using MSE-based loss. The results obtained from BPNN, where EF metrics of the model trained with MAE-based loss on both val\_id and val\_ood\_both are higher, imply that descriptor-based models may experience a significant performance drop, when dealing with datasets containing a wide range of atom species such as OC20. Allegro trained with MAE-based loss may be reasonably trained when observing the results of val\_id, but fails to correctly predict energy and forces on OOD

Table C.3: Numerical results (EF metrics) of MLFF models, visualized in Figures 1 and C.3.

Model	Loss	SiN			HfO			OC20		
		Train	Test	OOD	Train	Test	OOD	train	val_id	val_ood_both
BPNN	MSE-based	393.3	466.2	566.0	219.2	220.5	232.7	477.2	495.0	701.1
	MAE-based	291.4	316.2	459.4	206.4	210.2	227.4	984.2	991.3	991.5
DPA-1	MSE-based	335.0	347.1	472.0	232.7	234.0	257.1	127.8	128.9	184.6
	MAE-based	309.9	325.5	471.2	246.2	247.0	274.3	131.1	132.9	190.0
SchNet	MSE-based	192.0	471.2	820.7	171.5	417.4	506.4	145.3	147.2	219.8
	MAE-based	114.6	453.5	849.1	190.8	408.3	500.6	152.1	153.1	238.2
DimeNet++	MSE-based	84.3	119.8	188.3	37.1	48.0	69.7	102.1	88.2	131.9
	MAE-based	64.8	124.1	213.3	42.8	51.6	76.9	108.8	93.0	141.1
GemNet-T	MSE-based	144.5	164.9	184.9	43.0	45.8	62.7	113.4	102.4	130.9
	MAE-based	98.2	118.1	142.7	27.3	33.0	54.3	112.1	102.0	130.1
GemNet-dT	MSE-based	138.5	158.5	201.0	53.4	56.8	75.2	110.3	98.3	130.7
	MAE-based	87.9	120.4	171.1	36.3	41.9	67.2	103.4	91.4	124.5
NequIP	MSE-based	163.0	174.4	245.5	75.1	79.6	104.5	133.9	126.3	180.4
	MAE-based	157.9	170.4	245.1	79.8	82.6	117.8	140.7	135.0	200.4
Allegro	MSE-based	236.2	245.1	344.6	152.5	156.2	227.4	168.9	163.7	213.7
	MAE-based	209.0	225.4	349.3	152.7	156.0	360.5	169.4	169.6	2.1e+14
MACE	MSE-based	186.7	190.3	252.0	62.8	64.6	136.3	153.2	146.8	411.8
	MAE-based	147.1	155.1	217.7	65.7	66.2	92.2	164.5	158.5	218.4
SCN	MSE-based	138.2	153.5	230.0	50.4	56.9	81.1	-	-	-
	MAE-based	50.2	97.9	170.0	30.1	41.5	67.3	105.5	87.7	135.6

data of OC20, referred to as val\_ood\_both. It may be beneficial to analyze which factors of Allegro bring about such failure, because Allegro is compatible with the parallel computing of MD simulation tools such as LAMMPS while the other models cannot, which is important to enable large-scale simulations using MLFF models.

#### C.4.2 Results of Model Exploration Study

The numerical results of model exploration study, which are discussed in Section 4.3 and visualized in Figure 3, are listed in Table C.4, where the variation scales of each model and corresponding model sizes (*i.e.*, the number of parameters) are also included. The result of DPA-1 shows that DPA-1 can locate on the pareto-frontier in Figure 3. DPA-1, among the models evaluated in this study, is the only model that can be operated exclusively in DeePMD-kit [48], a TensorFlow-based framework. Therefore, the result of DPA-1 was excluded from Figure 1 because it may not be appropriate to directly compare that with the results of other models obtained from our PyTorch-based framework.

#### C.5 Data Scaling Effect

We additionally present the data scaling effect, which may be helpful to make a training strategy for MLFF researchers. We follow the experimental settings of previous works [22, 24], where the ReduceLROnPlateau scheduler is used to train MLFF models. Even though the training steps cannot be fixed due to this scheduler unlike our setting, we anticipate that the training might be stopped at similar update steps. Thus, in this data scaling experiments, MLFF models are trained by the fixed number of steps for model update; the setting is referred to as the equal budget setting. We randomly sample 20%, 40%, 60%, and 80% of data from the HfO training set (*i.e.*, 5.6k, 11.2k, 16.8k, and 22.4k), and trained models by 1000, 500, 334, and 250 epochs, correspondingly. We select six models (BPNN, GemNet-T, GemNet-dT, NequIP, Allegro, and MACE) and obtain their energy and force errors using the HfO test and OOD sets. Except the data size and epochs, the training recipe used in the MAE-based training experiments of Figure 1 is maintained.

The results are reported in Table C.5 and visualized in Figure C.4. For all models, as more training data are used, the training error increases while both the test error decreases, where the gap between these

Table C.4: Numerical results visualized in Figure 3 (HfO test set). Double, Base, Half, and Quarter indicate the range of the relative variation scale {2x, 1x, 0.5x, 0.25x}, correspondingly. <sup>†</sup>DPA-1 is specifically designed for use with DeePMD-kit, whereas the other models are implemented and evaluated within our framework.

Model	Variation scale	Number of Params.	EF metric	Inference time (ms)
BPNN	Double	1.92M	356.9	31.2
	Base	0.50M	356.1	31.5
	Half	0.13M	358.9	32.1
DPA-1 <sup>†</sup>	Base	6.14M	234.0	23.9
SchNet	Double	35.61M	743.6	24.3
	Base	9.09M	687.8	23.4
	Half	2.37M	627.0	23.1
	Quarter	0.64M	692.3	24.8
DimeNet++	Base	1.89M	81.8	62.5
	Half	0.48M	111.6	68.3
	Quarter	0.13M	145.2	68.9
GemNet-T	Base	1.89M	51.1	64.3
	Half	0.48M	67.6	62.1
	Quarter	0.13M	94.5	62.4
GemNet-dT	Double	9.14M	55.0	34.5
	Base	2.31M	64.7	29.5
	Half	0.59M	84.7	24.9
NequIP	Double	1.45M	113.9	62.1
	Base	0.36M	133.4	62.0
	Half	0.09M	160.0	61.3
	Quarter	0.02M	197.5	57.0
Allegro	Double	5.61M	243.1	48.7
	Base	1.40M	263.1	30.0
	Half	0.35M	289.9	27.4
MACE	Double	0.26M	86.1	35.6
	Base	0.11M	104.3	36.2
	Half	0.05M	128.3	35.9
SCN	Base	1.58M	62.2	327.2
	Half	0.47M	90.9	316.2
	Quarter	0.13M	183.3	332.2

two errors also decrease. The OOD error decreases in most cases. Therefore, our data scaling results suggest that training with fewer epochs and more data can be helpful to improve the generalization performance of MLFF models. Meanwhile, the difference between EF metrics using 80% and 100% training set is marginal, meaning that the performance improvement that can be gained by using more data seems to be saturated. Such observation implies that it is reasonable that our training set sampled from the raw dataset is sufficient for semiconductor MLFF benchmark.

We can opt for another data scaling experimental setup where the identical training epochs are used to train MLFF models; the setting is referred to as the equal epoch setting. In the equal epoch setting, as the training set size is reduced, the number of model update steps is also reduced. Considering that accuracy drop was clearly observed when using less data at the equal budget setting, we omit the experiments of the equal epoch setting because further accuracy drop is intuitively expected.

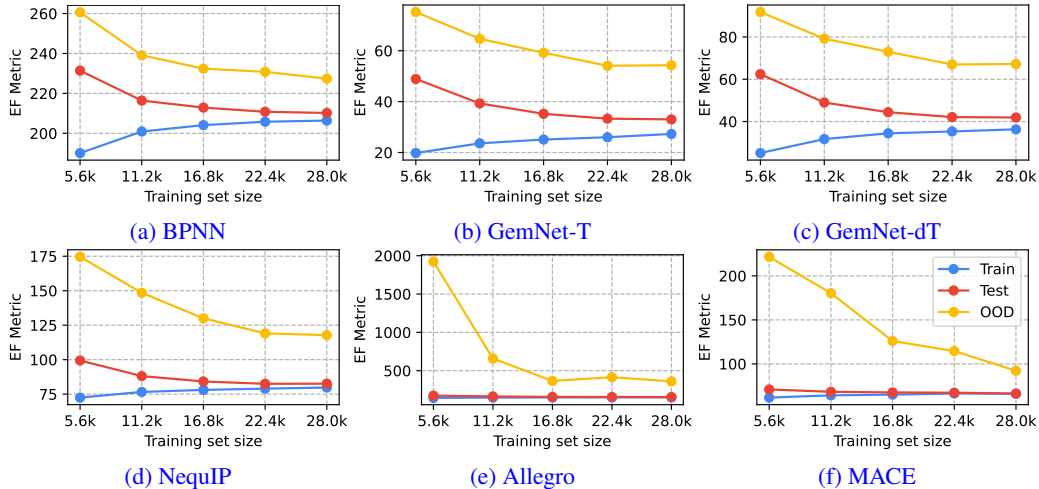


Figure C.4: Data scaling effect. The EF metrics are obtained from six MLFF models trained by using 20%, 40%, 60%, 80%, and 100% of the HfO training set. For all the points of Test and OOD curves, the test and OOD sets described in the main paper are identically used to obtain EF metrics. At each point of Train curves, the correspondingly sampled training set is used.

Table C.5: Numerical results visualized in Figure C.4.

Model	Training set size (ratio)	EF Metric ( $\downarrow$ )			Model	Training set size (ratio)	EF Metric ( $\downarrow$ )		
		Train	Test	OOD			Train	Test	OOD
BPNN	5.6k (20%)	190.1	231.4	260.6	NequIP	5.6k (20%)	72.5	99.4	174.6
	11.2k (40%)	200.9	216.4	239.1		11.2k (40%)	76.6	88.1	148.5
	16.8k (60%)	204.1	212.9	232.4		16.8k (60%)	78.1	84.2	130.0
	22.4k (80%)	205.8	210.8	230.8		22.4k (80%)	79.0	82.5	119.1
	28.0k (100%)	206.4	210.2	227.4		28.0k (100%)	79.8	82.6	117.8
GemNet-T	5.6k (20%)	19.8	48.9	75.2	Allegro	5.6k (20%)	144.7	174.6	1923.8
	11.2k (40%)	23.6	39.3	64.7		11.2k (40%)	148.9	163.8	658.3
	16.8k (60%)	25.1	35.2	59.2		16.8k (60%)	150.8	157.5	365.8
	22.4k (80%)	26.0	33.3	54.1		22.4k (80%)	152.0	156.4	415.3
	28.0k (100%)	27.3	33.0	54.3		28.0k (100%)	152.7	156.0	360.5
GemNet-dT	5.6k (20%)	25.1	62.4	91.8	MACE	5.6k (20%)	61.7	71.0	221.6
	11.2k (40%)	31.7	49.0	79.2		11.2k (40%)	64.3	68.3	180.4
	16.8k (60%)	34.4	44.4	73.0		16.8k (60%)	65.1	67.6	126.0
	22.4k (80%)	35.3	42.1	67.0		22.4k (80%)	66.4	67.3	114.6
	28.0k (100%)	36.3	41.9	67.2		28.0k (100%)	66.1	66.3	92.2

## D Details of Properties Benchmark

### D.1 Radial and Angular Distribution Functions

The dynamic indicators, namely RDF and ADF, are classified as such since they are derived from high-temperature MD simulation trajectories. In order to investigate the stability of models in different atomic environments induced by high turbulence or active movement caused by high thermal energy, we carefully selected high temperatures for SiN (1200 K) and HfO (1200 and 1800 K). RDF, also known as the pair correlation function, captures the changes in density relative to the distance from a chosen reference particle. On the other hand, ADF expands the scope of analysis beyond radial distances and centers on characterizing the angular distribution of particles surrounding a reference particle. The analysis of RDF and ADF plays a fundamental role in simulations as it enables the identification of structures, phases, and interactions. Moreover, it provides a deeper understanding of behaviors and reactions occurring within solids.

