

Evaluation of Foundational Machine Learned Interatomic Potentials for Migration Barrier Predictions

Achinthya Krishna Bheemaguli^{1,2}, Penghao Xiao^{3,*}, and Gopalakrishnan Sai Gautam^{2,*}

¹Department of Metallurgical and Materials Engineering, National Institute of Technology Karnataka, Surathkal 575025, India

²Department of Materials Engineering, Indian Institute of Science, Bengaluru, 560012, India

³Department of Physics and Atmospheric Science, Dalhousie University, Halifax B3H 4R2, Nova Scotia, Canada

*Email: penghao.xiao@dal.ca; saigautamg@iisc.ac.in

Abstract

Fast, and accurate prediction of ionic migration barriers (E_m) is crucial for designing next-generation battery materials that combine high energy density with facile ion transport. Given the computational costs associated with estimating E_m using conventional density functional theory (DFT) based nudged elastic band (NEB) calculations, we benchmark the accuracy in E_m and geometry predictions of five foundational machine learned interatomic potentials (MLIPs), which can potentially accelerate predictions of ionic transport. Specifically, we assess the accuracy of MACE-MP-0, Orb-v3, SevenNet, CHGNet, and M3GNet models, coupled with the NEB framework, against DFT-NEB-calculated E_m across a diverse set of battery-relevant chemistries and structures. Notably, MACE-MP-0 and Orb-v3 exhibit the lowest mean absolute errors in E_m predictions across the entire dataset and over data points that are not outliers, respectively. Importantly, Orb-v3 and SevenNet classify ‘good’ versus ‘bad’ ionic conductors with an accuracy of >82%, based on a threshold E_m of 500 meV, indicating their utility in high-throughput screening approaches. Notably, intermediate images generated by MACE-MP-0 and SevenNet provide better initial guesses relative to conventional interpolation techniques in >71% of structures, offering a practical route to accelerate subsequent DFT-NEB relaxations. Finally, we observe that accurate E_m predictions by MLIPs are not correlated with accurate (local) geometry predictions. Our work establishes the use-cases, accuracies, and limitations of foundational MLIPs in estimating E_m and should serve as a base for accelerating the discovery of novel ionic conductors for batteries and beyond.

1 Introduction

Developing next-generation batteries is essential for our transition into sustainable energy usage, given that the state-of-the-art lithium-ion batteries (LIBs), while already delivering excellent performance, are approaching their fundamental limits [1, 2], necessitating the discovery of novel beyond-LIB materials and chemistries. One key materials property that governs battery performance is the ionic diffusivity (D) of the

electroactive ion in electrodes and (solid) electrolytes. D is exponentially influenced by the energy barrier each ion must overcome, commonly referred to as the migration barrier (E_m), to hop from its initial lattice site to a symmetrically equivalent final site [3–5]. Each ionic hop is often mediated through vacancies in the lattice with the ion overcoming transition state(s) along the hop. D and E_m are thus related via the Arrhenius expression, $D = D_0 \exp(-E_m/k_B T)$ [6].

Accordingly, materials with a low E_m – in both electrodes and (solid) electrolytes exhibit higher ionic conductivity and enable faster charge/discharge rates [7]. In particular, emerging multivalent battery chemistries, such as Mg or Ca-based systems that promise higher volumetric energy densities [8], often suffer from poor rate performance [9–14]. Some Na-ion cathodes such as maricite-Na(Mn/Fe)PO₄, phosphate alluaudite-Na_xMnFe₂(PO₄)₃ and sulfate sodium superionic conductors (NaSICONs) that offer lower costs compared to LIB cathodes also suffer from poor rate performance [15,16]. Therefore, understanding and minimizing the E_m in candidate materials is crucial for advancing the next-generation of high-performance batteries.

Experimental techniques such as quasi-elastic neutron scattering [17], electrochemical impedance spectroscopy [18], nuclear magnetic resonance measurements [19], and galvanostatic intermittent titration techniques [20] are commonly employed to study ion dynamics in solids [21]. However, these methods often require access to large-scale facilities, and can exhibit chemistry or material-specific constraints/requirements, limiting their accessibility. As a result, computational approaches, particularly density functional theory (DFT [22,23])-based nudged elastic band (NEB [24]) calculations are commonly used for computationally estimating E_m with reasonable precision [25]. While ab initio molecular dynamics (AIMD) simulations can also be used for estimating E_m , such simulations are computationally expensive since they require the sampling over large length and long time scales across different temperatures to provide reasonable E_m [26,27]. Moreover, AIMD simulations can be unreliable for systems exhibiting high E_m (i.e., the ‘true negatives’ among materials that can conduct ions) due to insufficient sampling of ion dynamics, resulting in DFT-NEB being the usual technique deployed for E_m predictions.

Calculating E_m using DFT-NEB requires an initial guess for the minimum energy path (MEP), which is typically constructed by linearly interpolating the coordinates of the initial and final configurations of the moving ion. Each ‘image’ generated by linear interpolation is subsequently connected via an auxiliary spring force. Note that the initial interpolated guess is often far from the true MEP, increasing the computational expense of DFT-NEB calculations and making them prone to convergence difficulties [25]. Alternative approaches such as ‘ApproxNEB’ [28], have been proposed to reduce computational intensity, with limited efficacy.

Recently, foundational machine learned interatomic potentials (MLIPs [29–31]), also referred to as universal potentials, have emerged as a new paradigm in computational materials science. The foundational MLIPs, pre-trained on large and diverse datasets, can generalize to a wide range of downstream tasks [31] and are transferable across different materials and property prediction tasks [32–35], unlike classical MLIPs or force-fields that are constrained by a specific chemistry or property. Thus, foundational MLIPs are attractive candidates for accelerating atomistic simulations, including NEB calculations, by potentially improving initial MEP guesses and reducing the need for extensive DFT-based refinement or optimization, which can enable high-throughput screening of materials based on their E_m . Indeed, a recent work by Kang et al. [36] proposed an alternate to traditional DFT-based NEB calculations for E_m estimations by using an MLIP for generating the potential energy surface on a spatial grid and extracting the MEP without the need for pre-defined NEB-images.

Several studies have benchmarked the performance of foundational MLIPs on diverse material properties [37–40], but not on predicting E_m in solids. For instance, the ‘Matbench discovery’ platform [41] provides a standardized framework for ranking universal potentials, but does not yet evaluate their integration with NEB workflows for E_m predictions. Other MLIP benchmarking studies include the work by Zhao et al. [42] that evaluated the MLIPs on transition state search for chemical reactions involving molecules. Bihani et al. [43] benchmarked the performance of equivariant MLIPs on their generalizability to higher temperature simulations and unseen compositions, while Mannan et al. evaluated the performance of universal potentials against experimental measurements of elastic properties and structural accuracy among minerals [39]. So far, there has been no benchmark of the performance of state-of-the-art universal potentials in predicting E_m across a wide range of (battery) chemistries and materials, especially by integrating them with NEB workflows.

Here, we assess the performance of foundational MLIPs, namely, MACE-MP-0 [44, 45], SevenNet [46, 47], Orb-v3 [48, 49], CHGNet [50], and M3GNet [51] in predicting E_m with NEB calculations. Using the dataset DFT-calculated E_m compiled and curated by Devi et. al. [52, 53] that spans a wide range of materials and compositions, we benchmark the E_m predictions of the foundational MLIPs against conventional DFT-NEB values at the generalized gradient approximation (GGA [54]) level of exchange-correlation accuracy for 574 migration paths. Additionally, we introduce a metric to assess the similarity of MLIP-NEB relaxed structures with the ground truth of DFT-NEB computed MEPs from our previous works [25, 55–57]. Finally, we examine the correlation between accuracies in E_m and geometry predictions by the MLIPs considered.

Notably, we find that M3GNet and CHGNet tend to underestimate E_m and exhibit a high degree of confidence in predicting low E_m over a narrow range of possible E_m values, while the other potentials exhibit no clear bias and deliver consistent accuracy over a wide range of E_m values. Importantly, we observe that Orb-v3 and SevenNet classify systems as ‘good’ ($E_m < 500$ meV) or ‘bad’ ionic conductors with $> 82\%$ accuracy. Performing an MLIP-NEB, using any of the potential considered, does result in improved interpolated paths representing the MEP in over 66% of cases, indicating their utility in high-throughput screening workflows. Significantly, we find no evident correlation between the accuracy of E_m and geometry predictions, with MLIPs yielding higher accuracy in E_m predictions for systems with low E_m values, while demonstrating better geometry predictions in systems with large E_m . We hope that our study establishes use-cases and quantifies the reliability of using foundational MLIPs in predicting E_m over a diverse set of chemistries and crystal structures, which in turn should accelerate materials discovery for novel battery applications and beyond.

