

---

# MLIPAudit: A benchmarking tool for Machine Learned Interatomic Potentials

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Machine-learned interatomic potentials (MLIPs) promise to significantly advance  
2 atomistic simulations by delivering quantum-level accuracy for large molecular  
3 systems at a fraction of the computational cost of traditional electronic structure  
4 methods. While model hubs and categorisation efforts have emerged in recent  
5 years, it remains difficult to consistently discover, compare, and apply these models  
6 across diverse scenarios. The field still lacks a standardised and comprehensive  
7 framework for evaluating MLIP performance. We introduce MLIPAudit, an open,  
8 curated and modular benchmarking suite designed to assess the accuracy of MLIP  
9 models across a variety of application tasks. MLIPAudit offers a diverse collection  
10 of benchmark systems, including small organic compounds, molecular liquids,  
11 proteins and flexible peptides, along with pre-computed results for a range of  
12 pre-trained and published models. MLIPAudit also provides tools for users to  
13 evaluate their models using the same standardised pipeline. A continuously updated  
14 leaderboard tracks performance across benchmarks, enabling direct comparison  
15 on downstream tasks. By providing a unified, transparent reference framework  
16 for model validation and comparison, MLIPAudit aims to foster reproducibility,  
17 transparency, and community-driven progress in the development of MLIPs for  
18 complex molecular systems. The library is available on GitHub and on PyPI 14  
19 under the Apache license 2.0.

## 20 1 Introduction

21 The accurate prediction of molecular and material properties is a cornerstone of scientific progress  
22 across disciplines, including drug discovery, functional material design, and process chemistry [1–3].  
23 Traditionally, this has been done using classical force fields, which enable efficient simulations of  
24 large systems relying on predefined functional forms and parameters derived from experiments or first-  
25 principles methods [4, 5]. Although computationally inexpensive, classical force fields often struggle  
26 to capture complex chemical interactions or generalise beyond the systems for which they were  
27 parametrised. At the other end of the spectrum, first-principles methods such as density functional  
28 theory (DFT) offer higher accuracy but at significantly greater computational cost, typically limiting  
29 their use to systems with fewer than a few hundred atoms [6, 7]. In recent years, machine-learned  
30 interatomic potentials (MLIPs) have emerged as a compelling middle ground. These models aim to  
31 retain the accuracy of first-principles methods while approaching the efficiency of classical force  
32 fields, by learning the potential energy surface directly from high-level electronic structure data  
33 [8–25].

34 Despite the rapid emergence of diverse MLIP architectures, which have significantly broadened the  
35 scope of atomistic simulations, the field continues to lack a standardised and rigorous framework for  
36 evaluating model performance in downstream applications. Many benchmarks focus on energy and  
37 force errors, which miss aspects like stability, transferability, and robustness. Recent works propose  
38 more holistic evaluations [11, 26–34], which we detail in the Literature Review section. However, all  
39 these studies highlight the need for consistent and reproducible evaluation protocols that go beyond  
40 basic error metrics, aiming to establish benchmarking practices that reflect real-world simulation  
41 demands. Therefore, a universally adopted, comprehensive benchmarking suite that can guide both  
42 model development and deployment remains an open challenge for the community.

43 To address this gap, we introduce MLIPAudit: an open, curated repository of benchmarks, reference  
44 datasets, and model evaluations for MLIP models applied (in its first version) to the analysis of small  
45 molecules, molecular liquids and biomolecules. MLIPAudit is designed to complement model-centric  
46 testing by shifting the focus to systematic validation and comparison. It provides:

- 47 • A diverse set of benchmark systems, including organic small molecules, flexible peptides,  
48 folded protein domains, molecular liquids and solvated systems.
- 49 • Pre-computed results for a range of published and pretrained MLIP models, enabling direct,  
50 fair comparisons.
- 51 • A continuously updated leaderboard, tracking performance across different tasks.
- 52 • A suite of tools for users to submit and evaluate their models within the same benchmarking  
53 pipeline. We support both Jax-based and Torch-based models, as long as they have an ASE  
54 [35, 36] calculator.

55 By providing a shared reference point for assessing accuracy, robustness, and generalisation, MLIPAu-  
56 dit aims to facilitate transparency, reproducibility, and community-wide progress in the development  
57 and deployment of MLIPs for complex molecular systems.