The radial distribution function,  $g(r)$  is an average density of particles within distance  $r$ , which is represented by

$$g(r) = \frac{dn_r}{4\pi\rho r^2 dr}, \quad (4)$$

where  $\rho$  indicates an average density of the system.

For compound AB, there could be three types of RDF results: A-A, A-B, and B-B. Mean Absolute Error (MAE) was calculated against the DFT reference for each case. The resulting three MAE values were then averaged to obtain the RDF error for each specific structure. Similarly, ADF could yield six types of results: A-A-A, B-B-B, A-B-A, B-A-B, A-A-B, and B-B-A. For each ADF case, the MAE was calculated by comparing it to the respective DFT reference. The resulting six MAE values were then averaged to obtain the ADF error for each specific structure.

To compare two distribution functions, we employ L1 distance between distribution functions. Thus an error metric of a RDF  $g(r)$  generated by a MLFF is defined by

$$E_{\text{RDF}} = d(g, g_{\text{DFT}}) = \frac{1}{R_c} \int_0^{R_c} |g(r) - g_{\text{DFT}}(r)| dr \quad (5)$$

where,  $R_c$  is equal to  $6\text{\AA}$ .

Similarly, an error of an ADF  $h(\theta)$  is defined by

$$E_{\text{ADF}} = d(h, h_{\text{DFT}}) = \frac{1}{\pi} \int_0^\pi |h(\theta) - h_{\text{DFT}}(\theta)| d\theta. \quad (6)$$

In both datasets, high-temperature simulations with various initial structures were conducted for a duration of 6 ps. Ground truth trajectories were computed using DFT. RDF and ADF were computed by averaging the last 3 ps of the simulation; a total of 50 snapshots were averaged, with an interval of 6 fs between each snapshot. The simulated structures employed in this study were supercells, encompassing a varying number of atoms. Specifically, the supercells comprised a range of 2,592 to 3,456 atoms for HfO and 1,296 to 2,835 atoms for SiN (refer to Table D.1 for additional information). This range of atom counts facilitates future scalability studies of MLFF. Additionally, users have the flexibility to increase supercell sizes as desired.

For the SiN, four different structures were simulated at 1200 K. Two structures with Si:N ratios of 3:4 and 1:1, which were within the range of ratios covered by the training dataset, were used to compute  $\text{RDF}^{\text{ID}}$  and  $\text{ADF}^{\text{ID}}$ . The  $\text{RDF}^{\text{ID}}$  value was calculated by averaging the RDF errors of the two structures, and likewise, the  $\text{ADF}^{\text{ID}}$  value was obtained by averaging the ADF errors of the two structures. Additionally, two structures with Si:N ratios of 1:2 and 3:2, which were outside the range of ratios present in the training dataset, were employed to calculate  $\text{RDF}^{\text{OOD}}$  and  $\text{ADF}^{\text{OOD}}$ .

A total of eight structures were simulated for the hafnium oxide case. Among these, five structures were simulated at 1200 and 1800 K (resulting in a total of 10 simulations) using the same stoichiometry (Hf:O = 1:2) as in the training set. These simulation trajectories were then used to calculate  $\text{RDF}^{\text{ID}}$  and  $\text{ADF}^{\text{ID}}$ . On the other hand, the remaining three structures were considered out-of-distribution structures with varying stoichiometries of Hf:O, specifically 1:1, 2:3, and 4:7. These three structures were simulated at 1200 K and then employed to calculate  $\text{RDF}^{\text{OOD}}$  and  $\text{ADF}^{\text{OOD}}$ .

### D.1.1 Understanding RDF & ADF: Insight, Motive, and Explanation

In simulation studies, understanding atomic position patterns can be complex. To solve this, researchers turn to RDF and ADF, two pivotal post-processing techniques, to differentiate between materials or states. RDF quantifies atomic or molecular density fluctuations based on distance, providing insights into localized structural properties. In contrast, ADF elucidates angular tendencies between particle triplets, offering a perspective on molecular shapes and bonding angles. Collectively, they act as benchmarks, verifying the real-world alignment of simulations.

Table D.1: Details of evaluated structures for RDF and ADF. Structures with index (idx.) indicate the same crystal family but have different lattice constants. The index numbering has been restarted for each dataset (e.g., Cubic (1) of SiN and Cubic (1) of HfO have different lattice constants).

	Distribution	Structure (idx.)	Hf:O (or Si:N)	# atoms ( $3 \times 3 \times 3$ )	Temp.
SiN	ID	Amorphous	1:1	1296	1200 K
		Triclinic	3:4	2835	1200 K
	OOD	Cubic (1)	1:2	2592	1200 K
		Cubic (2)	3:2	2720	1200 K
HfO	ID	Monoclinic	1:2	2592	1200 & 1800 K
		Tetragonal	1:2	2592	1200 & 1800 K
		Cubic (1)	1:2	2592	1200 & 1800 K
		Orthorhombic (1)	1:2	2592	1200 & 1800 K
		Orthorhombic (2)	1:2	2592	1200 & 1800 K
	OOD	Hexagonal (1)	1:1	3456	1200 K
		Hexagonal (2)	2:3	3240	1200 K
		Cubic (2)	4:7	2376	1200 K

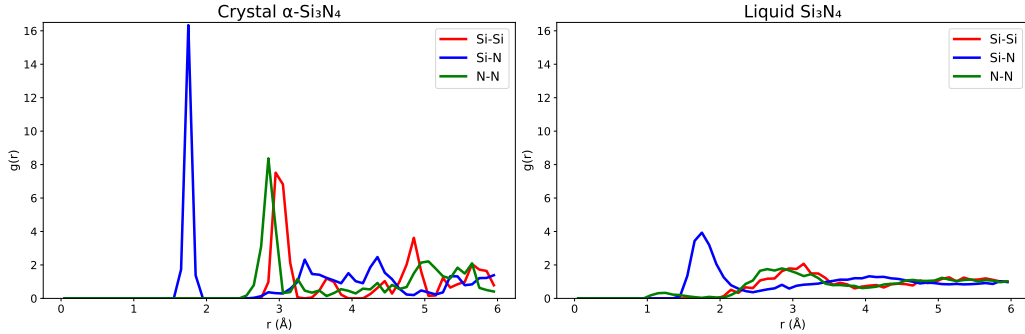


Figure D.1: RDFs of  $\alpha\text{-Si}_3\text{N}_4$  (left) and liquid- $\text{Si}_3\text{N}_4$  (right).

At the heart of MLFF models lies the ability to accurately predict energies and forces governing atoms and molecules. These predictions are crucial as they directly influence the movement, behavior, and interactions of atoms, which, in turn, dictate the patterns seen in RDF and ADF. When the RDF and ADF predicted by an MLFF closely align with those from DFT calculations, it's a strong indication that the model is capturing the essential physics and is making accurate energy and force predictions.

RDF and ADF can be used to capture phenomena such as phase transitions, system characteristics, and intermolecular dynamics. A prime example is the simulation of the silicon nitride phase shift from solid to liquid, where stark RDF modifications are observable in Figure D.1. What begins as distinct peaks in its crystalline form smoothes out during the melting process, signifying inherent structural shifts. Figure D.1 contrasts the RDF patterns of crystalline silicon during its initial state and post-melting phase, underlining the transformative effects of atomic rearrangements. Such configurations inherently relate to the balance of interatomic forces and the resultant energy landscape. When atoms approach each other under specific conditions, they might form bonds due to attraction or repel, settling into stable distances. This equilibrium gives rise to distinctive RDF and ADF profiles that mirror these atomic positions.

While RDF and ADF primarily help visualize atomic configurations, they also play a crucial role in assessing the fidelity of atomic paths predicted by MLFF. For example, In Figure D.3-(c), there is a noticeable anomaly in the RDF for Hf-Hf interactions on the OOD Hexagonal (2) structure when obtained using the MAE-based SchNet model. This anomaly is characterized by a pronounced peak at unusually short ' $r$ ' distances. This unlikely closeness between Hf-Hf atoms, especially when compared to ground truth data, not only challenges physical rationale but also signals potential weaknesses in the predictive capability of the MLFF model.



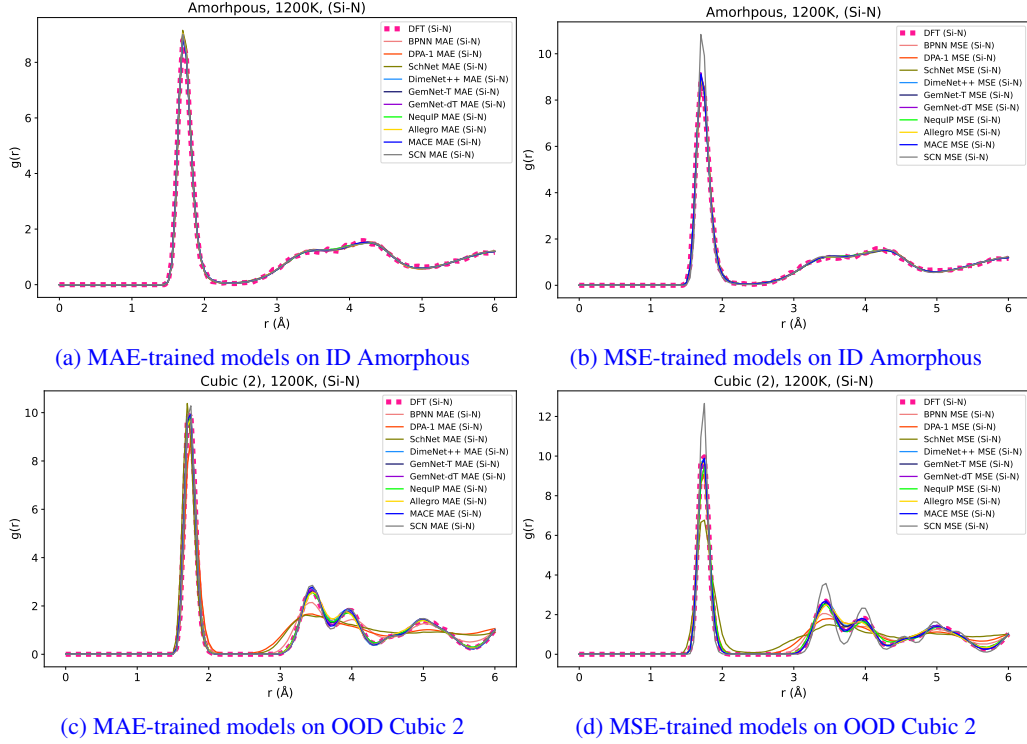


Figure D.2: RDF comparisons: DFT ground truth vs. model predictions for SiN in ID Amorphous and OOD Cubic 2 structures.

In our framework, we introduce straightforward evaluations of RDF and ADF for MLFF models by contrasting them against reference profiles from gold-standard DFT simulations (Table D.1). A close match between the MLFF-derived profiles and these references signifies model accuracy. Though a smaller error naturally indicates superior model performance, for visualization in the form of radar plots as shown in Figures D.13 and D.14, we inversely transform the RDF and ADF errors such that a value of 1 denotes optimal performance. The specifics of this transformation are detailed in Section D.4.

## D.2 Bulk Modulus and Equilibrium Volume

The bulk modulus and equilibrium volume derived from the Birch-Murnaghan equation of state (EoS) are critical in investigating thin film materials, including silicon nitride and hafnium oxide. The bulk modulus provides insights into the mechanical stability of these films, indicating their resistance to volume shifts under varying pressure conditions and their resilience against external stresses. Additionally, it characterizes the elastic properties of the films, impacting aspects like flexibility and hardness. Meanwhile, the equilibrium volume aids in understanding the relationship between the film's thickness and residual stress, enabling us to evaluate the film's state - strained or relaxed, and refine the growth procedure. Furthermore, both these parameters influence the interfacial properties such as adhesion strength and interface energy, integral for the performance and stability of thin film systems. Comprehension of these parameters is essential for engineering thin films with the required properties, for their potential application in electronics, optoelectronics, and microelectromechanical systems.