2 Methods

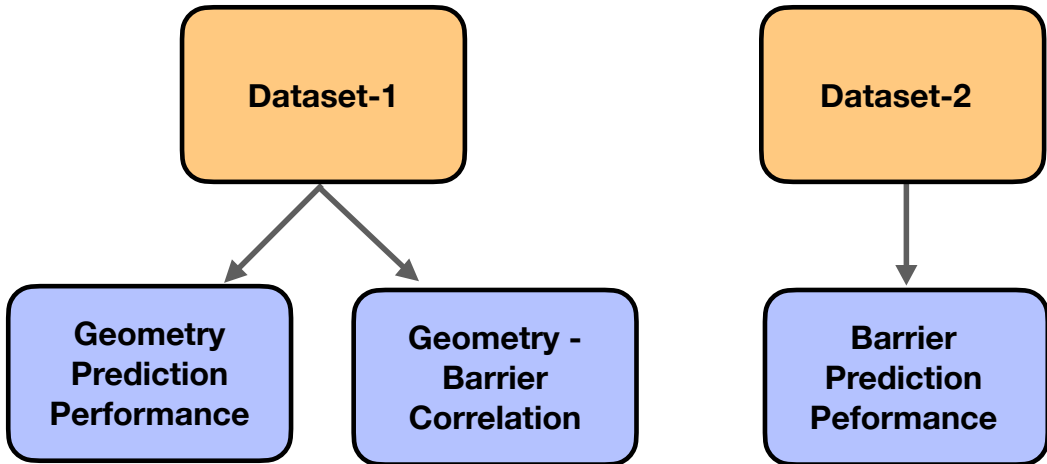


Figure 1: Overview of the methodology, indicating the use of two subsets of the E_m dataset that were created for examining geometry predictions, geometry-barrier correlations, and barrier predictions.

2.1 Datasets

We utilize two distinct subsets of the E_m dataset in this work. The smaller dataset, referred to as ‘Dataset-1’, comprises 60 DFT-NEB calculations, including multiple possible migration pathways for select structures. Each datapoint consists of GGA-calculated E_m and the relaxed structures of all intermediate images across a migration pathway. We constructed this dataset using the DFT-NEB results from our previous works [25, 55–57], encompassing crystal structures explored primarily as battery materials, such as layered structures, Weberites, spinels, olivines, perovskites, and NaSICONs. The E_m in Dataset-1 ranges from 0.06 eV to 2.88 eV.

The larger dataset, referred to as ‘Dataset-2’, is a subset of a literature-derived collection of E_m [52], which comprises of 621 DFT-calculated E_m and the initial and final configurations for each migration pathway. Among the 621 datapoints, we excluded systems exhibiting $E_m > 2.5$ eV, since such high E_m values would not correspond to any tangible rate performance under battery operating conditions. We also excluded systems that presented significant convergence difficulties during NEB calculations using any of the foundational MLIPs considered (~ 10 datapoints), so that a fair and quantitative comparison can be made across the MLIPs. Thus, the final subset that forms our Dataset-2 consists of 574 systems. The systems comprising both datasets are compiled in our GitHub repository, while Dataset-2 is also available as a json file on Zenodo.

2.2 Model Details

We used publicly available universal potentials that have demonstrated good performance for bulk properties, are constructed on graph-based neural network (GNN) architectures, and are compatible with the atomic simulation environment (ASE [58]). We leveraged ASE calculators to integrate the MLIPs with the NEB implementation available within ASE. The specific MLIPs we used include *i)* the MACE-MP-0 ‘large’

Table 1: Summary of the MLIPs used.

Model	Training data	Model type and key features
MACE-MP-0	MPtrj dataset	E(3)-equivariant GNN that captures many-body interactions
SevenNet-MF-ompa	MPtrj, OMat24, and sAlex	Equivariant GNN incorporating multifidelity learning with efficient parallelization
Orb-v3	MPtrj, OMat24, and Alex	Roto-equivariance inducing regularized GNN with analytical energy gradients (conservative forces) and (effectively) infinite neighbors
CHGNet	MPtrj dataset	GNN including magnetic moment inputs, thus incorporating information on atomic charges
M3GNet	MPtrj dataset	Includes three-body interactions within its GNN

foundation model, trained on approximately 1.6 million materials project [59] bulk-crystal relaxation trajectories (i.e., the ‘MPtrj’ dataset [60]) with maximal message equivariance ($L=2$), *ii*) the SevenNet-MF-ompa model, which incorporates multifidelity learning with a core architecture based on the neural equivariant interatomic potential (NequIP [61]) and is trained on MPtrj, ‘OMat24’ trajectories [62], and the subsampled Alexandria (sAlex [63]) datasets, *iii*) the Orb-v3-conservative-inf-omat that is trained on the MPtrj, OMat24 and Alexandria (Alex) dataset with learned forces being conservative by construction and effectively unlimited neighbor lists, *iv*) CHGNet v0.3.0, which is trained on MPtrj, and *v*) M3GNet MP-2021.2.8-EFS version, which is trained on MPtrj up to February 8, 2021. We used all models as provided, without any additional fine-tuning or hyperparameter optimizations. A summary of the specific details for each model is provided in Table 1.

2.3 NEB Calculations

Typically, DFT-NEB calculations employ linear interpolation (LI) of atomic coordinates between the initial and final endpoints for generating the initial guess for the images. Whereas, the image dependent pair potential (IDPP) interpolation technique, developed by Smidstrup and coworkers [64] utilizes a distance-matching objective to generate the initial guess for MEP. Given a preliminary benchmark of MACE-MP-0 based NEB calculations initialized with LI and IDPP interpolation, as detailed in Section S1 and Figure S1 of the supporting information (SI), we find IDPP interpolation to provide marginally better initial guesses for the eventual NEB calculations.

For NEB calculations of materials in Dataset-1 using all universal potentials, we generated seven intermediate images, mirroring the number used in the corresponding DFT-NEB calculations. The initial interpolated images were connected by a spring constant, $k = 5 \text{ eV}/\text{\AA}^2$, and we utilized the NEB implementation following the elastic band (EB, [65]) method with full spring force, given our benchmarking with MACE-MP-0 (see Section S1). We did not include the climbing image technique [24] in any of our MLIP-NEBs, as we did not see significantly different results with or without climbing image in our previous work [25]. We deemed the NEB converged when the band forces fell below $0.05 \text{ eV}/\text{\AA}$, while using the Broyden–Fletcher–Goldfarb–Shanno optimizer [66–69]. In the case of Dataset-2, we employed only three intermediate images for all foundational MLIPs considered, to reduce computational costs. Also, we used the set of optimized NEB parameters that were used for calculations on Dataset-1 (i.e., $k = 5 \text{ eV}/\text{\AA}^2$, IDPP interpolation, and the EB method) for all calculations involving Dataset-2.

2.4 Geometry Metrics

To quantitatively assess the similarity of local geometries between structures, we introduce the geometric similarity classification metric, θ for a given image structure between the endpoints and an averaged geometric similarity score, g , across a migration path. θ evaluates whether the local geometry of an intermediate image obtained via a MLIP-NEB calculation is a better approximation of the reference structure (i.e., DFT-NEB relaxed) compared to the image generated using simple LI. Thus, θ is useful to determine whether intermediate images relaxed using MLIP-NEBs provide a superior initial guess for DFT-NEB calculations than typical LI, based on local geometric features. We define θ using the following steps:

1. **Identify Nearest Neighbors and Pairwise Distances:** We identify the six nearest neighbors of the migrating ion using the Voronoi decomposition technique [70], as implemented in the pymatgen package [71]. Subsequently, we calculate all pairwise distances (d , using pymatgen) between the migrating ion (c) and its six neighbors (i, j, k, l, m, n), as well as among the neighbors themselves. The distances are calculated for structures relaxed/generated by DFT-NEB, MLIP-NEB, and LI. Thus, we compute, for any pair $\{x, y\} \subset \{i, j, k, l, m, n, c\}$ where $x \neq y$:

$$d_{xy}^{\text{DFT}}, d_{xy}^{\text{MLIP}}, d_{xy}^{\text{LI}}$$

2. **Calculate Absolute Errors in Pairwise Distances:** We compute the absolute difference between each pairwise distance in the MLIP-NEB relaxed structure and the LI structure with respect to the corresponding value in the DFT-NEB relaxed structure. These differences are stored in two 21-dimensional vectors:

$$\begin{aligned}\Delta d^{\text{MLIP}} &= |d_{xy}^{\text{DFT}} - d_{xy}^{\text{MLIP}}| \\ \Delta d^{\text{LI}} &= |d_{xy}^{\text{DFT}} - d_{xy}^{\text{LI}}|\end{aligned}$$

for all $\{x, y\} \subset \{i, j, k, l, m, n, c\}$, where $x \neq y$.

3. **Calculate Absolute Errors in Solid Angles:** Since two local geometries can have similar pairwise distances but differ in their angular orientations, we also consider the solid angles (Ω , calculated using pymatgen) subtended by each face of the Voronoi polyhedra formed by the six nearest neighbors.

$$\Omega_x^{\text{DFT}}, \Omega_x^{\text{MLIP}}, \Omega_x^{\text{LI}}, \text{ where } x \in \{a, b, c, d, e, f\}$$

In the above notation, two Ω , say $\Omega^{\text{DFT}}, \Omega^{\text{MLIP}}$, having the same x indicates that the polyhedral faces correspond to the same set of neighboring atoms. The absolute differences with the DFT-NEB relaxed structures are then calculated and stores as six-dimensional vectors:

$$\begin{aligned}\Delta \Omega^{\text{MLIP}} &= |\Omega_x^{\text{DFT}} - \Omega_x^{\text{MLIP}}| \\ \Delta \Omega^{\text{LI}} &= |\Omega_x^{\text{DFT}} - \Omega_x^{\text{LI}}|\end{aligned}$$

for all $x \in \{a, b, c, d, e, f\}$.

