Benchmarking Machine Learning Potentials for Crystal Structure Relaxation

Anonymous Author(s)

Affiliation Address email

Abstract

High-throughput materials discovery workflows require rapid and accurate relaxation of crystal structures to identify thermodynamically stable phases among thousands to millions of candidate structures. Yet current machine learning interatomic potential (MLIP) benchmarks focus predominantly on energy prediction rather than structure relaxation, creating a critical evaluation gap for models designed to accelerate optimization. Additionally, these benchmarks are trained on datasets consisting mainly of known stable or near-stable materials, thus failing to capture the challenges of unexplored chemical spaces. We address these limitations by introducing a benchmark that evaluates state-of-the-art MLIPs and a one-shot relaxation model on structure relaxation with crystals generated via a reinforcement learning pipeline. We compare energy lowering and average maximum force computed via DFT, as well as relaxation runtime. We also contrast direct force-prediction strategies against conservative energy-differentiation approaches to determine which paradigm delivers superior relaxation performance. Our results indicate that there is a clear disconnect between MLIP energy prediction and force convergence in relaxation, challenging current benchmarking approaches.

7 1 Introduction

2

3

5

6

7

8

10

11

12

13

14

15

16

19

20

21

22

23

24

25

27

28

29

30

31

34

35

The discovery of new materials with desired properties presents an important challenge in modern materials science, with applications including energy storage, catalysis, and electronics. Central to these discovery pipelines is structure relaxation: the process of optimizing atomic positions and structural parameters to energetically stable configurations. For crystalline materials, property predictions are sensitive to their underlying structure, and having crystals in a stable, relaxed state is crucial [13]. High-throughput screening of materials routinely involves relaxation of thousands to millions of candidate structures, making the efficiency and accuracy of structure optimization a critical bottleneck in materials discovery pipelines. Crystal structure relaxation is generally performed with density functional theory (DFT), which provides a first-principles estimate of the energy and forces of a system. However, DFT's computational cost scales cubically with the system size[6], making it difficult to massively screen the vast chemical space in short time. This trade-off between accuracy and throughput has long restricted the acceleration of materials discovery with computers. Machine learning interatomic potentials (MLIPs) have emerged as a promising solution to this computational bottleneck, offering the potential to achieve near-DFT accuracy at a fraction of the computational cost. MLIPs involve the mapping of atomic structures to the potential energy and forces, primarily via supervised learning architectures[5]. Recent advances in graph neural networks, equivariant architectures, and transformer-based models have demonstrated impressive performance on energy and force prediction tasks across diverse materials systems. MatBench Discovery [15] represents the most extensive benchmarking effort to date for assessing MLIPs in the prediction of

thermodynamic stability through energy and force-based metrics. However, because the benchmark 37 centers on classification tasks, such as convex hull distance estimation and regression errors, it does 38 not systematically evaluate MLIPs on the actual structure relaxation process of inorganic crystals 39 based on energy minimization. Consequently, rigorous comparative evaluation of MLIP-driven 40 geometry relaxations in inorganic systems remains limited. Further, the established benchmarks rely 41 on datasets derived from already-known stable or near-stable materials, such as the Materials Project 42 (MPTraj) derivatives used in WBM[16], creating an artificial selection bias toward well-explored chemical spaces. It is therefore unclear if MLIPs can generalize to unknown chemical systems.

To address this gap, we evaluate the performance of state-of-the-art MLIPs on crystal structures 45 generated through a reinforcement learning (RL) pipeline using CrystalGym[10]: CrystalGym is an 46 open-source RL environment for crystal generation based on the Gymnasium framework for training 47 policies with reward signals obtained directly from DFT. We evaluate 6 diverse architectures from 48 MatBench Discovery including universal graph neural network potentials (M3GNet[3], CHGNet[4], 49 MACE[2]), transformer-based approaches (EquiformerV2[12], eSEN-30M[8]), and direct structure prediction methods (DeepRelax[17]) on crystal structures obtained from trained CrystalGym RL policies. For evaluation, we compare the average formation energy difference and maximum atomic 52 force obtained using DFT simulations. We also compare the differences in formation energies relative 53 to the structure relaxed entirely with DFT. In practice, our energy-based analyses offer a more reliable 54 measure of relaxation accuracy than metrics such as RMSD, which only quantify structural deviations 55 from the reference state. The dataset of RL-generated crystals and their DFT-relaxed states will be 56 released upon publication.