The bulk modulus is a fundamental quantity that measures a material's resistance to compression and expansion. It quantifies the material's ability to withstand changes in volume when subjected to external pressure. Comparing the bulk moduli of different materials allows for the assessment of their relative compressibility and overall mechanical strength. Additionally, the equilibrium volume, obtained through the fitting of the Birch-Murnaghan EoS, represents the volume at which the material

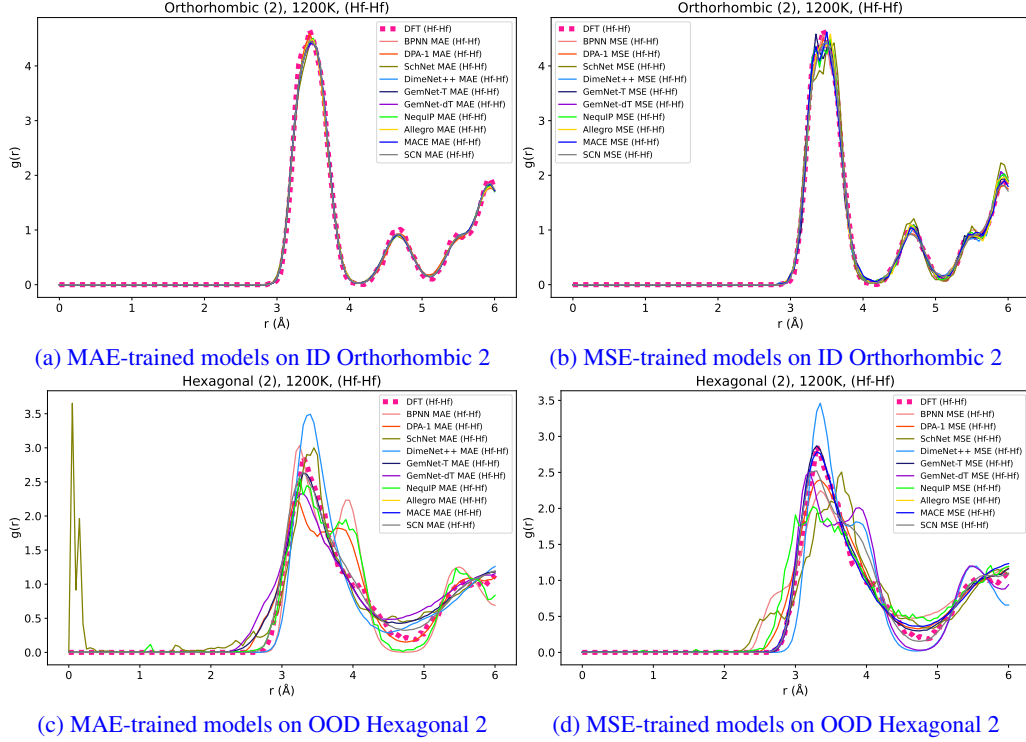


Figure D.3: RDF comparisons: DFT ground truth vs. model predictions for HfO in ID Orthorhombic 2 and OOD Hexagonal 2 structures.

possesses the minimum energy. Determining the equilibrium volume enables the identification of the most energetically stable configuration of the material, providing insights into its structural stability and phase transitions.

The bulk modulus and equilibrium volume can be obtained by fitting the Birch-Murnaghan EoS. The EoS establishes a relationship between the energy and volume of a solid material, and it can be expressed as follows:

$$E(V) = E_0 + \frac{9V_0B_0}{16} \left\{ \left[ \left( \frac{V_0}{V} \right)^{2/3} - 1 \right]^3 B'_0 + \left[ \left( \frac{V_0}{V} \right)^{2/3} - 1 \right]^2 \left[ 6 - 4 \left( \frac{V_0}{V} \right)^{2/3} \right] \right\}, \quad (7)$$

where  $V_0$  is the equilibrium volume,  $B_0$  is the bulk modulus, and  $B'_0$  is the derivative of the modulus with respect to pressure.

The internal energy versus volume data obtained from DFT calculations served as reference to assess the accuracy of MLFF models. In the case of HfO, five crystalline structures present in the training set were utilized to calculate  $B_0^{\text{ID}}$  and  $V_0^{\text{ID}}$ . Additionally, one crystal with different stoichiometry, not included in the training set, was employed to determine  $B_0^{\text{OOD}}$  and  $V_0^{\text{OOD}}$ . Similarly, in the case of SiN, five different crystalline structures with varying numbers of atoms were used in the EoS evaluation. Among these structures, three were included in the training set and employed to calculate  $B_0^{\text{ID}}$  and  $V_0^{\text{ID}}$ . The remaining two structures, which was not commonly found in the training set, was utilized to determine  $B_0^{\text{OOD}}$  and  $V_0^{\text{OOD}}$ . For more detailed information regarding the structures evaluated please refer to Table D.2. The values of  $B_0^{\text{ID}}$ ,  $B_0^{\text{OOD}}$ ,  $V_0^{\text{ID}}$ , and  $V_0^{\text{OOD}}$  were determined by calculating the arithmetic mean of the absolute percentage errors for the evaluated structures.

## D.2.1 Understanding the Equation of State: Insight, Motive, and Explanation

Our paper introduces the equation of state as an evaluation method, with metrics derived from the EoS, specifically the bulk modulus ( $B_0$ ) and equilibrium volume ( $V_0$ ), serving as the primary assessment

Table D.2: Details of evaluated structures for EoS.

	Distribution	Space Group	Hf-O (or Si-N) ratio	Number of atoms
SiN	ID	P31c	3:4	28
		P1	3:4	14
		P1	3:4	28
	OOD	Fd $\bar{3}$ m	3:4	14
		Fd $\bar{3}$ m	3:4	56
HfO	ID	P2 <sub>1</sub> /c	1:2	96
		P4 <sub>2</sub> /nmc	1:2	96
		Fm $\bar{3}$ m	1:2	96
		Pbca	1:2	96
		Pnma	1:2	96
	OOD	Fd $\bar{3}$ m	4:7	88

metrics. These metrics offer a way to assess how closely the predictions of the MLFF models align with the results from more established methods, like DFT.

The bulk modulus provides insights into a material’s resistance to compression. In simple terms, it helps us understand how much a material can be compressed. A higher value indicates that the material is less likely to change its volume under pressure. The equilibrium volume, on the other hand, tells us about the optimal volume where a material is most stable and uses energy most efficiently.

The core requirement of the MLFF model is its capability to simulate atomic interactions. When these interactions are observed at a larger scale, they manifest as properties such as the material’s resistance to compression and its most stable form. These properties are effectively represented by metrics like the bulk modulus and equilibrium volume. This suggests that a model which can accurately predict atomic interactions is also likely to be proficient at predicting these larger-scale properties.

An important validation of the MLFF model’s effectiveness is its alignment with DFT results. If our MLFF model’s predictions are consistent with DFT-calculated metrics, it indicates that our model has a good grasp on the intricate interactions between atoms. For ML experts, the EoS metrics serve a dual purpose. Not only do they offer a straightforward and clear measure of the model’s performance, but they also bridge the gap between intricate atomic simulations and tangible material properties, making the underlying science more accessible and understandable.

In the real experiment, when the mismatch of lattice constant ( $V_0^{1/3}$ ) exceeds over 5% between materials, it becomes very difficult to synthesize the materials. Thus, a reasonable MLFF model should predict volume accurately within this range. Since the GGA-level of DFT calculation already contains volume error of 2-3%[58], the error on the lattice constant should be less than 3%, which suggest that error on  $V_0$  ( $|V_{0,DFT} - V_{0,MLFF}|/V_{0,DFT}$ ) be less than 0.1.

Our framework offers an easy-to-use evaluation and reference data for a diverse range of solid structures (Table D.2). It not only automates the calculation of the bulk modulus and the equilibrium volume but also provides comparative DFT results as a gold standard. As depicted in these Figures D.4 and D.5, less effective MLFF models may struggle in accurately predicting atomic interactions when there’s a change in volume due to varying pressures. This misalignment in prediction can lead to deviations in the energy shift graph from the ideal DFT profile. Such discrepancies indicate the model’s challenges in capturing energy variations associated with changes in the internal atomic configurations, highlighting potential inaccuracies in forecasting the material’s mechanical properties. This approach gives us a broader perspective, enabling an assessment of MLFF’s capabilities beyond merely energy and force predictions, and demonstrates its application in real-world simulation scenarios.

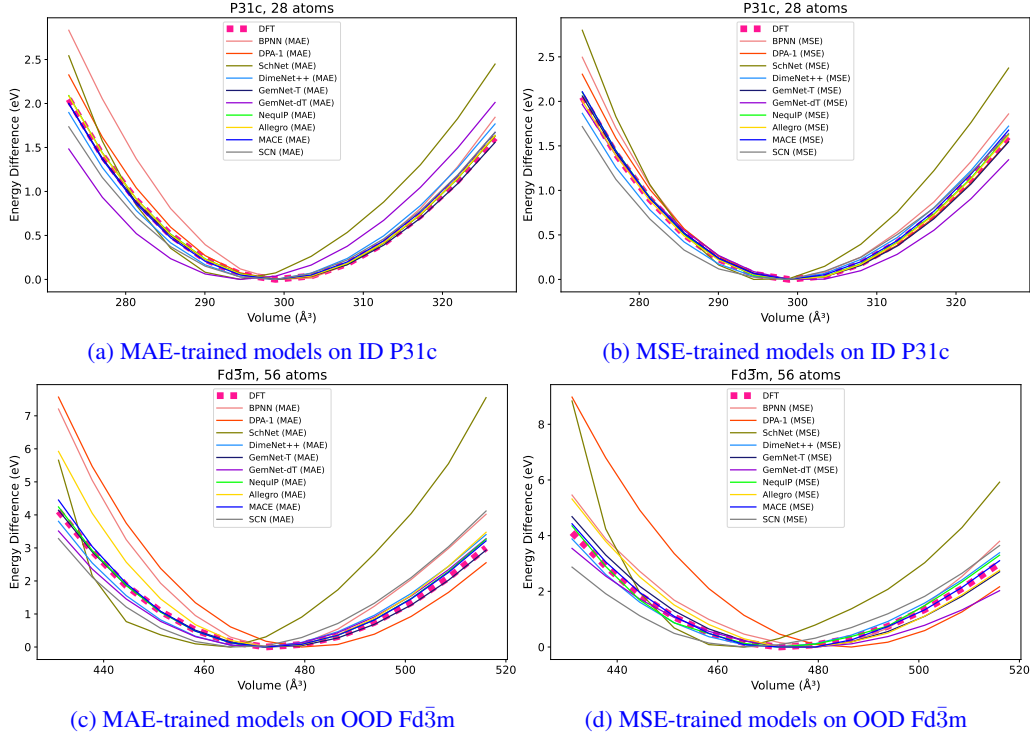


Figure D.4: EoS comparisons: DFT ground truth vs. model predictions for SiN in P31c and OOD Fd $\bar{3}m$  structures. Each curve has been shifted to reflect the energy difference relative to its minimum value.