We expect the local geometry of an MLIP-NEB relaxed structure to be a poorer approximation of the DFT-NEB relaxed structure than the corresponding LI structure, if at least one of the following conditions is met: (i) one of the 21 pairwise distances or 6 solid angles of the MLIP-NEB relaxed structure deviates

significantly more from the DFT-NEB geometry than the corresponding LI structure, or (ii) the average difference in pairwise distances or solid angles of the MLIP-NEB relaxed structure with the DFT-NEB reference is significantly higher compared to LI. To quantify these two conditions, we calculate δ , which represents the maximum value among the differences in the mean and maximum errors of distances and angles between the MLIP-NEB and LI structures:

$$\delta = \max \left(\overline{\Delta d}^{\text{MLIP}} - \overline{\Delta d}^{\text{LI}}, \quad \max(\Delta d^{\text{MLIP}}) - \max(\Delta d^{\text{LI}}), \right. \\ \left. \overline{\Delta \Omega}^{\text{MLIP}} - \overline{\Delta \Omega}^{\text{LI}}, \quad \max(\Delta \Omega^{\text{MLIP}}) - \max(\Delta \Omega^{\text{LI}}) \right)$$

Here, $\overline{\Delta d}$ and $\overline{\Delta \Omega}$ represent the mean of the absolute errors in distances and solid angles, respectively. $\max(\Delta d)$ and $\max(\Delta \Omega)$ represent the maximum absolute errors.

Finally, the metric θ would classify the structure as:

$$\theta = \begin{cases} \text{Good or } 1, & \delta < 0.01 \\ \text{Comparable or } 0, & 0.01 \leq \delta \leq 0.1 \\ \text{Bad or } -1, & \delta > 0.1 \end{cases} \quad (1)$$

Thus, δ quantifies the difference between the deviations of the MLIP-NEB and LI structures with respect to the DFT-NEB reference based on key local geometric features. A smaller (ideally negative) δ value signifies that the MLIP-NEB structure exhibits consistently lower errors, indicating it's a better approximation of the true DFT-NEB pathway. Conversely, larger (more positive) δ suggests that LI performed as well or even better than the MLIP-NEB for at least one of the local geometric attributes. Therefore, we numerically represent the 'good', 'comparable' and 'bad' structure as 1, 0 and -1 with θ . Finally, for a given system containing i intermediate images, we define g as,

$$g = \frac{\sum_i \theta_i}{\sum_i i}$$

In the case where all the i image local geometries are better (worse) predicted by MLIP-NEB compared to LI, g will take the value of 1 (-1).

3 Results

3.1 Barrier Prediction Performance

Figure 2 presents a comparison of E_m predictions on Dataset-2 across different foundational MLIPs (x -axis) with their corresponding DFT-NEB calculated E_m (y -axis). Green circles, yellow squares, pink pluses, blue triangles, and orange rhombuses represent E_m predictions using MACE-MP-0, SevenNet, Orb-v3, CHGNet, and M3GNet, respectively. Individual parity plot for each MLIP considered is compiled in Figures S2-S6 of the SI. Overall, MACE-MP-0 demonstrates the best performance with an MAE of 0.310 eV, while M3GNet records the highest MAE of 0.349 eV. The MAEs of Orb-v3, CHGNet, and SevenNet are 0.336, 0.343, and 0.344 eV, respectively. To provide a numerical context to the MAEs, DFT-NEB calculations typically carry a ~ 60 meV error in their predictions, and a change of 60 meV in E_m at 298 K corresponds to an order-of-magnitude change in D [72]. However, the calculation of MAEs is influenced by extreme outliers

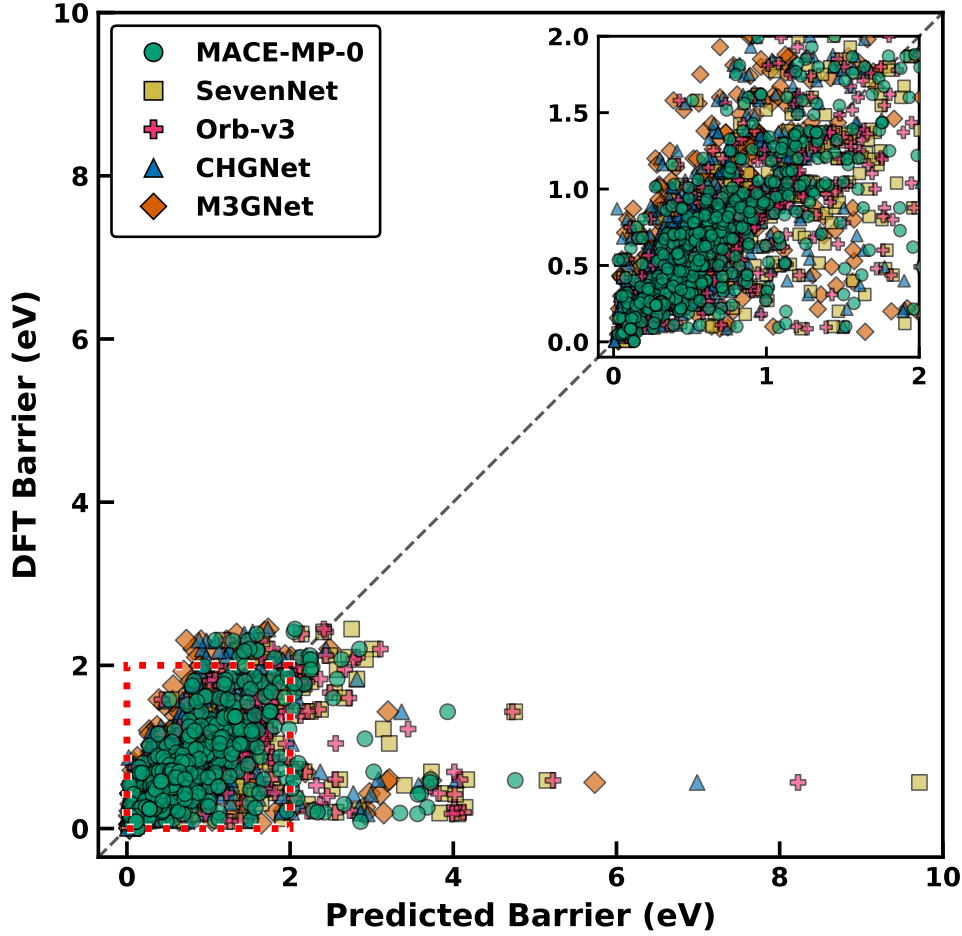


Figure 2: Parity plot of migration barrier predicted by various MLIPs against DFT-NEB, with the dotted black line indicating the parity line. Inset shows the parity plot for a smaller range of DFT/predicted values (0-2 eV).

that affect all MLIPs.

To obtain a more representative picture of MLIP performance, we exclude 17 systems that act as common outliers across all MLIPs, with each outlier exhibiting absolute errors exceeding 1 eV. Notably, excluding the common outliers also reveals a similar performance hierarchy as with retaining the entire dataset: MACE-MP-0 emerges with the best MAE of 0.239 eV, followed closely by Orb-v3 with 0.245 eV. The remaining MLIPs, namely SevenNet, CHGNet, and M3GNet show MAEs of 0.251, 0.275, and 0.290 eV, respectively.

Besides accuracy, we analyze the distribution of datapoints relative to the ideal parity line to determine whether the MLIPs exhibit systematic prediction biases (i.e., under- or over-estimation of E_m). Interestingly, we observe MACE-MP-0, SevenNet, and Orb-v3 to demonstrate a relatively balanced prediction behavior with fairly symmetric distributions of under and over-estimated datapoints. Represented as (number of under-estimated datapoints, number of over-estimated datapoints) pairs, MACE-MP-0, SevenNet, and Orb-

v3 exhibit distributions of (299,275), (244,330), and (242,332), respectively. In contrast, CHGNet and M3GNet show a bias toward under-estimating barriers, with under-estimated datapoints accounting for 73.1% and 78.2% of all predictions, respectively. Represented as (under-estimated, over-estimated) pairs, CHGNet and M3GNet exhibit distributions of (420,154) and (449,125), respectively.

To further understand individual MLIP capabilities, we examined each potential’s performance after excluding the outliers specific to each potential (i.e., systems with absolute errors > 1 eV as predicted by a given potential) to gain insight into the ‘best-case’ scenario of E_m predictions. Notably, despite having 37 outliers, Orb-v3 achieves the lowest MAE of 0.198 eV on its remaining (non-outlier) predictions. With 35 outliers, MACE-MP-0 is a close second with an MAE of 0.202 eV, while SevenNet, with 37 outliers, displays an MAE of 0.203 eV. CHGNet and M3GNet show higher MAE values of 0.248 eV and 0.257 eV with 31 and 36 outliers, respectively. Thus, we find that Orb-v3 can achieve higher accuracies on systems that it describes well while MACE-MP-0 achieves a better balance of both low errors and fewer outliers compared to other MLIPs.

3.2 Predictions Over Different Barrier Ranges

To understand the performance of MLIPs across various ranges of E_m , we divided the 574 datapoints into seven equal-sized distributions based on their DFT-NEB E_m values, as illustrated in Figure 3. While the x -axis in Figure 3 represents DFT-NEB E_m ranges, with bar widths indicating the span of E_m values within each bin, the y -axis shows the percentage of predictions within each bin that achieve absolute errors < 0.1 eV (signifying an “acceptable” degree of accuracy). The exact DFT barrier range of the data points present in each bin can be found in Table S1 of the SI. Trends in Figure 3 indicate that all MLIPs struggle with high E_m predictions, with only small percentage of systems exhibiting the acceptable accuracy in the highest barrier range (~ 1.31 - 2.50 eV). Specifically, the percentage of predictions with acceptable accuracy in the highest E_m range are 20.7% for Orb-v3, 18.3% for MACE-MP-0, 14.6% for SevenNet, 6.1% for M3GNet, and 3.7% for CHGNet.