Background 2

59

81

82

87

MLIP Architectures

CHGNet[4] is a graph-based MLIP that integrates site-specific magnetic moments (as proxies for 60 charge information) into its pretrained universal potential to capture both atomic positions and electronic orbital occupancy. MACE[2] enhances expressivity and efficiency using higher-body equivariant message passing, drastically reducing the number of layers needed while maintaining 63 fast, accurate interatomic force predictions. M3GNet is a materials graph neural network that incorporates explicit three-body interactions, atomic coordinates, and full lattice tensors, enabling 65 accurate tensorial predictions (forces and stresses) across the periodic table via auto-differentiation. 66 Additionally, we include an iteration-free structure relaxation approach, DeepRelax, which is not 67 an MLIP. DeepRelax is a generative model capable of directly predicting relaxed crystal structures 68 69 without iterative energy minimization. Using a periodicity-aware equivariant GNN, it achieves 100fold speed improvements over iterative models while maintaining competitive accuracy. CHGNet, 70 MACE, and M3GNet have been trained on $\sim 1.3M$ structures from the MPTraj dataset, and are 71 conservative models: atomic forces are obtained by differentiating the system energy. DeepRelax, in 72 contrast, is trained on the X-MN-O dataset[11] derived from the MP database. EquiformerV2[12] 73 presents a significant advancement in equivariant transformers, using eSCN convolutions to scale 74 to higher degree representations along with attention re-normalization, separable S^2 activation, 75 and separable layer normalization. Another key feature in our study is the utilization of direct force prediction compared to conservative force prediction. eSEN[8] is a message-passing neural 77 network that processes atomic structures through alternating edgewise and node-wise operations, 78 with atoms embedded as multi-channel spherical harmonic representations, maintaining continuous 79 representations throughout the network and significantly improving energy conservation. 80

2.2 Current Benchmarking Challenges

Evaluating ML models for materials applications presents unique challenges. Bartel et al.[1] demonstrated that accurate formation energy prediction does not necessarily translate to reliable stability predictions. In our work, we intend to evaluate whether the models optimized for accurate energy pre-85 dictions can perform well in structural relaxation, particularly with new chemical systems. Together, we address a broader theme: that the utility of ML models in materials discovery depends not only on the energy estimation but also on their ability to capture the overall energy landscapes. Current MLIP benchmarks predominantly focus on energy-centric metrics, typically evaluating performance by comparing predicted energies to DFT reference calculations[7]. MatBench Discovery evaluates

relaxation performance using RMSD between predicted and DFT-relaxed structures, but it does not validate these results against energies computed with DFT.

92 3 Methods

93 3.1 Experimental Setup and Dataset

We evaluate crystal structure relaxation performance across six state-of-the-art ML models. Five 94 models (CHGNet, MACE, M3GNet, eSEN, and EquiformerV2 trained on the MPTraj dataset) utilize 95 the Atomic Simulation Environment (ASE) library for structure optimization, while DeepRelax 96 employs a direct structure prediction approach without iterative optimization. For ASE-based models, 97 we maintained consistent use of the FIRE optimizer to ensure fair comparison. For DFT calculations, 98 we use Quantum Espresso v7.3 [9], an open-source software suite for atomic simulations. We generated crystals by performing rollouts using different RL policies trained on the CrystalGym 100 environment[10]. The original set of tasks of CrystalGym was to train policies to design crystal 101 compositions on known crystal structures (from Materials Project) and optimize properties such 102 as bulk modulus, band gap, and density. As structure optimization is not included as part of the 103 original CrystalGym pipeline, the generated crystals are not relaxed. By choosing structures with 105 unique compositions, we obtain around 1000 crystals. Further background is provided in Appendix A. The distribution of formation energies of these crystals is shown in the Appendix 3. The presence of crystals with positive formation energies indicates that many of them are thermodynamically unfavorable. 108

3.2 Evaluation

109

126

We assess model performance using three complementary metrics that capture different aspects of relaxation quality and computational efficiency.