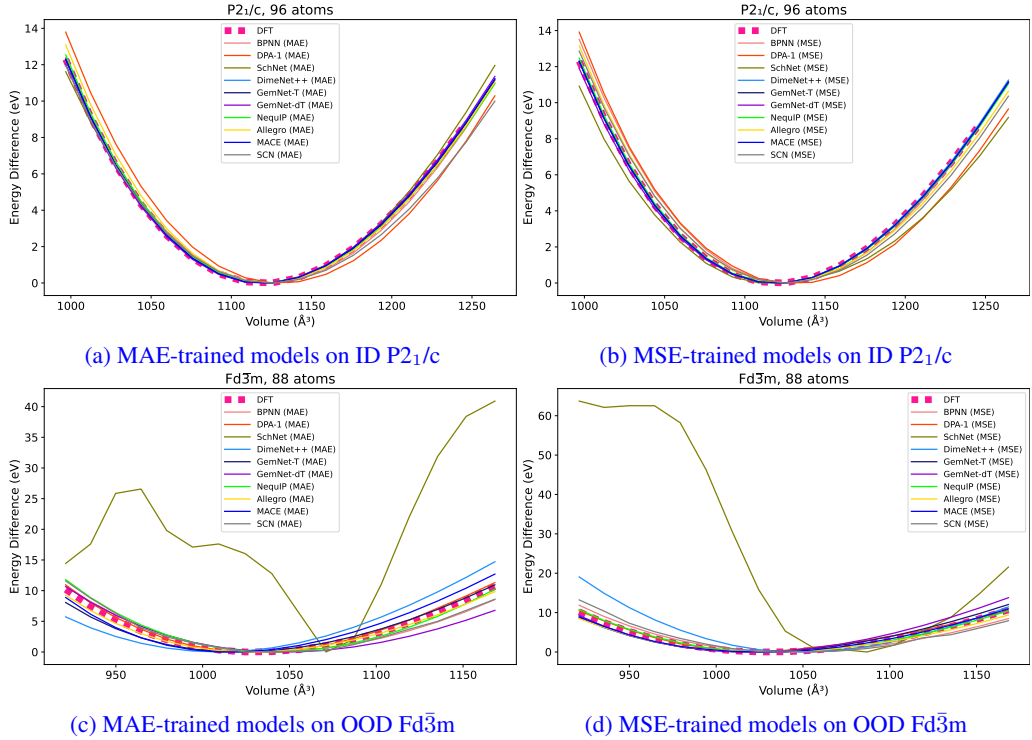


Figure D.5: EoS comparisons: DFT ground truth vs. model predictions for HfO in P2 $_1/c$  and OOD Fd $\bar{3}m$  structures. Each curve has been shifted to reflect the energy difference relative to its minimum value.

Table D.3: Details of PEC evaluations.

		Formula	Variable*	Range
SiN	two-body	Si-Si	distance (r)	1.000 $\sim$ 6.000 Å
		Si-N	distance (r)	1.000 $\sim$ 6.000 Å
		N-N	distance (r)	0.500 $\sim$ 6.000 Å
	many-body	Si <sub>3</sub> N	distance (r)	2.000 $\sim$ 6.000 Å
		SiN <sub>4</sub>	distance (r)	2.000 $\sim$ 6.000 Å
HfO	two-body	Hf-Hf	distance (r)	1.000 $\sim$ 6.000 Å
		Hf-O	distance (r)	1.000 $\sim$ 6.000 Å
		O-O	distance (r)	0.500 $\sim$ 6.000 Å
	many-body	HfO <sub>2</sub>	distance (r)	0.986 $\sim$ 6.786 Å
		HfO <sub>3</sub>	angle ( $\theta$ )	0 $\sim$ 140°
		HfO <sub>4</sub>	angle ( $\theta$ )	0 $\sim$ 180°

\*Please refer Figure D.6

### D.3 Potential Energy Curves

Potential curves for two-body interatomic interactions were computed by manipulating the distance between a pair of atoms. Additionally, specific MLFF models, which incorporate many-body terms, were assessed with an array of molecular structures that encompassed more than just two atoms. These many-body potential energy curves were derived by adjusting particular distances or angles within these molecular structures. In these evaluations, the ground truth obtained from DFT calculations was used as a reference. For detailed information on these evaluations, please refer to Table D.3 and Figure D.6. In the case of SiN many-body structures, we employed equilateral triangles (Si<sub>3</sub>N) and tetrahedron (SiN<sub>4</sub>) as they represent the atomic environment of each element with their first nearest neighbors in crystals. Additionally, the many-body structures of HfO were prepared by referencing a study on monohafnium oxide clusters, which is relevant for understanding defect sites in HfO thin films [59]. PEC evaluation has the potential to be a valuable tool in the solid-state domain, particularly in scenarios where certain chemical reactions or electrochemical properties are expected to take place in rare atomic environments.

#### D.3.1 Understanding PEC: Insight, Motive, and Explanation

MLFFs have emerged as a significant tool in computational materials science. Through data-driven processes, they can precisely predict energies and forces within molecular configurations. Their performance is considered reliable when they align well with PECs from high-fidelity DFT calculations.

PECs have historically been essential for empirical interatomic potential model fitting. They have served as the cornerstone of empirical methodologies, used to understand energy changes based on varying atomic configurations. Given this foundational role of PECs, their alignment with modern MLFFs becomes crucial.

The alignment of an MLFF with high-fidelity DFT-derived PECs verifies its ability to capture complex interactions and its competence in modeling forces. Notably, forces are derived from the gradients of the molecular energy landscape. An MLFF's ability to predict energy landscapes, as validated by its congruence with DFT-derived PECs, ensures that the resulting molecular simulations are trustworthy. In summation, the alignment between MLFFs and DFT-derived PECs, especially across diverse molecular configurations, demonstrates the adaptability and effectiveness of the MLFF model.

A distinguishing feature of our framework is its user-friendly evaluation capability across a broad spectrum of molecular structures (Figure D.6), especially concerning PECs. While our main focus has been on condensed-phase systems, handling sparse molecular systems poses inherent challenges.

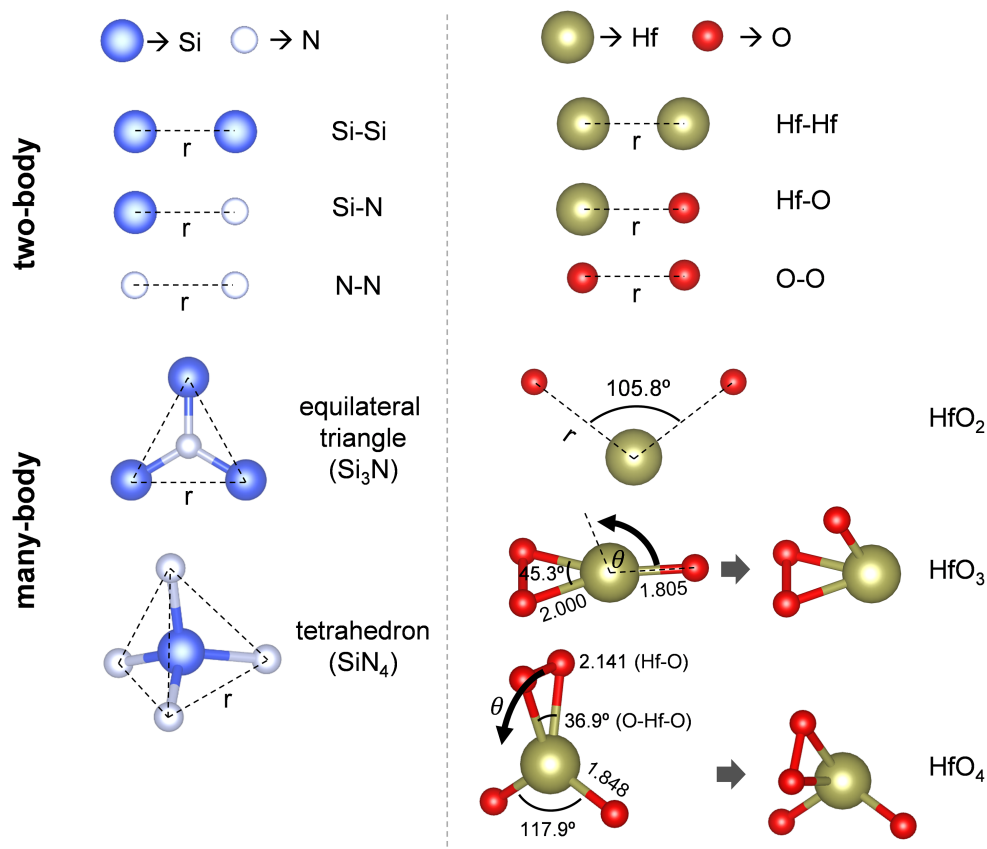


Figure D.6: Molecules used for PEC evaluations.

These sparse systems act as outliers for models predominantly trained on condensed-phase datasets. Successfully obtaining accurate PECs for such outliers underlines the robustness and versatility of the MLFF approach.

### D.3.2 PEC Evaluation

Figures D.9 to D.10 visualize the PEC evaluation results of models for SiN data, while Figures D.11 to D.12 depict the results for HfO. PEC evaluation poses a significant challenge as it involves assessing sparse atomic environments that are rarely encountered by models trained on condensed-phase datasets. While it is not essential for models used in condensed-phase material research to excel in PEC evaluation, the evaluation results provided valuable insights into the disparities between the two datasets. The SiN data included a wide range of atomic compositions and environments, including Si-only and N-only structures. In contrast, the HfO data exhibited a less diverse range. Consequently, the PEC evaluation results for SiN showed a closer resemblance to the DFT ground truth curve compared to the results for HfO.

To explain in a slightly more physics-oriented manner, given that most models predominantly train on low-energy data, predicting the high-energy PE surface becomes a challenging extrapolation task for these models. However, predicting this high-energy PE regime, usually occurring when the atoms get close, is important for stable MD simulation.

Two distinct methods can address this discrepancy in PE surface prediction. The first, termed as a data-centered approach, emphasizes the incorporation of more data from the high-energy regime, necessitating careful DFT calculations. The second approach, so called  $\delta$ -learning, tries to combine classical force-field model and MLFF models [60, 61, 62]. With this method, an initial approximation of the PE surface is derived using the classical force-field model, followed by the refinement of the



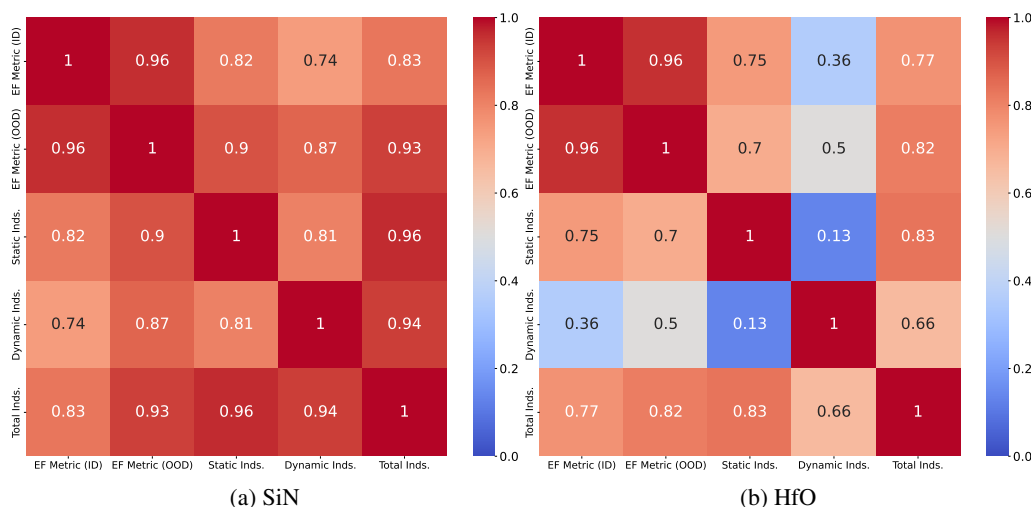


Figure D.7: Pearson correlation matrix illustrating the relationships between "EF Metric (ID)", "EF Metric (OOD)", "Static Inds.", "Dynamic Inds.", and "Total Inds" across two datasets: (a) SiN and (b) HfO. The "Static Inds." and "Dynamic Inds." for each model were computed by summing the indicator scores from both ID and OOD cases. The "Total Inds." represents the combined sum of the "Static Inds." and "Dynamic Inds." scores.