Importantly, we identify a “sweet spot” of E_m values where all MLIPs perform reasonably well. For example, in the low-barrier range (~ 0.0025 - 0.25 eV), more than 50% of predictions achieve acceptable accuracy across all MLIPs. Within this range, CHGNet shows the highest success rate (59.8%), followed by M3GNet and SevenNet (both 58.5%), while MACE-MP-0 and Orb-v3 achieve 53.7% and 57.3%, respectively. Additionally, Orb-v3 and SevenNet achieve their best performance (i.e., highest fraction of predicted datapoints with acceptable accuracy) in the 0.25-0.36 eV range, achieving 62.2% and 61% acceptable predictions, respectively. MACE-MP-0 performs best in the slightly higher 0.36-0.50 eV range with 57.8% accuracy. Meanwhile, CHGNet and M3GNet perform best in the lowest E_m range (~ 0.0025 - 0.25 eV) with 59.8% and 58.5% accuracy, respectively.

While all MLIPs show declining accuracy with increasing E_m , Orb-v3 exhibits the slowest degradation, maintaining better performance across a broader range of E_m values compared to other potentials. Thus, we find that ‘simpler’ graph models such as CHGNet and M3GNet demonstrate superior performance for materials with intrinsically low E_m values but lack consistency in their predictions over a wider range of E_m . On the other hand, increasing complexity among the graph models, such as in Orb-v3 or MACE-MP-0 allows for a more robust performance across a wide range of E_m values while sacrificing ‘peak’ performance for materials with low E_m , making them better suited for E_m predictions in novel materials. This variation in the performance of ‘simple’ and ‘complex’ MLIPs also reveals the general trade-off between building specialized and generalized models in the field of machine learning.

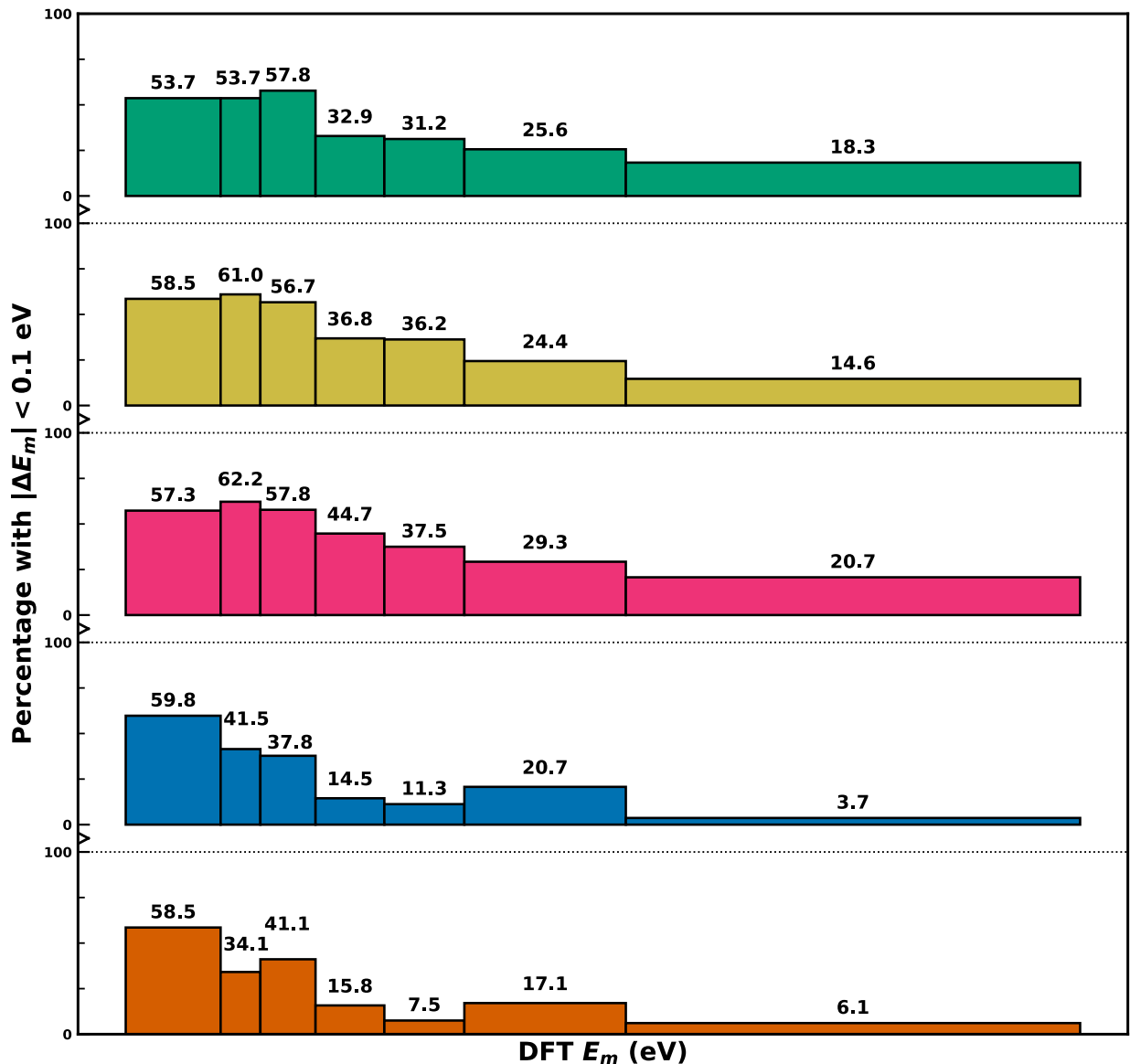


Figure 3: Barrier prediction performance of various MLIPs across different DFT-calculated E_m ranges. The dotted line and kink denote a change in the models, which are, from top to bottom: MACE-MP-0, SevenNet, Orb-v3, CHGNet, and M3GNet. Each bin contains an equal number of data points with the width corresponding to the range of DFT-calculated E_m within the bin. The height of each bin (as indicated by the numerical annotation on each bin) within each model represents the percentage of data points whose E_m values are predicted within an absolute error of 0.1 eV.

3.3 Barrier Classification Performance

To quantify the ability of the MLIPs considered to classify a material as a ‘good’ versus a ‘bad’ ionic conductor, which can be used for high-throughput identification of promising candidates, we present the confusion matrices for all MLIPs in Figure 4. Each potential in Figure 4 is represented using a distinct color, such as green (MACE-MP-0), yellow (SevenNet), pink (Orb-v3), blue (CHGNet), and orange (M3GNet). For the classification task, we use a threshold E_m of 500 meV, i.e., materials that exhibit a calculated/predicted