Formation energy reduction measures the difference in DFT-computed formation energies between initial and model-relaxed structures. This metric evaluates whether models successfully identify energetically favorable relaxation trajectories, with larger reductions indicating superior thermodynamic optimization performance. Additionally, we include the formation energy difference between relaxed structures from DFT and the MLIPs (Figure 1b) to highlight cases where MLIPs achieve better formation energy lowering than DFT.

Maximum atomic force represents the largest force magnitude across all atoms in the relaxed structure. This metric directly quantifies structural equilibrium quality, as well-relaxed crystals should exhibit near-zero atomic forces. Lower maximum forces indicate superior convergence to local energy minima.

Optimization runtime measures the computational time required for structure relaxation, providing essential insights into model efficiency for high-throughput materials discovery applications. This metric is particularly crucial for evaluating the practical deployment potential of different ML approaches as DFT proxies.

4 Results and Discussion

Our comprehensive evaluation reveals distinct performance profiles with clear trade-offs between computational efficiency and relaxation quality across three key metrics.

Formation Energy Reduction: While all models except for DeepRelax outperform DFT in formation energy reduction (Figure 1b), CHGNet and EquiformerV2 demonstrate superior thermodynamic optimization, achieving the most significant formation energy improvements compared to DFT (Figure 1a).

Force Convergence: eSEN and EquiformerV2 emerge as clear leaders, achieving force magnitudes >0.100 eV/Å, better than the next-best performer (M3GNet) and demonstrating superior ability to locate well-converged local minima (Figure 1d). Despite achieving the best energy statistics, CHGNet exhibits the second-worst force convergence, reinforcing that energy accuracy does not guarantee effective structural relaxation.

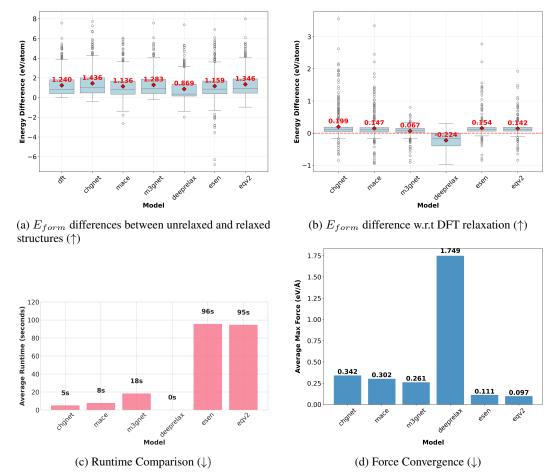


Figure 1: Evaluation of ML models for relaxation: (a) formation energy reduction achieved during relaxation (\uparrow), (b) formation energy difference between MLIP and DFT relaxation (\uparrow), (c) runtime of relaxation (\downarrow), (d) maximum force magnitudes in relaxed structures (\downarrow).

Computational Efficiency: DeepRelax offers relaxation times orders of magnitude faster than iterative approaches, but this speed comes at substantial accuracy costs in both formation energy and force convergence (Figure 1c). eSEN and EquiformerV2 provide superior accuracy but require significantly longer computation times with numerous outliers.

Overall Assessment: While CHGNet achieves the best formation energy, its force convergence is worse by a much larger margin compared to eSEN and eQV2. MACE emerges as an attractive middle ground, offering reasonable relaxation quality with moderate computational requirements, making it suitable for workflows requiring balanced throughput and accuracy, such as materials discovery pre-screening steps.