Table D.4: Raw scores of models for indicators using SiN dataset and trained with MAE-based loss. [Lower values indicate better performance.](#)

	EF mtr. <sup>ID</sup>	EF mtr. <sup>OOD</sup>	RDF <sup>ID</sup>	RDF <sup>OOD</sup>	ADF <sup>ID</sup>	ADF <sup>OOD</sup>	V <sub>0</sub> <sup>ID</sup>	V <sub>0</sub> <sup>OOD</sup>	B <sub>0</sub> <sup>ID</sup>	B <sub>0</sub> <sup>OOD</sup>
BPNN	3.2e+2	4.6e+2	4.7e-2	1.5e-1	3.5e-2	1.1e-1	5.4e-1	1.9e-1	2.1e+1	5.1e+1
DPA-1	3.3e+2	4.7e+2	5.0e-2	1.9e-1	3.4e-2	2.2e-1	1.6e-1	1.5e+0	6.0e+0	2.5e+1
SchNet	4.5e+2	8.5e+2	6.0e-2	1.1e+0	3.9e-2	4.3e-1	8.1e-1	1.9e+0	4.8e+1	9.3e+1
DimeNet++	1.2e+2	2.1e+2	4.2e-2	1.1e-1	3.2e-2	9.2e-2	2.8e-1	5.7e-1	1.4e+0	2.7e+0
GemNet-T	1.2e+2	1.4e+2	4.0e-2	3.3e-2	3.0e-2	2.1e-2	1.5e-1	1.1e-1	3.3e+0	1.2e+0
GemNet-dT	1.2e+2	1.7e+2	4.1e-2	6.9e-2	3.1e-2	5.9e-2	1.2e+0	5.8e-1	1.7e+0	3.2e+0
NequIP	1.7e+2	2.5e+2	4.1e-2	7.0e-2	3.1e-2	4.1e-2	8.3e-2	1.4e-1	1.4e+0	5.2e+0
Allegro	2.3e+2	3.5e+2	4.4e-2	4.2e-1	3.3e-2	2.3e-1	2.5e-2	8.2e-2	2.7e+0	2.8e+1
MACE	1.6e+2	2.2e+2	4.0e-2	4.4e-2	3.0e-2	3.5e-2	7.1e-2	1.2e-1	4.5e-1	6.9e+0
SCN	9.8e+1	1.7e+2	4.1e-2	5.7e-2	3.0e-2	4.4e-2	4.5e-1	1.3e+0	5.7e+0	7.9e+0

954 residual portion using an MLFF model. Both strategies present potential solutions when confronting  
955 inconsistencies in the PE surface predictions.

#### 956 D.4 Calculating Indicator Score: Formula and Details

957 For all simulation indicators and numerical metrics, lower values indicate better results than higher  
958 values. To draw radar plots for model comparison, an inverse linear transformation were applied to  
959 map them between 0 and 1, where a score of 1 indicates the best performance. Minimum observed  
960 error was considered as the perfect score of 1, serving as the ideal benchmark for each metric. In  
961 order to maintain reasonable values, we meticulously selected maximum thresholds for each metric.  
962 For detailed information regarding the transformation rule and maximum thresholds, please refer  
963 to Table D.8. The raw scores of the models prior to this mapping are summarized in Tables D.4  
964 to D.7. The details of calculating the raw score for each indicator are explained in Sections D.1 to D.3,  
965 along with evaluated structures and a description of simulation environments. The performance  
966 comparison among models based on indicators after the reverse mapping to scores ranging from 0 to  
967 1 is visualized using radar plots in Figures D.13 and D.14.

Table D.5: Raw scores of models for indicators using SiN dataset and trained with MSE-based loss. Lower values indicate better performance.

	EF mtr. <sup>ID</sup>	EF mtr. <sup>OOD</sup>	RDF <sup>ID</sup>	RDF <sup>OOD</sup>	ADF <sup>ID</sup>	ADF <sup>OOD</sup>	V <sub>0</sub> <sup>ID</sup>	V <sub>0</sub> <sup>OOD</sup>	B <sub>0</sub> <sup>ID</sup>	B <sub>0</sub> <sup>OOD</sup>
BPNN	4.7e+2	5.7e+2	5.0e-2	1.7e-1	3.3e-2	1.4e-1	4.7e-1	7.5e-1	1.2e+1	3.2e+1
DPA-1	3.5e+2	4.7e+2	5.0e-2	1.2e-1	3.5e-2	1.5e-1	3.6e-1	2.5e+0	2.0e+0	3.6e+1
SchNet	4.7e+2	8.2e+2	5.7e-2	1.1e+0	3.5e-2	4.9e-1	6.8e-1	1.8e+0	5.1e+1	1.5e+2
DimeNet++	1.2e+2	1.9e+2	4.5e-2	2.1e-1	3.1e-2	7.0e-2	1.9e-1	4.8e-1	5.0e-1	3.0e+0
GemNet-T	1.6e+2	1.8e+2	4.5e-2	3.6e-2	3.0e-2	2.7e-2	3.3e-1	5.3e-1	1.9e+0	1.2e+0
GemNet-dT	1.6e+2	2.0e+2	5.0e-2	1.1e-1	3.1e-2	7.2e-2	7.4e-1	8.0e-1	1.1e+1	2.3e+1
NequIP	1.7e+2	2.5e+2	4.5e-2	8.4e-2	3.1e-2	5.1e-2	1.6e-1	2.7e-1	2.4e+0	6.5e+0
Allegro	2.5e+2	3.4e+2	4.9e-2	1.4e-1	3.2e-2	1.1e-1	2.3e-1	7.2e-1	3.9e-1	9.8e+0
MACE	2.0e+2	2.6e+2	4.6e-2	3.7e-2	3.2e-2	2.4e-2	1.9e-1	6.3e-2	3.2e+0	6.0e+0
SCN	1.5e+2	2.3e+2	6.7e-2	1.2e-1	3.5e-2	1.3e-1	7.7e-1	1.3e+0	7.4e+0	5.8e+0

Table D.6: Raw scores of models for indicators using HfO dataset and trained with MAE-based loss. Lower values indicate better performance. N/A indicates an interrupted simulation for the dynamic indicator due to abnormal energy changes.

	EF mtr. <sup>ID</sup>	EF mtr. <sup>OOD</sup>	RDF <sup>ID</sup>	RDF <sup>OOD</sup>	ADF <sup>ID</sup>	ADF <sup>OOD</sup>	V <sub>0</sub> <sup>ID</sup>	V <sub>0</sub> <sup>OOD</sup>	B <sub>0</sub> <sup>ID</sup>	B <sub>0</sub> <sup>OOD</sup>
BPNN	2.1e+2	2.3e+2	4.9e-2	1.3e-1	4.3e-2	1.1e-1	3.4e-1	1.1e+0	7.5e+0	1.0e+1
DPA-1	2.5e+2	2.7e+2	7.7e-2	1.3e-1	8.1e-2	9.9e-2	7.7e-1	6.0e-2	8.2e+0	7.5e+0
SchNet	4.1e+2	5.0e+2	5.1e-2	2.4e-1	4.5e-2	1.5e-1	8.0e-1	3.8e+0	7.2e+1	3.8e+2
DimeNet++	5.2e+1	7.7e+1	2.9e-2	6.7e-2	2.2e-2	4.3e-2	9.5e-2	2.2e+0	9.2e-1	3.6e+0
GemNet-T	3.3e+1	5.4e+1	2.9e-2	1.2e-1	2.2e-2	8.5e-2	1.6e-2	9.2e-1	3.0e-1	7.0e+0
GemNet-dT	4.2e+1	6.7e+1	2.9e-2	6.7e-2	2.1e-2	8.5e-2	6.8e-2	1.6e+0	6.5e-1	2.0e+1
NequIP	8.3e+1	1.2e+2	3.3e-2	1.9e-1	3.1e-2	1.5e-1	1.4e-1	9.0e-1	1.9e+0	3.5e+0
Allegro	1.6e+2	3.6e+2	4.7e-2	N/A	4.7e-2	N/A	3.0e-1	1.7e-1	2.6e+0	8.3e+0
MACE	6.6e+1	9.2e+1	3.2e-2	N/A	2.6e-2	N/A	5.9e-2	1.6e-1	5.8e-1	8.7e+0
SCN	4.2e+1	6.7e+1	2.8e-2	6.4e-2	2.1e-2	5.6e-2	4.6e-1	8.6e-1	5.3e+0	5.2e+0

Table D.7: Raw scores of models for indicators using HfO dataset and trained with MSE-based loss. Lower values indicate better performance. N/A indicates an interrupted simulation for the dynamic indicator due to abnormal energy changes.

	EF mtr. <sup>ID</sup>	EF mtr. <sup>OOD</sup>	RDF <sup>ID</sup>	RDF <sup>OOD</sup>	ADF <sup>ID</sup>	ADF <sup>OOD</sup>	V <sub>0</sub> <sup>ID</sup>	V <sub>0</sub> <sup>OOD</sup>	B <sub>0</sub> <sup>ID</sup>	B <sub>0</sub> <sup>OOD</sup>
BPNN	2.2e+2	2.3e+2	4.5e-2	9.1e-2	4.5e-2	7.6e-2	6.3e-1	1.2e+0	3.6e+0	6.3e+0
DPA-1	2.3e+2	2.6e+2	7.1e-2	1.3e-1	7.4e-2	8.9e-2	6.3e-1	6.4e-1	5.9e+0	2.8e+0
SchNet	4.2e+2	5.1e+2	9.0e-2	4.1e-1	9.2e-2	1.5e-1	1.1e+0	6.6e+0	1.2e+2	7.7e+2
DimeNet++	4.8e+1	7.0e+1	3.2e-2	1.3e-1	2.5e-2	1.4e-1	6.3e-2	1.9e+0	1.0e+0	3.9e+1
GemNet-T	4.6e+1	6.3e+1	3.4e-2	6.5e-2	2.6e-2	5.1e-2	4.4e-2	9.5e-1	2.9e-1	1.5e+0
GemNet-dT	5.7e+1	7.5e+1	3.2e-2	1.6e-1	2.5e-2	1.3e-1	1.7e-1	1.1e+0	1.5e+0	1.3e+1
NequIP	8.0e+1	1.0e+2	3.8e-2	1.4e-1	3.5e-2	8.5e-2	5.8e-2	2.7e-1	1.1e+0	8.1e-1
Allegro	1.6e+2	2.3e+2	4.8e-2	N/A	4.7e-2	N/A	3.3e-1	1.9e-1	9.4e-1	1.4e+1
MACE	6.5e+1	1.4e+2	3.3e-2	1.0e-1	2.6e-2	7.0e-2	3.1e-2	5.4e-1	3.7e-1	3.1e+0
SCN	5.7e+1	8.1e+1	2.9e-2	8.0e-2	2.2e-2	7.1e-2	3.2e-1	1.5e+0	2.2e+0	3.7e+0

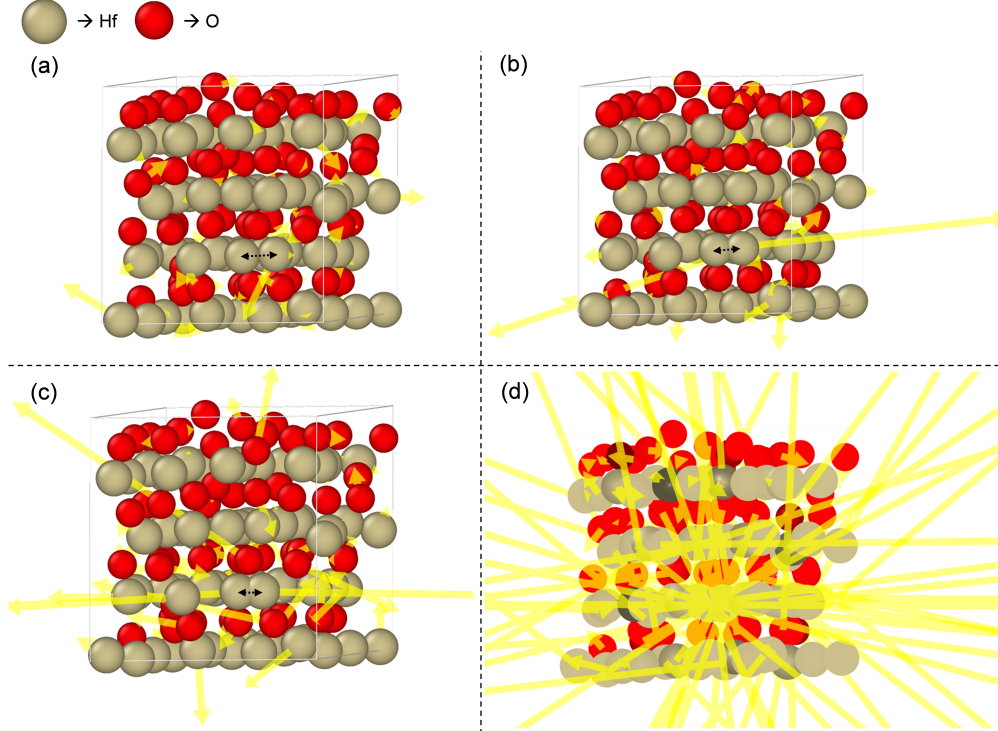


Figure D.8: Hf collision during the evaluation of Allegro in dynamic indicators. (a), (b), (c), and (d) are presented in chronological order, with each atom's relative force magnitude and direction indicated by a yellow arrow.