Confusion Matrix

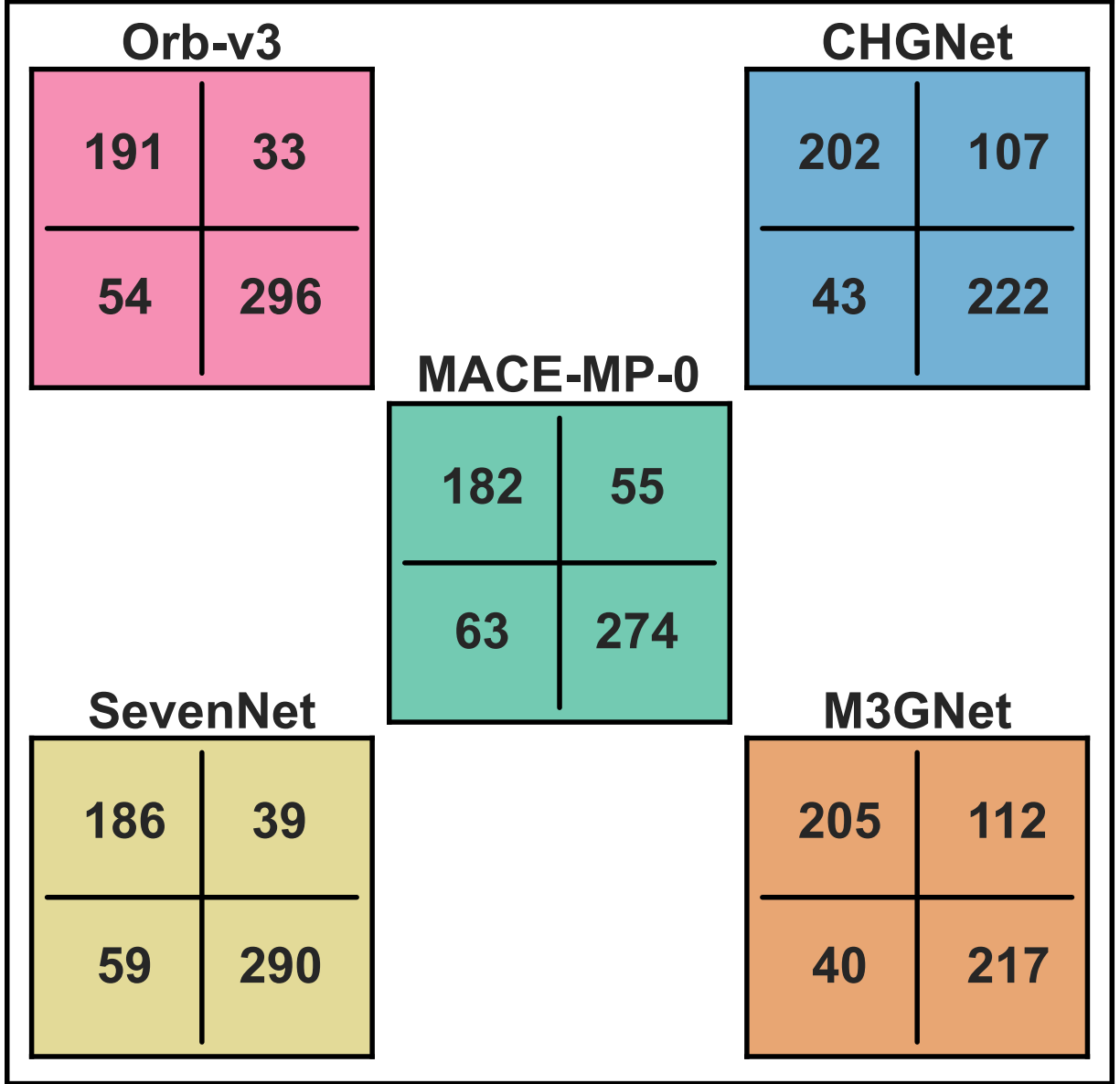


Figure 4: Confusion matrices for barrier prediction across different models. Each matrix corresponds to a specific model and is structured such that the upper-left, upper-right, lower-left, and lower-right cells represent the counts of true positives (TP), false positives (FP), false negatives (FN), and true negatives (TN), respectively. A prediction is considered a TP (TN) if both the DFT-computed and model-predicted E_m are less than (greater than or equal to) 500 meV.

$E_m < 500$ meV are labeled good ionic conductors, while materials that show higher values of E_m are labeled bad ionic conductors. Within each confusion matrix, the true positive (TP), the true negative (TN), the false positive (FP) and the false negative (FN) numbers are listed on the top left, bottom right, top right, and bottom left cells, respectively.

From Figure 4, we observe that Orb-v3 achieves the highest combined number of TP and TN, correctly

classifying 487 out of 574 systems (i.e., an accuracy of 84.84%). In comparison, M3GNet yields the lowest TP+TN count of 422 systems (73.52%). SevenNet, MACE-MP-0 and CHGNet correctly classify 476 (82.93%), 456 (79.44%), and 424 (73.87%) systems, respectively. These results highlight Orb-v3 as the most reliable model for distinguishing good and poor ionic conductors, followed closely by SevenNet (accuracy > 80%), making both reliable for high-throughput classification tasks.

3.4 Geometry Prediction Performance

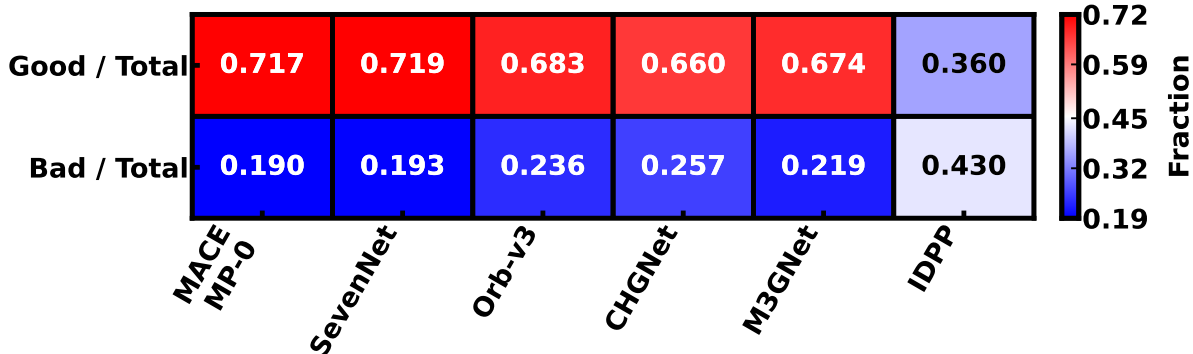


Figure 5: Performance of MLIPs on local geometry prediction. Each entry in the heatmap represents a performance fraction for a given MLIP with the last column corresponding to IDPP. The top (bottom) row shows the fraction of structures classified as ‘good’ (‘bad’) to the total number of structures. The heatmap color bar varies from red (high fractions) to blue (low fractions).

It is important for any MLIP to not only get the E_m correct but also get the underlying geometries that constitute the MEP (and hence yield the E_m value) correct for the MLIP to be truly accurate. Thus, we quantify the performance of the MLIPs considered in their predictions of local geometry of the intermediate image structures in Dataset-1 (using θ of Equation 1) as a heatmap in Figure 5. Note that we performed an NEB calculation with EB, $k = 5 \text{ eV}/\text{\AA}^2$, and IDPP interpolation with each MLIP and for each material in Dataset-1 to generate the statistics displayed in Figure 5. For each MLIP (x -axis), we denote the fraction of structures with ‘good’ (top row) and ‘bad’ (bottom) local geometries in Figure 5. Ideally, the MLIPs should exhibit a high (low) fraction of structures with good (bad) local geometry. Since IDPP (without any subsequent MLIP-based relaxation) also behaves like a potential for generating a guess for the MEP, we include IDPP’s statistics in Figure 5.

Among all the MLIPs, SevenNet exhibits the highest fraction of good geometries (0.719), indicating that it frequently generates accurate local geometries. On the other hand, MACE-MP-0 exhibits the lowest fraction of bad geometries (0.190), indicating that it frequently avoids generating inaccurate structures. The difference between the fraction of good and bad geometry predictions for both MACE-MP-0 and SevenNet are similar (0.527 and 0.526, respectively), indicating that both models perform equally well in generating good local geometries.

Other MLIPs show poorer geometry predictions, with Orb-v3, M3GNet, and CHGNet displaying good (bad) fractions of 0.683 (0.236), 0.674 (0.219), and 0.660 (0.257), respectively, with CHGNet showing the smallest difference between the good and bad geometry fractions (0.403). Thus, MACE-MP-0 and SevenNet show significantly better local geometry predictions upon relaxation with NEB compared to Orb-v3, M3GNet

and CHGNet, while all MLIPs provide better initial guesses to the MEP than LI in at least 66% of structures (i.e., intermediate images). Also, we note that IDPP generated structures are statistically much farther from DFT than MLIPs, with LI being better than IDPP in 43% of the cases. Given our definition of θ and the specific systems present in Dataset-1, we find that IDPP does not make a significant difference in enhancing the initial guess for the MEP as compared to LI across all MLIPs.