5 Conclusion

This work presents the first systematic evaluation of state-of-the-art MLIPs for crystal structure relaxation using unseen RL-generated materials, highlighting trade-offs between computational efficiency and relaxation quality. Our analysis reveals a critical disconnect between energy accuracy and relaxation performance: CHGNet achieves the best formation energy reduction but has the second-worst force convergence. This challenges energy-centric benchmarking and shows that energy prediction alone is insufficient for evaluating relaxation. Rather than identifying a single "best" model, our multi-metric evaluation emphasizes balancing runtime, formation energy reduction, and force convergence according to application needs. Future work should expand dataset size and model diversity to more comprehensively assess MLIPs as DFT proxies across chemical systems.

References

- 158 [1] Christopher J. Bartel, Amalie Trewartha, Qi Wang, Alexander Dunn, Anubhav Jain, and Gerbrand Ceder. A critical examination of compound stability predictions from machinelearned formation energies. *npj Computational Materials*, 6(1):97, Jul 2020.
- [2] Ilyes Batatia, David P Kovacs, Gregor Simm, Christoph Ortner, and Gábor Csányi. Mace:
 Higher order equivariant message passing neural networks for fast and accurate force fields.
 Advances in neural information processing systems, 35:11423–11436, 2022.
- [3] Chi Chen and Shyue Ong. A universal graph deep learning interatomic potential for the periodic
 table. April 2022.
- 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.
- [5] L. Fiedler, K. Shah, M. Bussmann, and A. Cangi. Deep dive into machine learning density functional theory for materials science and chemistry. *Phys. Rev. Mater.*, 6:040301, Apr 2022.
- [6] Lenz Fiedler, Normand A. Modine, Steve Schmerler, Dayton J. Vogel, Gabriel A. Popoola, Aidan P. Thompson, Sivasankaran Rajamanickam, and Attila Cangi. Predicting electronic structures at any length scale with machine learning. *npj Computational Materials*, 9(1):115, Jun 2023.
- 175 [7] 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, Mar 2025.
- 178 [8] Xiang Fu, Brandon M. Wood, Luis Barroso-Luque, Daniel S. Levine, Meng Gao, Misko
 179 Dzamba, and C. Lawrence Zitnick. Learning smooth and expressive interatomic potentials for
 180 physical property prediction, 2025.
- [9] Paolo Giannozzi, Stefano Baroni, Nicola Bonini, Matteo Calandra, Roberto Car, Carlo Cavazzoni, Davide Ceresoli, Guido L Chiarotti, Matteo Cococcioni, Ismaila Dabo, et al. Quantum espresso: a modular and open-source software project for quantum simulations of materials.
 Journal of physics: Condensed matter, 21(39):395502, 2009.
- [10] Prashant Govindarajan, Mathieu Reymond, Antoine Clavaud, Mariano Phielipp, Santiago Miret,
 and Sarath Chandar. Crystalgym: A new benchmark for materials discovery using reinforcement
 learning. In AI for Accelerated Materials Design ICLR 2025, 2025.
- 188 [11] Sungwon Kim. A structure translation model for crystal compounds release for manuscript acceptance. 2023.
- 190 [12] Yi-Lun Liao, Brandon Wood, Abhishek Das, and Tess Smidt. Equiformerv2: Improved equivariant transformer for scaling to higher-degree representations. *arXiv preprint arXiv:2306.12059*, 2023.
- 193 [13] Artem R. Oganov, Chris J. Pickard, Qiang Zhu, and Richard J. Needs. Structure prediction drives materials discovery. *Nature Reviews Materials*, 4(5):331–348, May 2019.
- [14] Gianluca Prandini, Antimo Marrazzo, Ivano E Castelli, Nicolas Mounet, and Nicola Marzari.
 Precision and efficiency in solid-state pseudopotential calculations. npj computational materials,
 4 (1): 72, 2018.
- [15] Janosh Riebesell, Rhys E. A. 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, Jun
 2025.
- [16] Hai-Chen Wang, Silvana Botti, and Miguel A. L. Marques. Predicting stable crystalline compounds using chemical similarity. *npj Computational Materials*, 7(1):12, Jan 2021.
- Ziduo Yang, Yi-Ming Zhao, Xian Wang, Xiaoqing Liu, Xiuying Zhang, Yifan Li, Qiujie Lv,
 Calvin Yu-Chian Chen, and Lei Shen. Scalable crystal structure relaxation using an iteration-free deep generative model with uncertainty quantification. March 2024.