Table D.8: The maximum threshold ( $TH_{\max}$ ) used for mapping a metric  $x$ . The transformation rule is as follows:  $x' = \frac{TH_{\max} - x}{TH_{\max} - x_{\min}}$ , where  $x'$  represents the transformed score,  $x$  is the original score,  $TH_{\max}$  is the maximum threshold, and  $x_{\min}$  represents the minimum value among the model evaluation results, where achieving the minimum signifies the most desirable outcome for each metric.

Metric	$x_{\min}$		$TH_{\max}$
	HfO	SiN	
EF Metric	3.3e+1	9.8e+1	8.5e+2
RDF	2.8e-2	3.3e-2	4.5e-1
ADF	2.1e-2	2.1e-2	4.5e-1
Bulk Modulus ( $B_0$ )	1.9e-1	3.9e-1	5.0e+1
Equilibrium Vol. ( $V_0$ )	1.6e-2	2.5e-2	3.0e+0

## 968 D.5 Empirical Model Analysis

969 **BPNN** requires hand-crafted features but has simplest overall structure, which makes the model  
 970 relatively fast and less accurate than recent SOTA models.

971 **DPA-1** has relatively low accuracy on the EF metric, but shows better performance on simulation  
 972 indicators.

973 **SchNet** lies on the pareto-frontier as the fastest model, but the speed gain may not be enough to  
 974 compensate overall accuracy drop compared to recent fast models such as Allegro or GemNet-dT.

975 **DimeNet++** and **GemNet-T** are models with similar base structure. They show a similar overall  
 976 tendency: generally high accuracy but less accurate on the  $V_0$  indicator, and having relatively slow  
 977 inference speed. Overall, GemNet-T, which lies on the most accurate side of pareto-frontier, seems  
 978 slightly better predictive performance than DimeNet++.

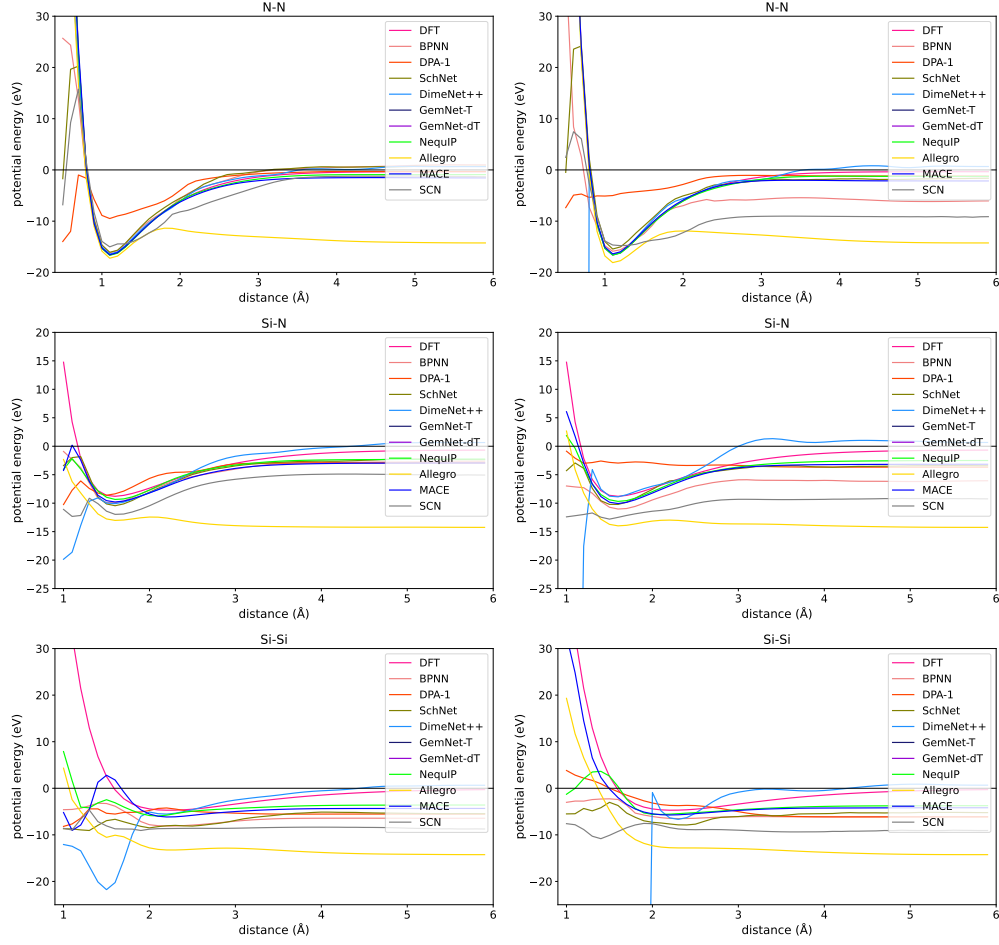


Figure D.9: Comparison of two-body PECs for SiN using MLFF models: models trained with MAE-based loss (left) and MSE-based loss (right).

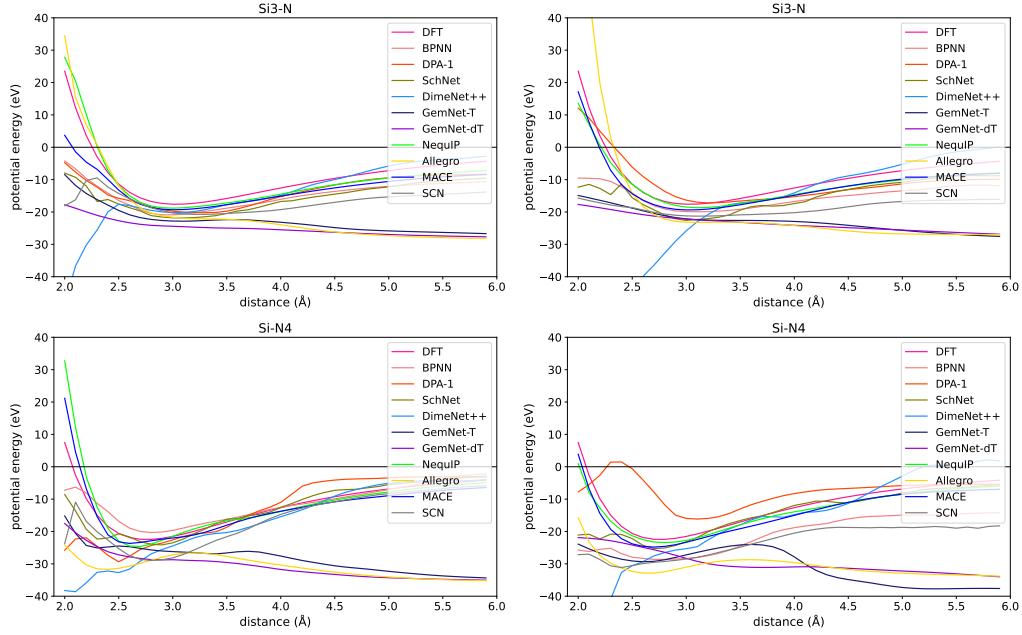


Figure D.10: Comparison of many-body PECs for SiN using MLFF models: models trained with MAE-based loss (left) and MSE-based loss (right).

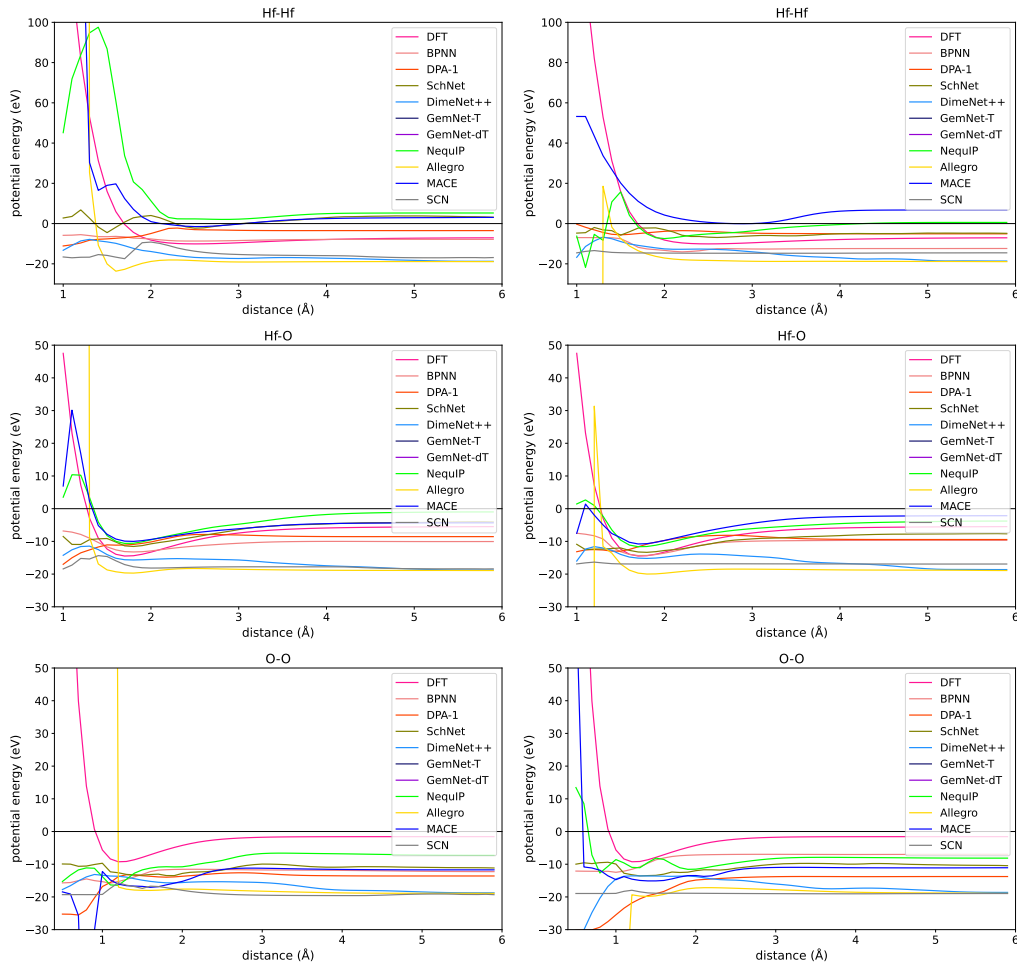


Figure D.11: Comparison of two-body PECs for HfO using MLFF models: models trained with MAE-based loss (left) and MSE-based loss (right).

979 **GemNet-dT** employs direct force prediction, which does not require backpropagation to compute  
 980 forces from energy. As a result, it provides a high inference speed with low accuracy drop on ID  
 981 samples, but not on OOD samples.

982 **NequIP** and **MACE** show high accuracy on both EF metrics and simulation indicators. However,  
 983 MACE is 1.7x faster than NequIP with similar accuracy. Compared to GemNet-dT, MACE has  
 984 equivalent EF accuracy similar prediction results on energy and force with a slightly slower inference  
 985 speed, but it works better on simulation indicators, and on OOD samples.

986 **Allegro** has a high inference speed and decent EF accuracy, however, it show perform unstable on  
 987 MD simulations with OOD samples.