3.5 Geometry-Barrier Correlation

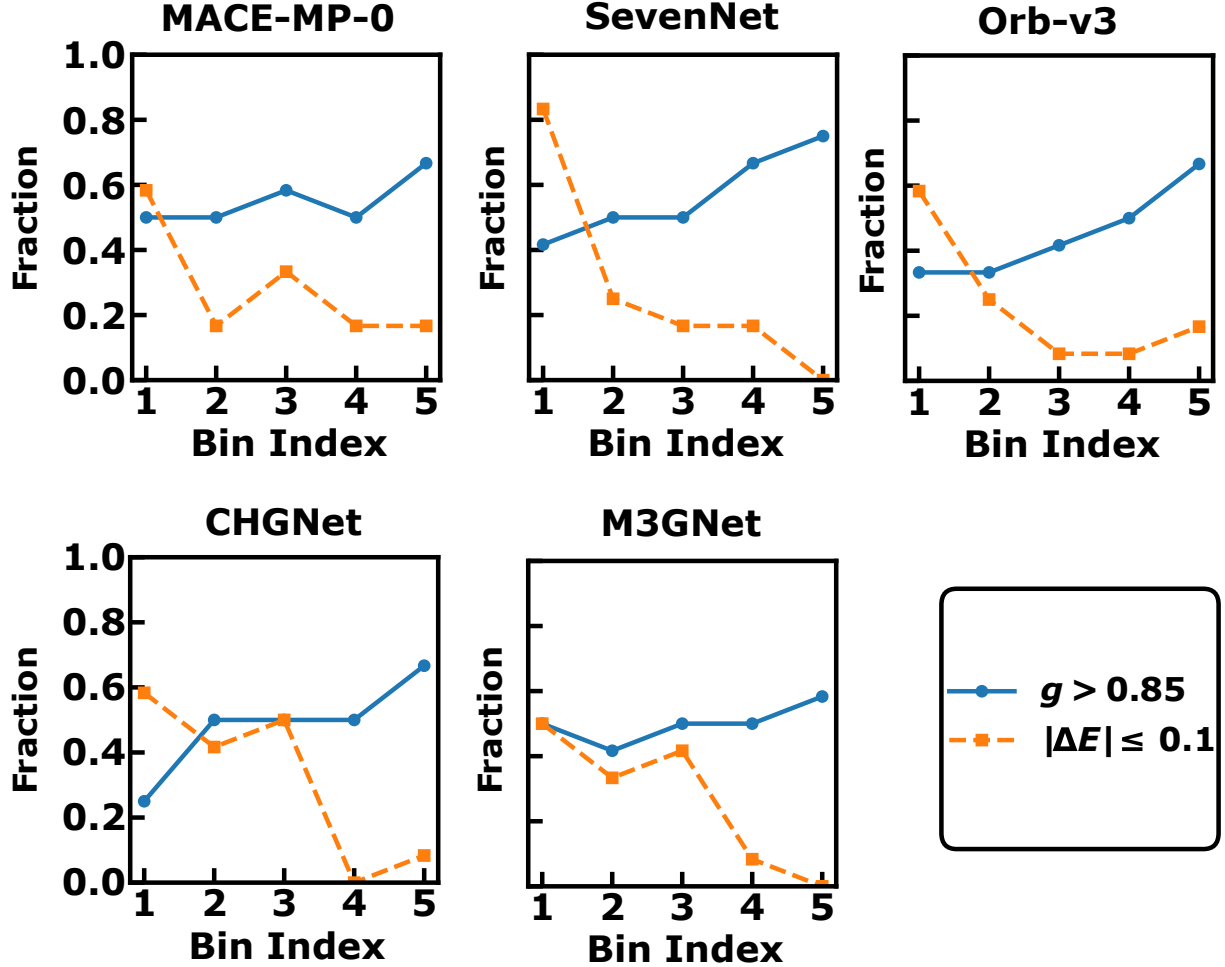


Figure 6: Correlation between E_m and local geometry prediction. Blue circles correspond to the fraction of data points within a given bin that has $g > 0.85$, while the orange rectangles represent the fraction of data points in the same bin which have an absolute error in the E_m prediction, $\Delta E \leq 0.1$ eV.

To investigate whether there is a correlation between geometry and E_m prediction performance, i.e., whether a precise E_m prediction by a given MLIP is due to its precise geometry prediction, we divide Dataset-1 into five bins based on their DFT-NEB E_m , ensuring that each bin contains an equal number of data points. Bins with indices 1, 2, 3, 4, and 5 correspond to E_m ranges of [0.058,0.64], (0.64,0.97], (0.97,1.24], (1.24,1.81], and (1.81,2.88] eV, respectively. Note that the number of bins in Figure 6 is different from that of Figure 3 since the datasets considered for both figures are different. Within each bin, we

estimate the fraction of migration paths with ‘good’ local geometry ($g > 0.85$, blue points) and low absolute errors in E_m ($\Delta E \leq 0.1$ eV, orange points), and plot the statistics for each MLIP in Figure 6. Note that $g > 0.85$ signifies cases where the predicted geometry is better than unrelaxed LI for at least six out of the seven intermediate images.

Overall, Figure 6 reveals the absence of any positive correlation between barrier and geometry prediction performance, and more strikingly, an inverse relationship. For example, all models perform poorly in predicting high E_m (bin 5), which is consistent with our observations in Figure 4. However, all models also achieve their best geometry predictions for bin 5. In other words, the best geometry predictions are coincident with the worst E_m predictions. The geometry prediction success rates within bin-5 are 66.7%, 75.0%, 66.7%, 66.7%, and 58.3% for MACE-MP-0, SevenNet, Orb-v3, CHGNet, and M3GNet, respectively, and the corresponding E_m prediction success rates (i.e., $\Delta E \leq 0.1$ eV) are 16.7%, 0%, 16.7%, 8.3%, and 0%, respectively.

To further assess geometry-barrier correlation, we examine instances where MLIPs perform well in both metrics. Note that, we term a model to exhibit a ‘good performance in both metrics’ if both the fractions in a given bin are > 0.5 . Only two potentials show this good performance, and only in a single bin (Bin-1), namely, MACE-MP-0 with a success rate of 58.3% in E_m prediction and 50.0% in geometry prediction, and M3GNet with a 50.0% success rate for both metrics. SevenNet, Orb-v3, and CHGNet do not achieve this good performance in any bin. Moreover, we find no consistent pattern across all bins and all MLIPs and no instances where good E_m predictions coincide with good geometry prediction. Instead, the data suggests that these two performance metrics are largely independent, and that a good E_m prediction does not necessarily arise from a good local geometry prediction (and vice-versa).

4 Discussion

Given the critical role of E_m in battery materials design and the high computational costs associated with DFT-NEB calculations, we have evaluated the performance of foundational MLIPs including MACE-MP-0, SevenNet, Orb-v3, CHGNet, and M3GNet, for E_m predictions upon integration with the NEB framework over two data subsets containing E_m and structural data (Figure 1). Specifically, we investigated *i*) the ability of MLIPs to predict E_m accurately, *ii*) the likelihood of generating MLIP-NEB-relaxed image geometries that are close to the ground truth (DFT-NEB), and *iii*) whether any correlation exists between the accuracy in E_m prediction and geometry relaxation.

Analyzing E_m predictions across the entire Dataset-2, we find that MACE-MP-0 exhibits the lowest MAE (Figure 2), followed in order by Orb-v3, CHGNet, SevenNet, and M3GNet. On excluding outliers that are common to all models, we observe SevenNet to exhibit a slightly lower MAE than CHGNet, with the rest of the performance order being the same. Interestingly, when assessing each model independently after removing their respective outliers, Orb-v3 demonstrates the best MAE of 0.198 eV, marginally outperforming MACE-MP0 (0.202 eV), with the other models exhibiting larger errors (0.203-0.257 eV). Thus, Orb-v3 provides the best prediction errors for E_m , among MLIPs considered, in systems with a robust description of the corresponding potential energy surface.

Notably, during endpoint relaxations for Orb-v3, 153 systems failed to converge within the threshold forces over 1000 optimization steps despite attempting multiple optimization algorithms. However, to maintain consistency with the other models in the study, we did not modify the obtained results and included the Orb-v3 results as is. Nevertheless, more robust relaxation strategies with extended optimization steps may enable

Orb-v3 to achieve better accuracies in E_m predictions, potentially surpassing the other MLIPs considered across the entire Dataset-2.

We observe that M3GNet and CHGNet exhibit a systematic bias toward underestimating E_m , whereas MACE-MP-0, SevenNet, and Orb-v3 do not display such a tendency. A more granular analysis (Figure 3) reveals that all models struggle with accurately predicting high E_m values. Among them, Orb-v3 shows a relatively slow decay in prediction accuracy as the E_m value increases. Interestingly, the simpler models CHGNet and M3GNet outperform their more complex counterparts within a very narrow range of low E_m but exhibit a rapid decline in performance as the range expands.

Using a threshold E_m of 500 meV to categorize structures as ‘good’ or ‘bad’ conductors of ions (Figure 4), we find that all MLIPs are able to identify good conductors with reasonable accuracy ($>73\%$). Orb-v3 and SevenNet display the highest accuracies in classifying good (or bad) conductors, with $\sim 85\%$ and $\sim 83\%$ accuracy, respectively, making them highly suitable for high-throughput screening of candidate battery materials.

Our study on Dataset-1 indicates that MLIP-NEB relaxations tend to produce image geometries that are as close as (or closer to) DFT-NEB structures than those obtained through simple LI or IDPP interpolation in the majority ($\sim 66\%$, Figure 5) of cases. Among the considered models, MACE-MP-0 and SevenNet stand out in geometry predictions, relaxing to geometries that are worse than LI or IDPP ones in only 19% of migration paths, suggesting that employing MACE-MP-0 or SevenNet NEB-relaxed images as initial guesses for DFT-NEB calculations could significantly accelerate convergence and reduce computational costs. Note that although our metric, θ (see Equation 1), captures critical local geometric features, it can be improved further to decisively quantify local structural similarity.

Finally, when simultaneously evaluating the likelihood of accurate barrier prediction and better geometry initialization (Figure 6), we observe no evident correlation between the two among all MLIPs considered. Thus, we find that accurate barrier predictions do not necessarily imply better geometry predictions, and vice versa. One possible explanation for this counterintuitive trend is that for systems with low E_m , the potential energy surfaces are likely ‘flat’ with variations in local geometries, indicating that even large errors in local bond distances or local bond angles made by the MLIPs do not significantly change the predicted E_m , thus leading to accurate E_m even with inaccurate geometries. On the other hand, for systems with large E_m , the potential energy surfaces should exhibit ‘deep’ minima associated with the ‘stable’ sites occupied by the migrating ion, signifying that even small errors in predicting local bond distances or angles by the MLIPs can cause large errors in the E_m predictions, thus resulting in inaccurate E_m even with mostly accurate geometries.

5 Conclusion

Given the importance of accurate and swift predictions of E_m in materials for battery applications (and beyond), we systematically evaluated five foundational MLIPs (MACE-MP-0, SevenNet, Orb-v3, CHGNet, and M3GNet) via integration with the NEB framework across a diverse set of chemistries and materials relevant for batteries. We benchmarked the barrier prediction accuracy against DFT-NEB calculated values for all MLIPs, and found MACE-MP-0 to achieve the lowest MAE (0.310 eV), while Orb-v3 demonstrated significantly better performance (MAE of 0.198 eV) when evaluated on data points that were not outliers. Further, we assessed the capability of the MLIPs to classify materials as good ($E_m < 500$ meV) or bad ($E_m \geq 500$ meV) ionic diffusers, and observed Orb-v3 and SevenNet to accurately classify $>82\%$ of migration paths,

making them suitable for high-throughput screening applications. Based on our novel geometric similarity metric, we demonstrated that MLIP-NEB relaxations produce image structures that are closer to DFT-NEB calculated references than LI in over 66% of cases. Finally, we discovered no direct correlation between barrier prediction accuracy and the similarity of the MLIP-NEB relaxed geometry to the DFT-NEB reference. We hope that our work establishes use-cases, accuracies, and limitations in using foundational MLIPs for predicting and/or accelerating E_m calculations, which should result in better discovery of novel ionic conductors with applications as electrodes and (solid) electrolytes in batteries and other related technologies.