## 58 2 Literature Review

59 MLIP Audit aims to expand the existing methods and tools for benchmarking MLIPs. To put this  
60 work in context, we summarise current efforts for MLIP benchmarking here.

61 **Static regression metrics:** The first and most fundamental level of MLIP evaluation involves the  
62 use of standard regression metrics to quantify a model’s ability to reproduce the reference quantum-  
63 mechanical (QM) data it was trained on. The most common benchmarks in this category are the  
64 root-mean-square-error (RMSE) and mean-absolute-error (MAE) calculated for energies and atomic  
65 forces on a held-out validation dataset [37]. These benchmarks are routinely reported with the release  
66 of new MLIP models, and state-of-the-art models achieve high accuracy on these tests. Although  
67 benchmarks for atomic energies and forces are a necessary baseline for the interpolation accuracy of  
68 the models, they are insufficient to estimate their practical utility. This is demonstrated, for example,  
69 by Gonzales et al. [38], who found that three models with very similar force validation error show  
70 significant variation in performance on a structural relaxation task.

71 **Assessment of physical and chemical behaviour:** Recent MLIP benchmarks generally accompany  
72 model releases and assess performance on physical and chemical properties using QM or experimental  
73 data, typically tailored to specific use cases. For models trained on small organic molecules, standard  
74 tests include dihedral scans, conformer selection, vibrational frequencies, and interaction energies  
75 [32, 39, 40]. Biomolecular benchmarks cover backbone sampling, water properties, and folding  
76 dynamics [32, 40, 41], while models trained on reactivity data are evaluated on their ability to  
77 reproduce product, reactant, and transition state geometries, as well as reaction pathways via string or  
78 NEB methods [33, 42].

79 Comparative studies have also emerged, evaluating multiple MLIPs across diverse benchmarks. Fu et  
80 al. [27] propose a suite spanning organic molecules, peptides, and materials, and find that models

81 with low force errors may still perform poorly on simulation-based metrics like energy conservation  
82 and sampling. Similarly, Liu et al. [43] report discrepancies in atom dynamics and rare events, even  
83 for models with strong regression accuracy. These findings reflect a growing consensus that static  
84 error metrics alone are insufficient for evaluating MLIPs, and that dynamic and simulation-based  
85 benchmarks are increasingly essential.

86 **Standardised benchmarks:** While a great variety of benchmarks for accurate physical and chemical  
87 properties can be collected from individual model releases and MLIP evaluation studies, a need  
88 remains for standardised benchmarks that can be used to compare models on a level playing field and  
89 get a holistic view of their utility regarding practical tasks.

90 This gap is addressed by leaderboards and standardised frameworks. MLIP Arena [26] is a leaderboard  
91 based on a benchmark platform focused on physical awareness, stability, reactivity, and predictive  
92 power. The framework comprises a small but well-selected suite of benchmarks that address known  
93 problems like data leakage, transferability, and overreliance on specific errors. Matbench Discovery  
94 [44] features a leaderboard and evaluation framework that is easily extendable to additional models  
95 and focused exclusively on materials science. MOFSimBench [45] is a standardised benchmark  
96 specialised on metal-organic frameworks that highlights simulation metrics and bulk properties.  
97 MLIPX [46] provides a framework with a user-centric perspective, providing a set of reusable recipes  
98 that allow users to compose benchmarks for their specific tasks.

99 These standardised frameworks are valuable tools to evaluate and compare MLIP models. However,  
100 they are limited to a specific domain of application, employ a small number of benchmarks or require  
101 development by the MLIP user.

### 102 3 MLIPAudit Benchmarks

103 To enable a rigorous and meaningful evaluation of MLIP models, MLIPAudit includes a curated and  
104 modular suite of benchmarks that span a range of molecular systems and complexity levels (Figure  
105 1). These benchmarks are designed to capture both general-purpose and domain-specific challenges  
106 faced by MLIPs in industrial applications. Benchmark subsets each emphasise different aspects  
107 of model performance, such as elemental molecular dynamics stability, non-covalent interactions,  
108 conformational ranking of small organic compounds, or sampling of rotamers in biomolecules. A  
109 description of the rationale for each benchmark on the different categories is given in Appendix  
110 A, including: (i) general systems designed for molecular dynamics stability and scaling, (ii) small  
111 molecules relevant to materials chemistry, (iii) molecular liquids, and (iv) biomolecules.