988 **SCN** achieves high overall accuracy, but has a significantly slow inference speed which is 5x slower  
 989 than GemNet-T, and 9x slower than MACE.

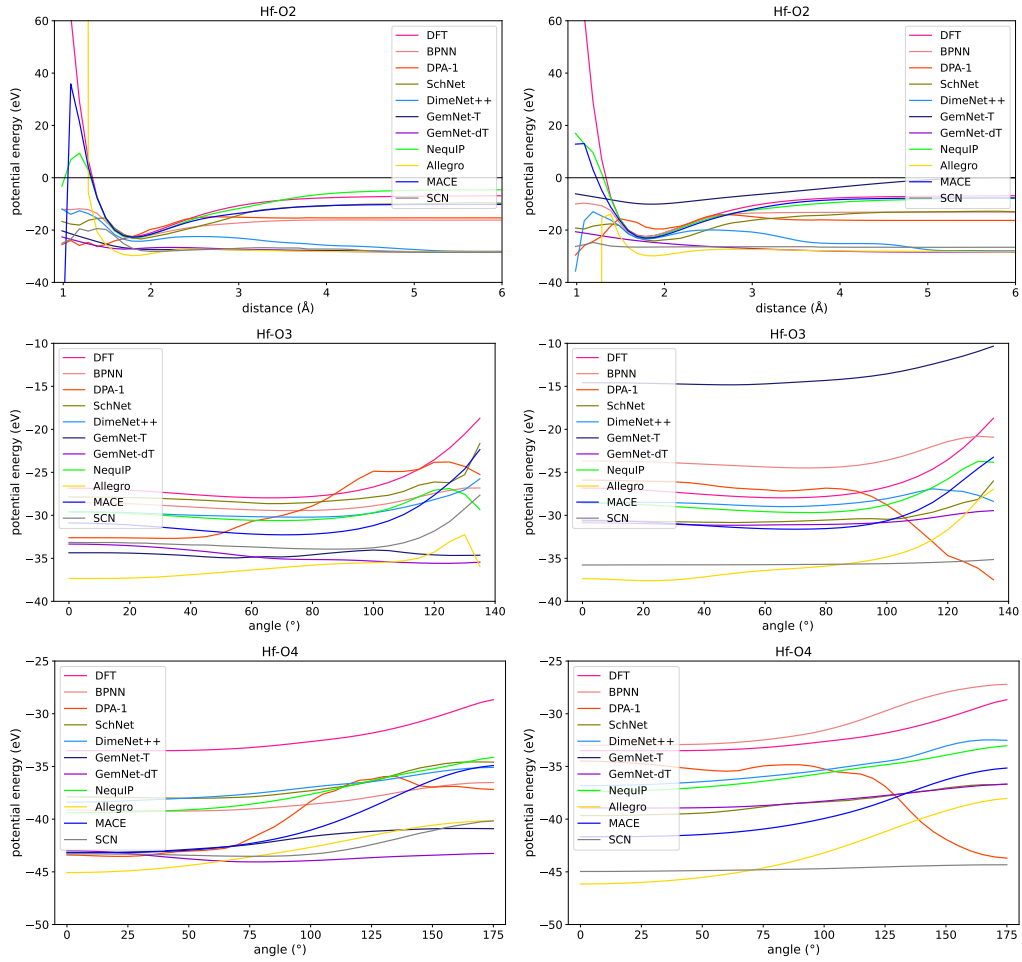
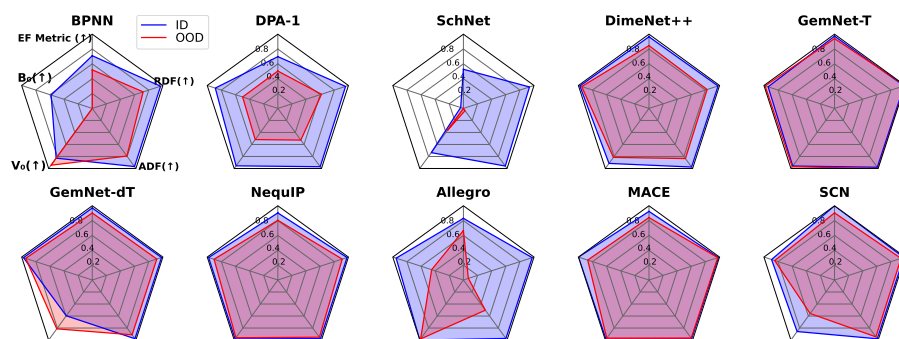
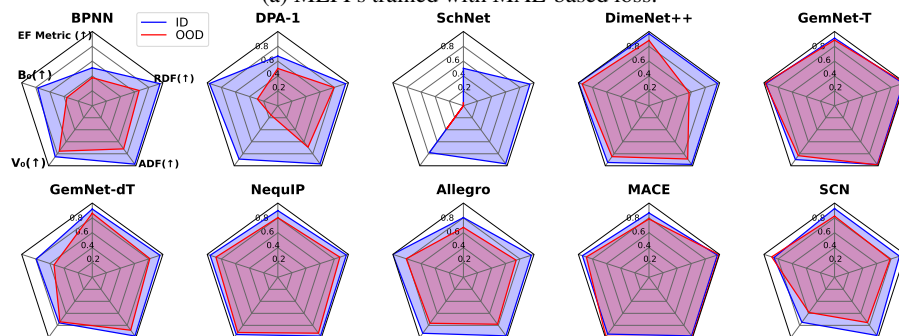


Figure D.12: Comparison of many-body PECs for HfO using MLFF models: models trained with MAE-based loss (left) and MSE-based loss (right).



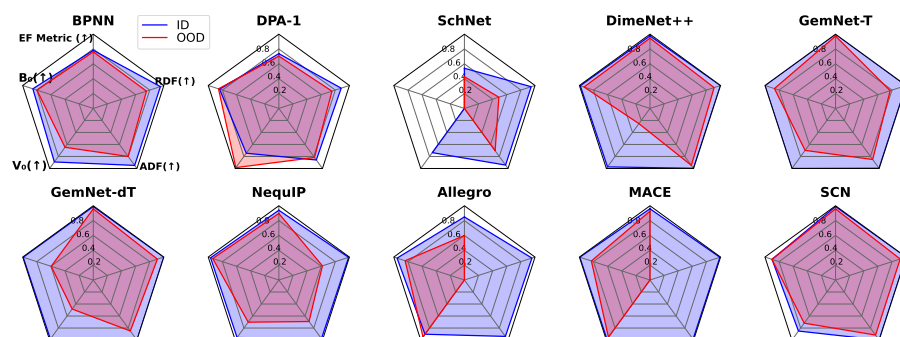


(a) MLFFs trained with MAE-based loss.

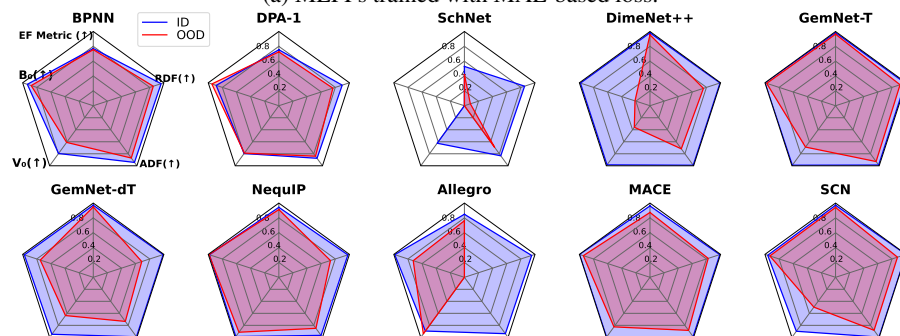


(b) MLFFs trained with MSE-based loss.

Figure D.13: Comprehensive comparison of models on the SiN dataset, based on EF metric and simulation metrics. **Higher values indicate better performance.** The red and blue plots represent the results for ID and OOD, respectively.



(a) MLFFs trained with MAE-based loss.



(b) MLFFs trained with MSE-based loss.

Figure D.14: Comprehensive comparison of models on the HfO dataset, based on EF metric and simulation metrics. **Higher values indicate better performance.** The red and blue plots represent the results for ID and OOD, respectively.

## References

- [1] Seung Soo Kim, Soo Kyeom Yong, Whayoung Kim, Sukin Kang, Hyeon Woo Park, Kyung Jean Yoon, Dong Sun Sheen, Seho Lee, and Cheol Seong Hwang. Review of semiconductor flash memory devices for material and process issues. *Advanced Materials*, 2022.
- [2] Roberto Bez, Paolo Fantini, and Agostino Pirovano. Historical review of semiconductor memories. In *Semiconductor Memories and Systems*, pages 1–26. Elsevier, 2022.
- [3] Ndubuisi G Orji, Mustafa Badaroglu, Bryan M Barnes, Carlos Beitia, Benjamin D Bunday, Umberto Celano, Regis J Kline, Mark Neisser, Yaw Obeng, and AE Vladar. Metrology for the next generation of semiconductor devices. *Nature Electronics*, 1(10):532–547, 2018.
- [4] Koji Nakamae. Electron microscopy in semiconductor inspection. *Measurement Science and Technology*, 32(5):052003, 2021.
- [5] Dylan M Anstine and Olexandr Isayev. Machine learning interatomic potentials and long-range physics. *The Journal of Physical Chemistry A*, 127(11):2417–2431, 2023.
- [6] Kyuhyun Lee, Dongsun Yoo, Wonseok Jeong, and Seungwu Han. Simple-nn: An efficient package for training and executing neural-network interatomic potentials. *Computer Physics Communications*, 242:95–103, 2019.
- [7] Emir Kocer, Tsz Wai Ko, and Jörg Behler. Neural network potentials: A concise overview of methods. *Annual Review of Physical Chemistry*, 73:163–186, 2022.
- [8] Sina Stocker, Johannes Gasteiger, Florian Becker, Stephan Günnemann, and Johannes T Margraf. How robust are modern graph neural network potentials in long and hot molecular dynamics simulations? *Machine Learning: Science and Technology*, 3(4):045010, 2022.
- [9] Xiang Fu, Zhenghao Wu, Wujie Wang, Tian Xie, Sinan Ketten, Rafael Gomez-Bombarelli, and Tommi S. Jaakkola. Forces are not enough: Benchmark and critical evaluation for machine learning force fields with molecular simulations. *Transactions on Machine Learning Research*, 2023.
- [10] Aidan P Thompson, Laura P Swiler, Christian R Trott, Stephen M Foiles, and Garritt J Tucker. Spectral neighbor analysis method for automated generation of quantum-accurate interatomic potentials. *Journal of Computational Physics*, 285:316–330, 2015.
- [11] Albert P Bartók, Mike C Payne, Risi Kondor, and Gábor Csányi. Gaussian approximation potentials: The accuracy of quantum mechanics, without the electrons. *Physical Review Letters*, 104(13):136403, 2010.
- [12] Albert P Bartók, Risi Kondor, and Gábor Csányi. On representing chemical environments. *Physical Review B*, 87(18):184115, 2013.
- [13] Jörg Behler and Michele Parrinello. Generalized neural-network representation of high-dimensional potential-energy surfaces. *Physical Review Letters*, 98(14):146401, 2007.
- [14] Yaolong Zhang, Ce Hu, and Bin Jiang. Embedded atom neural network potentials: Efficient and accurate machine learning with a physically inspired representation. *The Journal of Physical Chemistry Letters*, 10(17):4962–4967, 2019.
- [15] Yaolong Zhang, Junfan Xia, and Bin Jiang. Reann: A pytorch-based end-to-end multi-functional deep neural network package for molecular, reactive, and periodic systems. *The Journal of Chemical Physics*, 156(11):114801, 2022.
- [16] Kristof Schütt, Pieter-Jan Kindermans, Huziel Enoc Saucedo Felix, Stefan Chmiela, Alexandre Tkatchenko, and Klaus-Robert Müller. Schnet: A continuous-filter convolutional neural network for modeling quantum interactions. *Advances in Neural Information Processing Systems*, 30, 2017.