6 Data availability

Both datasets and associated python scripts used in this work are compiled and available for free at our GitHub repository. Dataset-2 is available as a json file on Zenodo.

7 Acknowledgements

G.S.G. acknowledges financial support from the Science and Engineering Research Board (SERB) of the Department of Science and Technology, Government of India, under sanction number IPA/2021/000007. A.K.B. thanks the Ministry of Human Resource Development, Government of India, for financial assistance. The authors gratefully acknowledge the super-computing facility offered by ACENET and the Digital Research Alliance of Canada. The authors acknowledge the computational resources of the super computer ‘PARAM Pravega’ provided by the super computer education and research centre (SERC) at IISc. The authors also thank the Jülich Supercomputing Centre (at Forschungszentrum Jülich), Germany for the use of the ‘JURECA’ supercomputer, under projects ‘hpc-prf-emdft’ and ‘hpc-prf-desal’.

8 Competing Interests

The authors declare no competing financial or non-financial interests.

References

- [1] Matthew Li, Jun Lu, Zhongwei Chen, and Khalil Amine. 30 years of lithium-ion batteries. *Advanced materials*, 30(33):1800561, 2018.
- [2] M Stanley Whittingham. Ultimate limits to intercalation reactions for lithium batteries. *Chemical Reviews*, 114(23):11413–11413, 2014.
- [3] Myounggu Park, Xiangchun Zhang, Myoungdo Chung, Gregory B Less, and Ann Marie Sastry. A review of conduction phenomena in li-ion batteries. *Journal of power sources*, 195(24):7904–7929, 2010.
- [4] John Christopher Bachman, Sokseiha Muy, Alexis Grimaud, Hao-Hsun Chang, Nir Pour, Simon F Lux, Odysseas Paschos, Filippo Maglia, Saskia Lupart, Peter Lamp, et al. Inorganic solid-state electrolytes for lithium batteries: mechanisms and properties governing ion conduction. *Chemical reviews*, 116(1):140–162, 2016.

- [5] Eibar Flores, Christian Wölke, Peng Yan, Martin Winter, Tejs Vegge, Isidora Cekic-Laskovic, and Arghya Bhowmik. Learning the laws of lithium-ion transport in electrolytes using symbolic regression. *Digital Discovery*, 1(4):440–447, 2022.
- [6] George H Vineyard. Frequency factors and isotope effects in solid state rate processes. *Journal of Physics and Chemistry of Solids*, 3(1-2):121–127, 1957.
- [7] Anton Van der Ven, Zhi Deng, Swastika Banerjee, and Shyue Ping Ong. Rechargeable alkali-ion battery materials: theory and computation. *Chemical reviews*, 120(14):6977–7019, 2020.
- [8] Pieremanuele Canepa, Gopalakrishnan Sai Gautam, Daniel C Hannah, Rahul Malik, Miao Liu, Kevin G Gallagher, Kristin A Persson, and Gerbrand Ceder. Odyssey of multivalent cathode materials: open questions and future challenges. *Chemical reviews*, 117(5):4287–4341, 2017.
- [9] Yuki Oriksa, Titus Masese, Yukinori Koyama, Takuya Mori, Masashi Hattori, Kentaro Yamamoto, Tet-suya Okado, Zhen-Dong Huang, Taketoshi Minato, Cédric Tassel, et al. High energy density rechargeable magnesium battery using earth-abundant and non-toxic elements. *Scientific reports*, 4(1):5622, 2014.
- [10] Qirong Liu, Haitao Wang, Chunlei Jiang, and Yongbing Tang. Multi-ion strategies towards emerging rechargeable batteries with high performance. *Energy Storage Materials*, 23:566–586, 2019.
- [11] Yanliang Liang, Hui Dong, Doron Aurbach, and Yan Yao. Current status and future directions of multivalent metal-ion batteries. *Nature Energy*, 5(9):646–656, 2020.
- [12] Ryan D Bayliss, Baris Key, Gopalakrishnan Sai Gautam, Pieremanuele Canepa, Bob Jin Kwon, Saul H Lapidus, Fulya Dogan, Abdullah A Adil, Andrew S Lipton, Peter J Baker, et al. Probing mg migration in spinel oxides. *Chemistry of Materials*, 32(2):663–670, 2019.
- [13] Gopalakrishnan Sai Gautam, Xiaoqi Sun, Victor Duffort, Linda F Nazar, and Gerbrand Ceder. Impact of intermediate sites on bulk diffusion barriers: Mg intercalation in $\text{mg}_2\text{mo}_3\text{o}_8$. *Journal of Materials Chemistry A*, 4(45):17643–17648, 2016.
- [14] Dong Wang, Zhenyu Zhang, Yue Hao, Hongxing Jia, Xing Shen, Baihua Qu, Guangsheng Huang, Xiaoyuan Zhou, Jingfeng Wang, Chaohe Xu, et al. Challenges and progress in rechargeable magnesium-ion batteries: materials, interfaces, and devices. *Advanced Functional Materials*, 34(51):2410406, 2024.
- [15] Lin Zhu, Jia-Ying Xie, Guo-Miao Zhou, De-An Zhang, and An Du. First principles investigation of voltage, structure, ionic and electronic conduction of olivine and maricite namnpo_4 . *Solid State Ionics*, 398:116274, 2023.
- [16] Debolina Deb and Gopalakrishnan Sai Gautam. Critical overview of polyanionic frameworks as positive electrodes for na-ion batteries. *Journal of Materials Research*, 37(19):3169–3196, 2022.
- [17] Bettina Schwaighofer, Miguel A Gonzalez, Mark R Johnson, John SO Evans, and Ivana Radosavljevic Evans. Ionic mobility in energy materials: Through the lens of quasielastic neutron scattering. *Chemistry of Materials*, 37(10):3575–3593, 2025.
- [18] Shangshang Wang, Jianbo Zhang, Oumaïma Gharbi, Vincent Vivier, Ming Gao, and Mark E Orazem. Electrochemical impedance spectroscopy. *Nat. Rev. Methods Primers*, 1(1), June 2021.

- [19] Raphaële J Clément, Peter G Bruce, and Clare P Grey. manganese-based p2-type transition metal oxides as sodium-ion battery cathode materials. *Journal of The Electrochemical Society*, 162(14):A2589, 2015.
- [20] Stephen Dongmin Kang and William C Chueh. Galvanostatic intermittent titration technique reinvented: Part i. a critical review. *Journal of The Electrochemical Society*, 168(12):120504, 2021.
- [21] Paul Heitjans and Jörg Kärger. *Diffusion in condensed matter: methods, materials, models*. Springer Science & Business Media, 2006.
- [22] Pierre Hohenberg and Walter Kohn. Inhomogeneous electron gas. *Physical review*, 136(3B):B864, 1964.
- [23] Walter Kohn and Lu Jeu Sham. Self-consistent equations including exchange and correlation effects. *Physical review*, 140(4A):A1133, 1965.
- [24] Graeme Henkelman, Blas P Uberuaga, and Hannes Jónsson. A climbing image nudged elastic band method for finding saddle points and minimum energy paths. *The Journal of chemical physics*, 113(22):9901–9904, 2000.
- [25] Reshma Devi, Baltej Singh, Pieremanuele Canepa, and Gopalakrishnan Sai Gautam. Effect of exchange-correlation functionals on the estimation of migration barriers in battery materials. *npj Computational Materials*, 8(1):160, 2022.
- [26] Daan Frenkel and Berend Smit. *Understanding molecular simulation: from algorithms to applications*. Elsevier, 2023.
- [27] Xingfeng He, Yizhou Zhu, Alexander Epstein, and Yifei Mo. Statistical variances of diffusional properties from ab initio molecular dynamics simulations. *npj Computational Materials*, 4(1):18, 2018.
- [28] Ziqin Rong, Daniil Kitchaev, Pieremanuele Canepa, Wenxuan Huang, and Gerbrand Ceder. An efficient algorithm for finding the minimum energy path for cation migration in ionic materials. *The Journal of chemical physics*, 145(7), 2016.
- [29] Alexandre Duval, Simon V Mathis, Chaitanya K Joshi, Victor Schmidt, Santiago Miret, Fragkiskos D Malliaros, Taco Cohen, Pietro Lio, Yoshua Bengio, and Michael Bronstein. A hitchhiker’s guide to geometric gnns for 3d atomic systems. *arXiv preprint arXiv:2312.07511*, 2023.
- [30] Oliver T Unke, Stefan Chmiela, Huziel E Sauceda, Michael Gastegger, Igor Poltavsky, Kristof T Schutt, Alexandre Tkatchenko, and Klaus-Robert Muller. Machine learning force fields. *Chemical Reviews*, 121(16):10142–10186, 2021.
- [31] Junyoung Choi, Gunwook Nam, Jaesik Choi, and Yousung Jung. A perspective on foundation models in chemistry. *JACS Au*, 5(4):1499–1518, 2025.
- [32] Fabian L Thiemann, Niamh O’neill, Venkat Kapil, Angelos Michaelides, and Christoph Schran. Introduction to machine learning potentials for atomistic simulations. *Journal of Physics: Condensed Matter*, 37(7):073002, 2024.
- [33] Xiang Fu, Zhenghao Wu, Wujie Wang, Tian Xie, Sinan Keten, Rafael Gomez-Bombarelli, and Tommi Jaakkola. Forces are not enough: Benchmark and critical evaluation for machine learning force fields with molecular simulations. *arXiv preprint arXiv:2210.07237*, 2022.