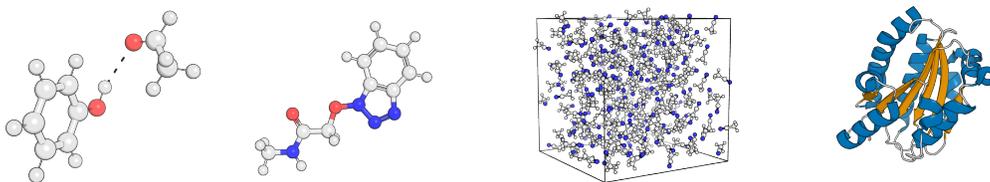


Figure 1: Representative molecular systems spanning increasing levels of structural and environmental complexity, from isolated dimers and drug-like molecules, to condensed-phase molecular liquids and folded biomolecules.

112 We have evaluated the performance of the three graph-based MLIPs provided in the open-source mlip  
113 library [25]: MACE [9], NequIP [11], and ViSNet [41]. All three models were trained on a subset of  
114 the SPICE2 dataset [47], which includes 1,737,896 molecular structures across 15 elements (B, Br, C,  
115 Cl, F, H, I, K, Li, N, Na, O, P, S, Si). From now on, MACE-SPICE2, NequIP-SPICE2 and ViSNet-  
116 SPICE2. Training protocols and dataset curation details are available in [25]. We trained versions  
117 of each of these models (MACE-t1x, NequIP-t1x, ViSNet-t1x) using 10% (randomly sampled) of

118 the original t1x dataset [48], containing a total of one million structures and four elements (H, C, N,  
 119 O). Additionally, we have trained two versions of ViSNet using different subsets of SPICE2 and 1  
 120 million datapoints from t1x (both taken from the OpenMolecules dataset - OMOL [42]), respectively,  
 121 ViSNet-SPICE2(charged)-t1x, ViSNet-SPICE2(neutral)-t1x (When not specified, the neutral version  
 122 is used). The mlipaudit library also supports Torch-based models as long as they have been  
 123 wrapper in an ASE Calculator class [35, 36]. For completeness, we have evaluated a non-exhaustive  
 124 subset of Torch-based models using their original implementation, namely: MACE-OFF [32], MACE-  
 125 MP [9], and UMA-Small [34]. Two comments on these are worth raising: (1) runtime are not optimal  
 126 for these models as they rely on ASE instead of JAXMD for simulations, (2) MACE-MP is trained  
 127 for materials and at a different level of DFT theory. It is therefore not well suited for the benchmarks  
 128 presented in MLIPAudit. We nonetheless added it as it is largely considered a reference model in the  
 129 community and as results provide some interesting insights.

130 To ensure fair and consistent comparison across models, we define a composite score  $S_m \in [0, 1]$   
 131 that averages soft-thresholded, normalised benchmark metric scores, rewarding models that approach  
 132 DFT-level accuracy. Only benchmarks compatible with a model’s element set are included, ensuring  
 133 broad applicability without penalising for unsupported systems. Though readers should note that  
 134 unless all benchmarks are completed, aggregate scores should be caveated. For full details, see  
 135 Appendix B.

136 For each benchmark, a set of test cases has been curated (Appendix C, Table 4). As public datasets  
 137 increase, it becomes increasingly challenging to ensure zero overlap between the training data and the  
 138 relevant chemistry that one needs to include to ensure the relevance and reliability of the benchmarks.  
 139 In Appendix C-Table 5, we disclose the overlap between the MLIPAudit test cases per benchmark  
 140 and the training set for the presented internal models. In most cases, the overlap is either zero or  
 141 under 10 %. But, for the conformer selection benchmark, for which two molecules (adenosine and  
 142 efivarez) from the Wiggle150 [49] dataset were present in the model’s training set. We do not provide  
 143 this information for external open source models. In the following, we will discuss the different  
 144 scores and how the overlap may impact ranking.