- [17] Johannes Gasteiger, Janek Groß, and Stephan Günnemann. Directional message passing for molecular graphs. *International Conference on Learning Representations*, 2020.
- [18] Johannes Gasteiger, Shankari Giri, Johannes T Margraf, and Stephan Günnemann. Fast and uncertainty-aware directional message passing for non-equilibrium molecules. *Machine Learning for Molecules Workshop, Advances in Neural Information Processing Systems*, 2020.
- [19] Johannes Gasteiger, Florian Becker, and Stephan Günnemann. Gemnet: Universal directional graph neural networks for molecules. *Advances in Neural Information Processing Systems*, 34, 2021.
- [20] Chengxuan Ying, Tianle Cai, Shengjie Luo, Shuxin Zheng, Guolin Ke, Di He, Yanming Shen, and Tie-Yan Liu. Do transformers really perform badly for graph representation? *Advances in Neural Information Processing Systems*, 34, 2021.
- [21] Kristof Schütt, Oliver Unke, and Michael Gastegger. Equivariant message passing for the prediction of tensorial properties and molecular spectra. *International Conference on Machine Learning*, 2021.
- [22] Simon Batzner, Albert Musaelian, Lixin Sun, Mario Geiger, Jonathan P Mailoa, Mordechai Kornbluth, Nicola Molinari, Tess E Smidt, and Boris Kozinsky. E(3)-equivariant graph neural networks for data-efficient and accurate interatomic potentials. *Nature Communications*, 13(1):2453, 2022.
- [23] Albert Musaelian, Simon Batzner, Anders Johansson, Lixin Sun, Cameron J Owen, Mordechai Kornbluth, and Boris Kozinsky. Learning local equivariant representations for large-scale atomistic dynamics. *Nature Communications*, 14(1):579, 2023.
- [24] Ilyes Batatia, David Peter Kovacs, Gregor N. C. Simm, Christoph Ortner, and Gabor Csanyi. Mace: Higher order equivariant message passing neural networks for fast and accurate force fields. *Advances in Neural Information Processing Systems*, 35, 2022.
- [25] Larry Zitnick, Abhishek Das, Adeesh Kolluru, Janice Lan, Muhammed Shuaibi, Anuroop Sriram, Zachary Ulissi, and Brandon Wood. Spherical channels for modeling atomic interactions. *Advances in Neural Information Processing Systems*, 35, 2022.
- [26] Justin S Smith, Olexandr Isayev, and Adrian E Roitberg. Ani-1, a data set of 20 million calculated off-equilibrium conformations for organic molecules. *Scientific Data*, 4(1):1–8, 2017.
- [27] Anders S Christensen and O Anatole Von Lilienfeld. On the role of gradients for machine learning of molecular energies and forces. *Machine Learning: Science and Technology*, 1(4):045018, 2020.
- [28] Dávid Péter Kovács, Cas van der Oord, Jiri Kucera, Alice EA Allen, Daniel J Cole, Christoph Ortner, and Gábor Csányi. Linear atomic cluster expansion force fields for organic molecules: beyond rmse. *Journal of Chemical Theory and Computation*, 17(12):7696–7711, 2021.
- [29] Lowik Chanussot, Abhishek Das, Siddharth Goyal, Thibaut Lavril, Muhammed Shuaibi, Morgane Riviere, Kevin Tran, Javier Heras-Domingo, Caleb Ho, Weihua Hu, et al. Open catalyst 2020 (oc20) dataset and community challenges. *ACS Catalysis*, 11(10):6059–6072, 2021.
- [30] Richard Tran, Janice Lan, Muhammed Shuaibi, Brandon M Wood, Siddharth Goyal, Abhishek Das, Javier Heras-Domingo, Adeesh Kolluru, Ammar Rizvi, Nima Shoghi, et al. The open catalyst 2022 (oc22) dataset and challenges for oxide electrocatalysts. *ACS Catalysis*, 13(5):3066–3084, 2023.

- [31] Yunxing Zuo, Chi Chen, Xiangguo Li, Zhi Deng, Yiming Chen, Jörg Behler, Gábor Csányi, Alexander V. Shapeev, Aidan P. Thompson, Mitchell A. Wood, and Shyue Ping Ong. Performance and cost assessment of machine learning interatomic potentials. *The Journal of Physical Chemistry A*, 124(4):731–745, 2020.
- [32] Volker L Deringer, Miguel A Caro, and Gábor Csányi. A general-purpose machine-learning force field for bulk and nanostructured phosphorus. *Nature Communications*, 11(1):5461, 2020.
- [33] Ganesh Sivaraman, Anand Narayanan Krishnamoorthy, Matthias Baur, Christian Holm, Marius Stan, Gábor Csányi, Chris Benmore, and Álvaro Vázquez-Mayagoitia. Machine-learned interatomic potentials by active learning: amorphous and liquid hafnium dioxide. *npj Computational Materials*, 6(1):104, 2020.
- [34] Wen Zhao, Hu Qiu, and Wanlin Guo. A deep neural network potential for water confined in graphene nanocapillaries. *The Journal of Physical Chemistry C*, 126(25):10546–10553, 2022.
- [35] Venkatesh Botu, Rohit Batra, James Chapman, and Rampi Ramprasad. Machine learning force fields: construction, validation, and outlook. *The Journal of Physical Chemistry C*, 121(1):511–522, 2017.
- [36] Vladimir A Gritsenko. Electronic structure of silicon nitride. *Physics-Uspekhi*, 55(5):498, 2012.
- [37] Zhenhai Li, Tianyu Wang, Jiajie Yu, Jialin Meng, Yongkai Liu, Hao Zhu, Qingqing Sun, David Wei Zhang, and Lin Chen. Ferroelectric hafnium oxide films for in-memory computing applications. *Advanced Electronic Materials*, 8:2200951, 2022.
- [38] Georg Kresse and Jürgen Furthmüller. Efficiency of ab-initio total energy calculations for metals and semiconductors using a plane-wave basis set. *Computational Materials Science*, 6(1):15–50, 1996.
- [39] Georg Kresse and Jürgen Furthmüller. Efficient iterative schemes for ab initio total-energy calculations using a plane-wave basis set. *Physical Review B*, 54(16):11169, 1996.
- [40] Georg Kresse and Jürgen Hafner. Ab initio molecular dynamics for liquid metals. *Physical Review B*, 47(1):558, 1993.
- [41] John P Perdew, Kieron Burke, and Matthias Ernzerhof. Generalized gradient approximation made simple. *Physical Review Letters*, 77(18):3865, 1996.
- [42] Anubhav Jain, Shyue Ping Ong, Geoffroy Hautier, Wei Chen, William Davidson Richards, Stephen Dacek, Shreyas Cholia, Dan Gunter, David Skinner, Gerbrand Ceder, et al. Commentary: The materials project: A materials genome approach to accelerating materials innovation. *APL materials*, 1(1), 2013.
- [43] Stefano Curtarolo, Wahyu Setyawan, Gus LW Hart, Michal Jahnatek, Roman V Chepulskii, Richard H Taylor, Shidong Wang, Junkai Xue, Kesong Yang, Ohad Levy, et al. Aflow: An automatic framework for high-throughput materials discovery. *Computational Materials Science*, 58:218–226, 2012.
- [44] Muratahan Aykol, Shyam S Dwaraknath, Wenhao Sun, and Kristin A Persson. Thermodynamic limit for synthesis of metastable inorganic materials. *Science Advances*, 4(4):eaq0148, 2018.
- [45] Sungwoo Kang, Wonseok Jeong, Changho Hong, Seungwoo Hwang, Youngchae Yoon, and Seungwu Han. Accelerated identification of equilibrium structures of multicomponent inorganic crystals using machine learning potentials. *npj Computational Materials*, 8(1):108, 2022.
- [46] Changho Hong, Jeong Min Choi, Wonseok Jeong, Sungwoo Kang, Suyeon Ju, Kyeongpung Lee, Jisu Jung, Yong Youn, and Seungwu Han. Training machine-learning potentials for crystal structure prediction using disordered structures. *Physical Review B*, 102(22):224104, 2020.

- [47] Johannes Gasteiger, Muhammed Shuaibi, Anuroop Sriram, Stephan Günnemann, Zachary Ward Ulissi, C. Lawrence Zitnick, and Abhishek Das. Gemnet-OC: Developing graph neural networks for large and diverse molecular simulation datasets. *Transactions on Machine Learning Research*, 2022.
- [48] Han Wang, Linfeng Zhang, Jiequn Han, and Weinan E. DeePMD-kit: A deep learning package for many-body potential energy representation and molecular dynamics. *Computer Physics Communications*, 228:178–184, 2018.
- [49] Duo Zhang, Hangrui Bi, Fu-Zhi Dai, Wanrun Jiang, Linfeng Zhang, and Han Wang. Dpa-1: Pretraining of attention-based deep potential model for molecular simulation. *arXiv preprint arXiv:2208.08236*, 2022.
- [50] Francis Birch. Finite elastic strain of cubic crystals. *Physical Review*, 71(11):809, 1947.
- [51] Paul J Steinhardt, David R Nelson, and Marco Ronchetti. Bond-orientational order in liquids and glasses. *Physical Review B*, 28(2):784, 1983.
- [52] Wolfgang Lechner and Christoph Dellago. Accurate determination of crystal structures based on averaged local bond order parameters. *The Journal of Chemical Physics*, 129(11):114707, 2008.
- [53] Ask Hjorth Larsen, Jens Jørgen Mortensen, Jakob Blomqvist, Ivano E Castelli, Rune Christensen, Marcin Dułak, Jesper Friis, Michael N Groves, Bjørk Hammer, Cory Hargus, et al. The atomic simulation environment—a python library for working with atoms. *Journal of Physics: Condensed Matter*, 29(27):273002, 2017.
- [54] Shyue Ping Ong, William Davidson Richards, Anubhav Jain, Geoffroy Hautier, Michael Kocher, Shreyas Cholia, Dan Gunter, Vincent L Chevrier, Kristin A Persson, and Gerbrand Ceder. Python materials genomics (pymatgen): A robust, open-source python library for materials analysis. *Computational Materials Science*, 68:314–319, 2013.
- [55] Ilyes Batatia, Simon Batzner, Dávid Péter Kovács, Albert Musaelian, Gregor NC Simm, Ralf Drautz, Christoph Ortner, Boris Kozinsky, and Gábor Csányi. The design space of e(3)-equivariant atom-centered interatomic potentials. *arXiv preprint arXiv:2205.06643*, 2022.
- [56] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [57] Sashank J. Reddi, Satyen Kale, and Sanjiv Kumar. On the convergence of adam and beyond. *International Conference on Learning Representations*, 2018.
- [58] Jinwoo Park, Byung Deok Yu, and Suklyun Hong. Van der waals density functional theory study for bulk solids with bcc, fcc, and diamond structures. *Current Applied Physics*, 15(8):885–891, 2015.
- [59] Hua-Jin Zhai, Wen-Jie Chen, Shu-Juan Lin, Xin Huang, and Lai-Sheng Wang. Monohafnium oxide clusters hfo n—and hfo n (n= 1–6): Oxygen radicals, superoxides, peroxides, diradicals, and triradicals. *The Journal of Physical Chemistry A*, 117(6):1042–1052, 2013.
- [60] Oliver T Unke, Stefan Chmiela, Michael Gastegger, Kristof T Schütt, Huziel E Sauceda, and Klaus-Robert Müller. Spookynet: Learning force fields with electronic degrees of freedom and nonlocal effects. *Nature communications*, 12(1):7273, 2021.
- [61] Raghunathan Ramakrishnan, Pavlo O Dral, Matthias Rupp, and O Anatole Von Lilienfeld. Big data meets quantum chemistry approximations: the  $\delta$ -machine learning approach. *Journal of chemical theory and computation*, 11(5):2087–2096, 2015.
- [62] Mihail Bogojeski, Leslie Vogt-Maranto, Mark E Tuckerman, Klaus-Robert Müller, and Kieron Burke. Quantum chemical accuracy from density functional approximations via machine learning. *Nature communications*, 11(1):5223, 2020.