- [34] Ryan Jacobs, Dane Morgan, Siamak Attarian, Jun Meng, Chen Shen, Zhenghao Wu, Clare Yijia Xie, Julia H Yang, Nongnuch Artrith, Ben Blaiszik, et al. A practical guide to machine learning interatomic potentials—status and future. *Current Opinion in Solid State and Materials Science*, 35:101214, 2025.
- [35] Suyeon Ju, Jinmu You, Gijin Kim, Yutack Park, Hyungmin An, and Seungwu Han. Application of pretrained universal machine-learning interatomic potential for physicochemical simulation of liquid electrolytes in li-ion batteries. *Digital Discovery*, 2025.
- [36] Hanwen Kang, Tenglong Lu, Zhanbin Qi, Jiandong Guo, Sheng Meng, and Miao Liu. Fasttrack: a fast method to evaluate mass transport in solid leveraging universal machine learning interatomic potential. *AI for Science*, 1(1):015004, 2025.
- [37] Bruno Focassio, Luis Paulo M. Freitas, and Gabriel R Schleder. Performance assessment of universal machine learning interatomic potentials: Challenges and directions for materials’ surfaces. *ACS Applied Materials & Interfaces*, 17(9):13111–13121, 2024.
- [38] Antoine Loew, Dewen Sun, Hai-Chen Wang, Silvana Botti, and Miguel AL Marques. Universal machine learning interatomic potentials are ready for phonons. *npj Computational Materials*, 11(1):178, 2025.
- [39] Sajid Mannan, Carmelo Gonzales, Vaibhav Bihani, Kin Long Kelvin Lee, Nitya Nand Gosvami, Santiago Miret, and NM Anoop Krishnan. Evaluating universal interatomic potentials for molecular dynamics of real-world minerals. In *AI for Accelerated Materials Design-ICLR 2025*, 2025.
- [40] Haochen Yu, Matteo Giantomassi, Giuliana Materzanini, Junjie Wang, and Gian-Marco Rignanese. Systematic assessment of various universal machine-learning interatomic potentials. *Materials Genome Engineering Advances*, 2(3):e58, 2024.
- [41] Janosh Riebesell, Rhys EA Goodall, Philipp Benner, Yuan Chiang, Bowen Deng, Gerbrand Ceder, Mark Asta, Alpha A Lee, Anubhav Jain, and Kristin A Persson. A framework to evaluate machine learning crystal stability predictions. *Nature Machine Intelligence*, 7(6):836–847, 2025.
- [42] Qiyuan Zhao, Yunhong Han, Duo Zhang, Jiaxu Wang, Peichen Zhong, Taoyong Cui, Bangchen Yin, Yirui Cao, Haojun Jia, and Chenru Duan. Harnessing machine learning to enhance transition state search with interatomic potentials and generative models. *Advanced Science*, 12:e06240, 2025.
- [43] Vaibhav Bihani, Sajid Mannan, Utkarsh Pratiush, Tao Du, Zhimin Chen, Santiago Miret, Matthieu Micoulaut, Morten M Smedskjaer, Sayan Ranu, and NM Anoop Krishnan. Egraffbench: evaluation of equivariant graph neural network force fields for atomistic simulations. *Digital Discovery*, 3(4):759–768, 2024.
- [44] Ilyes Batatia, Philipp Benner, Yuan Chiang, Alin M Elena, Dávid P Kovács, Janosh Riebesell, Xavier R Advincula, Mark Asta, Matthew Avaylon, William J Baldwin, et al. A foundation model for atomistic materials chemistry. *arXiv preprint arXiv:2401.00096*, 2023.
- [45] 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-centred interatomic potentials. *Nature Machine Intelligence*, 7(1):56–67, 2025.

- [46] Yutack Park, Jaesun Kim, Seungwoo Hwang, and Seungwu Han. Scalable parallel algorithm for graph neural network interatomic potentials in molecular dynamics simulations. *Journal of chemical theory and computation*, 20(11):4857–4868, 2024.
- [47] Jaesun Kim, Jisu Kim, Jaehoon Kim, Jiho Lee, Yutack Park, Youngho Kang, and Seungwu Han. Data-efficient multifidelity training for high-fidelity machine learning interatomic potentials. *Journal of the American Chemical Society*, 147(1):1042–1054, 2024.
- [48] Mark Neumann, James Gin, Benjamin Rhodes, Steven Bennett, Zhiyi Li, Hitarth Choubisa, Arthur Hussey, and Jonathan Godwin. Orb: A fast, scalable neural network potential. *arXiv preprint arXiv:2410.22570*, 2024.
- [49] Benjamin Rhodes, Sander Vandenhaute, Vaidotas Šimkus, James Gin, Jonathan Godwin, Tim Duignan, and Mark Neumann. Orb-v3: atomistic simulation at scale. *arXiv preprint arXiv:2504.06231*, 2025.
- [50] Bowen Deng, Peichen Zhong, KyuJung Jun, Janosh Riebesell, Kevin Han, Christopher J Bartel, and Gerbrand Ceder. Chgnet as a pretrained universal neural network potential for charge-informed atomistic modelling. *Nature Machine Intelligence*, 5(9):1031–1041, 2023.
- [51] Chi Chen and Shyue Ping Ong. A universal graph deep learning interatomic potential for the periodic table. *Nature Computational Science*, 2(11):718–728, 2022.
- [52] Reshma Devi, Avaneesh Balasubramanian, Keith T Butler, and Gopalakrishnan Sai Gautam. A literature-derived dataset of migration barriers for quantifying ionic transport in battery materials. *Scientific Data*, 2025.
- [53] Reshma Devi, Avaneesh Balasubramanian, Keith T. Butler, and Gopalakrishnan Sai Gautam. Dft-neb-migration-barrier-dataset v1.0, 2025.
- [54] John P Perdew, Kieron Burke, and Matthias Ernzerhof. Generalized gradient approximation made simple. *Physical review letters*, 77(18):3865, 1996.
- [55] Dereje Bekele Tekliye, Ankit Kumar, Xie Weihang, Thelakkattu Devassy Mercy, Pieremanuele Canepa, and Gopalakrishnan Sai Gautam. Exploration of nasicon frameworks as calcium-ion battery electrodes. *Chemistry of Materials*, 34(22):10133–10143, 2022.
- [56] Dereje Bekele Tekliye and Gopalakrishnan Sai Gautam. Fluoride frameworks as potential calcium battery cathodes. *Journal of Materials Chemistry A*, 12(30):18993–19007, 2024.
- [57] Debolina Deb and Gopalakrishnan Sai Gautam. Exploration of oxyfluoride frameworks as na-ion cathodes. *Chemistry of Materials*, 36(24):11892–11904, 2024.
- [58] 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.
- [59] 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.

- [60] B Deng, Z Peichen, KJ Jun, R Janosh, K Han, CJ Bartel, and CJ Gerbrand. Materials project trajectory (mptrj) dataset. *Figshare*. <https://doi.org/10.6084/m9.figshare.23713842.v2>, 2023.
- [61] 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.
- [62] Luis Barroso-Luque, Muhammed Shuaibi, Xiang Fu, Brandon M Wood, Misko Dzamba, Meng Gao, Ammar Rizvi, C Lawrence Zitnick, and Zachary W Ulissi. Open materials 2024 (omat24) inorganic materials dataset and models. *arXiv preprint arXiv:2410.12771*, 2024.
- [63] Jonathan Schmidt, Tiago FT Cerqueira, Aldo H Romero, Antoine Loew, Fabian Jäger, Hai-Chen Wang, Silvana Botti, and Miguel AL Marques. Improving machine-learning models in materials science through large datasets. *Materials Today Physics*, 48:101560, 2024.
- [64] Søren Smidstrup, Andreas Pedersen, Kurt Stokbro, and Hannes Jónsson. Improved initial guess for minimum energy path calculations. *The Journal of chemical physics*, 140(21), 2014.
- [65] Esben L Kolsbjerg, Michael N Groves, and Bjørk Hammer. An automated nudged elastic band method. *The Journal of chemical physics*, 145:094107, 2016.
- [66] CG Broyden. The convergence of a class of double rank minimization algorithms. *Journal of the Institute of Mathematics and its Applications*, 6:75–90, 1970.
- [67] Roger Fletcher. A new approach to variable metric algorithms. *The computer journal*, 13(3):317–322, 1970.
- [68] Donald Goldfarb. A family of variable-metric methods derived by variational means. *Mathematics of computation*, 24(109):23–26, 1970.
- [69] David F Shanno. Conditioning of quasi-newton methods for function minimization. *Mathematics of computation*, 24(111):647–656, 1970.
- [70] M O’Keeffe. A proposed rigorous definition of coordination number. *Foundations of Crystallography*, 35(5):772–775, 1979.
- [71] 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.
- [72] Ziqin Rong, Rahul Malik, Pieremanuele Canepa, Gopalakrishnan Sai Gautam, Miao Liu, Anubhav Jain, Kristin Persson, and Gerbrand Ceder. Materials design rules for multivalent ion mobility in intercalation structures. *Chemistry of Materials*, 27(17):6016–6021, 2015.