### 145 3.1 Overall ranking

146 Table 1 highlights the generalisation capabilities of the top-performing models. In the following,  
 147 we will analyse separately external open-source models run using the original implementation from  
 148 our internal models. Some models did not complete all benchmarks; we refer you to Appendix A,  
 149 Table 7 for more information. Missing benchmarks can be due to the availability of elements in  
 150 the training set (essentially the models trained on t1x only) or runtime issues due to the reliance of  
 151 external models on ASE [35, 36].

Table 1: Overall MLIPAudit scores

Source	Rank	Model Name	Average Score	Benchmarks
External	1	UMA-Small	0.70	12/14
External	2	MACE-OFF	0.63	11/14
External	3	MACE-MP	0.41	9/14
Internal	1	ViSNet-SPICE2	0.70	14/14
Internal	2	NequIP-SPICE2	0.70	14/14
Internal	3	ViSNet-SPICE2-t1x	0.70	14/14
Internal	4	MACE-SPICE2	0.63	14/14
Internal	4	NequIP-t1x	0.10	4/14
Internal	5	MACE-t1x	0.10	4/14
Internal	6	ViSNet-t1x	0.10	4/14

152 For the external models, UMA-Small achieves the highest average score (0.70), completing 12/14  
 153 benchmarks, followed by MACE-OFF (0.63), completing 11/14 benchmarks. MACE-MP completes  
 154 9/14 and scores 0.41; we include this model on purpose as a test for the Physics the benchmarks,  
 155 as MACE-MP is trained the MPtrj dataset [50] and therefore specialised on crystalline matter and

156 not condensed matter. All internal models completed the 14 benchmarks. ViSNet-SPICE2-t1x and  
157 ViSNet-SPICE2 attain the strongest performance (0.70), closely followed by NequIP-SPICE2 (0.68)  
158 and MACE-SPICE2 (0.63). The models specifically trained on the t1x dataset [48] score lower (0.1)  
159 and cover only a subset of benchmarks (4/14), reflecting the impact of training data breadth and  
160 domain coverage. Models consistently performing well across domains underscore the benefits of  
161 comprehensive training and robust architectures. However, it is worth noting that model performance  
162 is reflective of training strategy, not solely the model architecture, and it should not be considered an  
163 assessment of the model architecture. It is also important to note that UMA-Small, MACE-OFF, and  
164 MACE-MP may include train-test overlaps, and therefore their scores could be artificially overstated.

## 165 3.2 Categorical ranking

166 In Appendix C-Table 6, we summarise the category-based ranking analysis, which further reveals the  
167 specialisation and limitations of each MLIP model across different chemical domains. In the General  
168 category, which tests for molecular dynamics stability, most models (internal and external) achieve  
169 perfect scores, indicating strong stability for different chemical entities in vacuum and in solution  
170 4. The picture becomes more differentiated in the Small-molecule benchmarks. For the external  
171 models, UMA-Small leads with a score of 0.56, followed by MACE-OFF (0.50) and MACE-MP  
172 (0.36). The ViSNet-SPICE2-t1x variant is the best internal model in this category (0.65). Among  
173 models trained purely on SPICE2 [47], ViSNet-SPICE2, NequIP-SPICE2, and MACE-SPICE2  
174 cluster closely together (0.52-0.51), demonstrating consistent performance across gas-phase and  
175 conformational tasks. In contrast, models trained primarily on the t1x dataset [48] exhibit lower  
176 performance (0.11-0.16), consistent with the dataset’s focus on reactive gas-phase chemistry rather  
177 than diverse molecular energetics or equilibrium conformational distributions. The Molecular-liquids  
178 category shows the strongest overall spread. Within the external models, UMA-Small achieves  
179 the highest score (0.98), followed by MACE-OFF (0.73). MACE-MP, trained on inorganic crystal  
180 trajectories, underperforms here (0.45), reflecting the domain shift between crystalline materials and  
181 molecular liquids. The internal models trained on SPICE2 perform similarly with scores around  
182 0.95-0.97. These results highlight that SPICE2-trained models, despite being built from largely  
183 gas-phase and small-molecule electronic-structure data, still transfer effectively to condensed-phase  
184 structure and energetics. Performance diverges further in the Biomolecule category, which probes  
185 larger solvated, flexible, and chemically complex systems. External and Internal models (except  
186 for models trained exclusively on t1x) score very high in this category, around 0.8-1.0. However,  
187 MACE-MP also scores high (0.79), which highlights that the length of the simulation is not enough  
188 to assess the dynamical behaviour of the systems. Simulation length is constrained by computational  
189 resources, as this is the most expensive benchmark to run (more details will follow). t1x-trained  
190 models again unsurprisingly trail behind, consistent with their lack of exposure to biomolecular  
191 chemistry. Overall, these results emphasise the importance of both training data diversity and domain  
192 alignment for robust generalisation across molecular and biomolecular environments, while also  
193 pointing to meaningful architectural and training-strategy differences even within closely related  
194 model families.