207 A CrystalGym Background

CrystalGym [10] is an open-source RL environment that allows training of online RL policies to 208 sequentially place atoms in a given crystal structure backbone, with the aim of designing crystals to match a desired property value. The reward (obtained at the terminal state) is a distance function of the property obtained from DFT with respect to the desired value. Several RL policies have 211 been trained in this environment with PPO, SAC, Rainbow, and DQN, optimizing for properties 212 such as band gap, density, and bulk modulus, and focusing on tasks of varying difficulty. For this 213 study, rollouts with different seeds were performed using the trained policies to generate a library of 214 ~ 1000 crystals, which was used for evaluation and benchmarking. Around 60% of the RL-generated 215 compositions were novel—i.e., not found in the Materials Project. 216

217 B DFT parameters

We performed DFT single-point SCF simulations using Quantum Espresso v7.1 [9], which is fully open-source. Solid-state pseudopotentials from SSSP version 1.3.0 [14] were used for the calculations. The settings used are listed below.

- 1. calculation
- scf for single-point calculations
- vc-relax for relaxation
- 224 2. nstep: 50 (for relaxation)
- 3. ecutwfc: 50

221

222

- 4. ecutrho: 400
- 5. occupations: smearing
- 6. degauss: 0.001
- 229 7. nspin: 1
- 8. electron_maxstep: 300
- 9. mixing_mode: plain
- 232 10. mixing_beta: 0.7
- 233 11. diagonalization: david
- 12. kpoints: Chosen automatically from Kpoint density.

B.1 Formation Energy Calculation

The formation energy per atom was calculated using the total energies of the crystals and their constituent elements. The total energies of the isolated elements (88 in the action space) were calculated by performing SCF calculations on the most stable elemental crystals (i.e., 0 formation energy) present in the Materials Project. For elements that do not have a stable elemental crystal (e.g. Lu) or those that have large number of atoms in the elemental crystal (e.g. P, Se), the total energies were calculated for a single atom inside a primary cubic cell of length 10. For a crystal with N atoms, the formation energy (per atom) calculation is defined as follows.

$$E_{form} = \left(\frac{E_{tot} - \sum_{i} \frac{N_i}{n_i} E_{tot}^i}{N}\right) \text{ (eV/atom)}$$
 (1)

Here, N_i is the number of atoms of the constituent element i present in the crystal, n_i is the number of atoms (sites) of i in the elemental crystal, and E_{tot}^i is the total energy of i in the most stable elemental crystal form.

246 C Experimental Details

247 C.1 RL-Generated Crystals

- 248 The RL-generated crystals are generated in a dataset-independent manner using DFT-based reward
- 249 signals starting from existing structures. While they do not necessarily resemble materials in Materials
- 250 Project, some of the compositions (with unique structures) exist in the database.

251 C.2 Data Preprocessing and Filtering

- Our evaluation dataset initially contained 1,005 crystal structures generated through reinforcement
- 253 learning using CrystalGym. However, computing DFT forces for model-relaxed structures revealed
- that a subset of predicted structures failed DFT calculations across different models. To ensure robust
- 255 statistical comparison, we applied intersection filtering, retaining only structures that successfully
- 256 completed DFT force calculations for all evaluated models. This filtering process reduced our final
- analysis set to 831 structures, ensuring consistent evaluation across all models and metrics.