## 195 3.3 Single benchmark highlighted results

### 196 3.3.1 Reactivity benchmarks

197 Internal models trained exclusively on SPICE2 (ViSNet-SPICE2, NequIP-SPICE2, MACE-SPICE2)  
198 perform notably badly in the reactivity task with scores below 0.1 (Table 2). It is worth noting  
199 that all internal models completed all test cases (100/100 for the nudge elastic band (NEB) bench-  
200 mark,  $\sim 12000/12000$  for the transition-state-theory (TST) benchmark), indicating that performance  
201 differences stem from modelling accuracy rather than lack of elements in the training set. These  
202 results suggest that, in the context of reactivity benchmarks, domain-specific training still offers a  
203 measurable edge, especially when accurate prediction of reaction energies or barriers is the primary  
204 objective. t1x trained models perform better in this category with scores ranging from 0.4-0.8 in  
205 the TST benchmark and 0.38-0.58 in the nudge-elastic-band (NEB) convergence benchmark, with

206 most notably the ViSNet-SPICE2-t1x (charged and neutral) lead this category with 0.8 and 0.58,  
 207 respectively.

Table 2: Reactivity Benchmarks Ranking

Source	Rank	Benchmark	Model Name	Score	Test Cases
External	1	Small Molecule Reactivity TST	UMA-Small	0.86	11961/11961
External	2	Small Molecule Reactivity TST	MACE-OFF	0.12	11961/11961
External	3	Small Molecule Reactivity TST	MACE-MP	0.05	11961/11961
Internal	1	Small Molecule Reactivity TST	ViSNET-SPICE2-t1x	0.77	11961/11961
Internal	2	Small Molecule Reactivity TST	NequIP-t1x	0.41	11961/11961
Internal	3	Small Molecule Reactivity TST	MACE-t1x	0.39	11961/11961
Internal	3	Small Molecule Reactivity TST	ViSNET-t1x	0.39	11961/11961
Internal	4	Small Molecule Reactivity TST	MACE-SPICE2	0.1	11961/11961
Internal	5	Small Molecule Reactivity TST	ViSNET-SPICE2	0.05	11961/11961
Internal	5	Small Molecule Reactivity TST	NequIP-SPICE2	0.05	11961/11961
Internal	1	Small Molecule Reactivity NEB	ViSNET-SPICE2-t1x	0.58	100/100
Internal	2	Small Molecule Reactivity NEB	NequIP-t1x	0.58	100/100
Internal	3	Small Molecule Reactivity NEB	MACE-t1x	0.44	100/100
Internal	3	Small Molecule Reactivity NEB	ViSNET-t1x	0.38	100/100
Internal	4	Small Molecule Reactivity NEB	MACE-SPICE2	0.1	100/100
Internal	4	Small Molecule Reactivity NEB	ViSNET-SPICE2	0.1	100/100
Internal	4	Small Molecule Reactivity NEB	NequIP-SPICE2	0.1	100/100

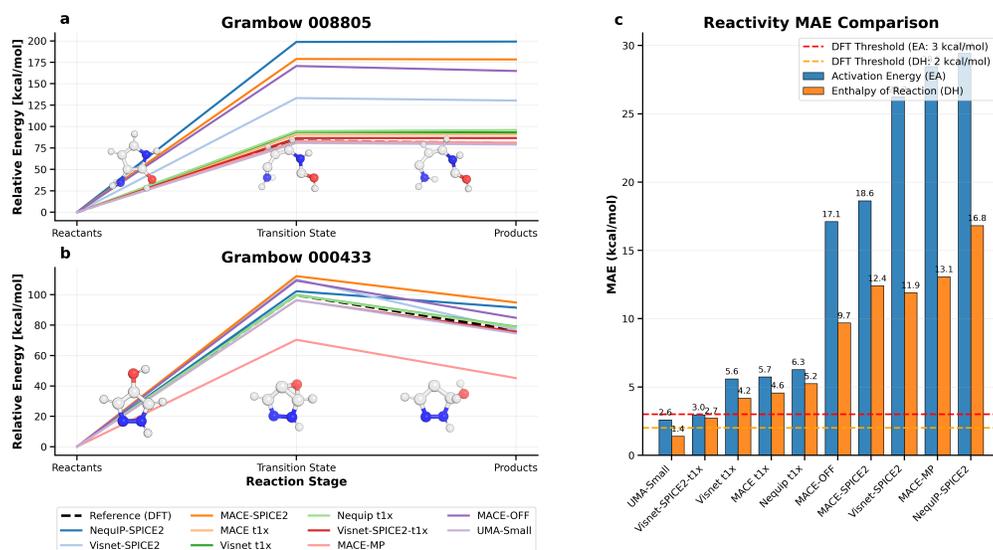


Figure 2: Reactivity benchmark performance. (a–b) Reaction energy profiles for two Grambow reactions (IDs 008805 and 000433) [51] MLIP predictions to DFT references. (c) MAEs for activation energies (EA) and reaction enthalpies across the benchmark.

208 As shown in Figure 2, all t1x-trained models outperform SPICE2 trained MLIPs (and SPICE1 in the  
 209 case of MACE-OFF), which show much larger errors, especially for activation energies.

210 From the external models, UMA-Small excels in the reactivity benchmark with a score of 0.86, with  
 211 MACE-OFF following behind with a score of 0.12. While remarkable, all our test-cases come from  
 212 the Grambow dataset [51], which is included in the t1x dataset [48], which is included in full in the  
 213 UMA-Small training data.

### 214 3.3.2 Molecular liquids benchmark: water radial distribution function

215 Having a closer look at the single benchmarks, the water radial distribution function (RDF) benchmark  
 216 provides a compelling illustration of the strengths of MLIPs over traditional force fields. As shown in

217 Figure 3, all five internal MLIP models, MACE-SPICE2, ViSNet-SPICE2, ViSNet-SPICE2(neutral)-  
 218 t1x, ViSNet-SPICE2(charged)-t1x, ViSNet-SPICE2 and NequIP-SPICE2, reproduce the experimental  
 219 RDF profile with high fidelity across the full radial range, accurately reproducing both the first  
 220 solvation shell peak and subsequent oscillations. And this is also true for the original implementations  
 221 of UMA-Small and MACE-OFF. In contrast, TIP3P and TIP4P [52], two of the most widely used  
 222 classical water models, show notable deviations, particularly in the overstructured and exaggerated  
 height of the first peak, a known artefact in rigid water models [53]. Notably, MACE-MP produces