258 C.3 Relaxation Parameters

- For DeepRelax we used the default checkpoint and parameters on the model GitHub. For ASE-based models, we maintained consistent optimization parameters to ensure fair comparison:
- 1. optimizer: Fast Inertial Relaxation Engine (FIRE)
- 262 2. cell filter: Unit Cell Filter
- 3. max optimization steps: 500
- 4. fmax: 0.01

265 D Compute Resources

For all tasks we used an Nvidia Quadro RTX 8000 to ensure fair comparison of runtimes.

267 E Supplementary Results

268 E.1 Additional Plots

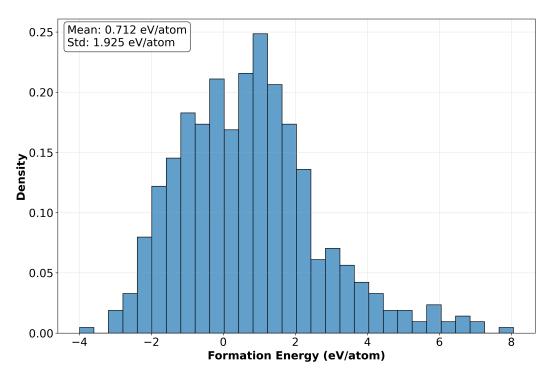
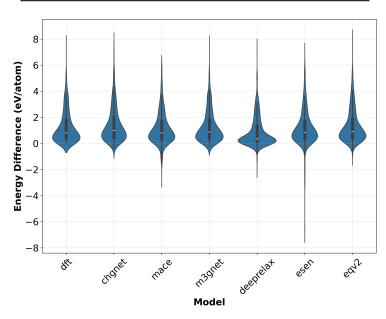


Figure 2: Histogram of unrelaxed formation energies

(a) Percentage of structures with positive-to-negative energy transitions

CHGNet	MACE	M3GNet	DeepRelax	eSEN	eqV2
28.8%	22.2%	24.0%	13.9%	25.1%	27.7%



(b) Distribution of Formation Energy Differences between Unrelaxed and Relaxed Structures

Figure 3:

Formation energy reduction achieved during relaxation: (a) Percentage of successful positive to negative formation energy changes across all structures, (b) distribution of formation energy differences between unrelaxed and relaxed structures.

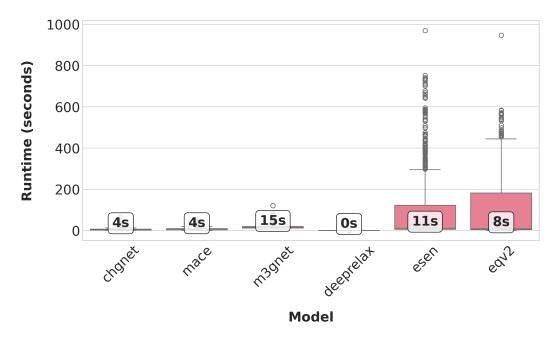


Figure 4: Runtime Box Plot

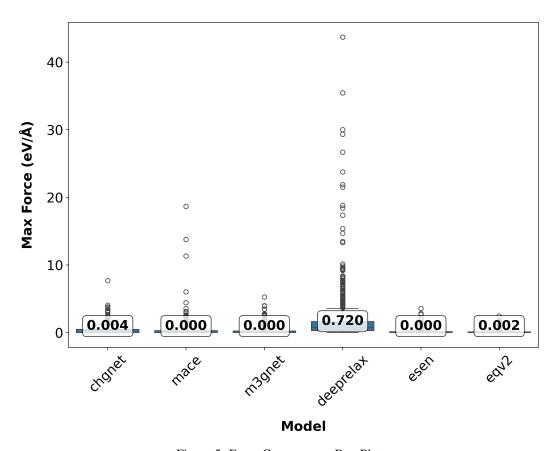


Figure 5: Force Convergence Box Plot