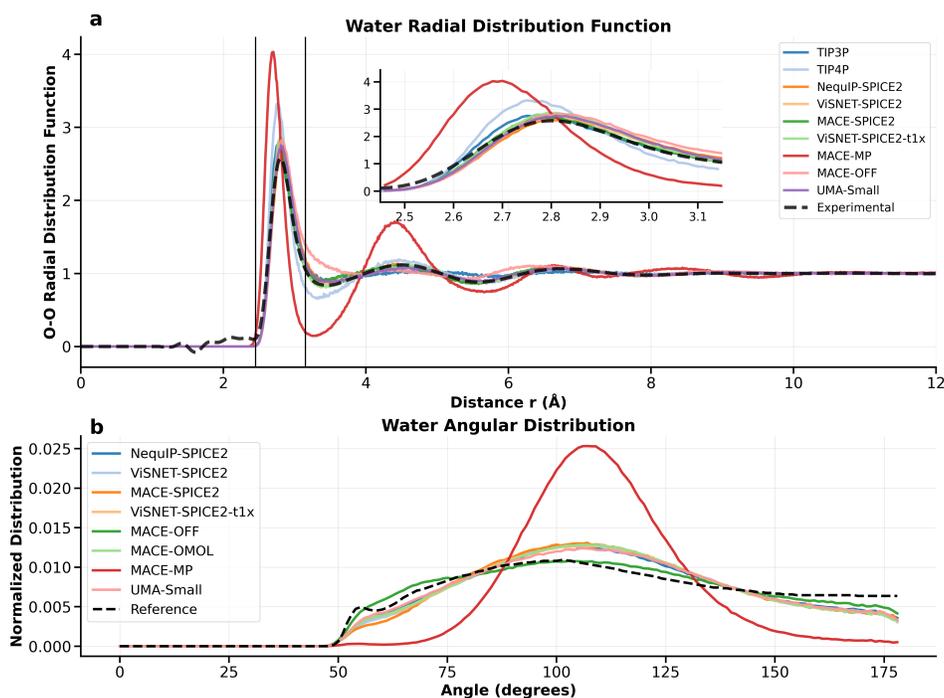


Figure 3: Water radial distribution function and angular distribution for the example models, compared with the experimental observable and two water classical forcefields TIP3P and TIP4P [52]

223 crystalline water even when simulated at 300 K, indicating that the model remains strongly biased  
 224 toward its crystal-structure training data despite liquid-phase simulation conditions. This behaviour is  
 225 evident in the radial distribution function (RDF): whereas liquid water shows a broadened first O–O  
 226 peak near 2.8 Å and damped oscillations characteristic of short-range order, crystalline (ice-like)  
 227 water exhibits sharp, well-defined peaks extending to long range, reflecting persistent translational  
 228 order. These qualitative differences are well-established in the literature [54].

230 This alignment between MLIP predictions and experimental data strongly supports the notion that  
 231 learned potentials, trained on accurate quantum data, can capture the subtle balance of hydrogen  
 232 bonding and thermal fluctuations that define liquid water structure, without the need for hand-tuned  
 233 parametrisation. This not only reflects the higher representational capacity of MLIPs but also  
 234 demonstrates their ability to generalise to bulk-phase properties, a capability that classical force fields  
 235 struggle to match without introducing complex polarisable terms or many-body corrections.

### 236 3.3.3 Small molecules benchmarks: dihedral scans

237 The dihedral scan benchmark highlights another area where MLIP models show outstanding agree-  
 238 ment with quantum reference data. As shown in Figure 4, the energy profiles predicted by all MLIP  
 239 models align nearly perfectly with DFT-calculated torsional energy curves across a representative  
 240 scan. This agreement is not only qualitative—preserving the positions and heights of barriers, but  
 241 also quantitatively precise, with RMSE values all well below the 1.0 kcal/mol DFT-level convergence  
 242 threshold. This strong performance is further reflected in the ranking table (Appendix C, Table 6),

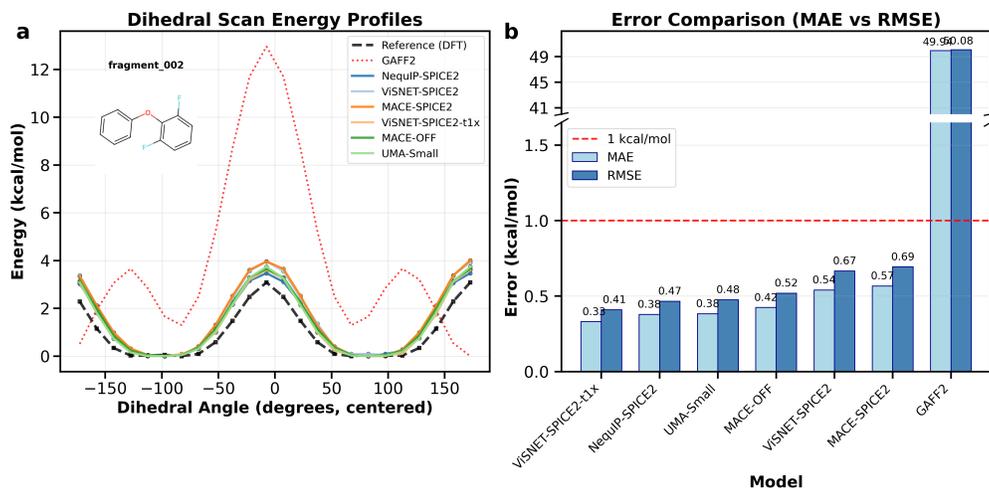


Figure 4: Dihedral scan benchmark. (a) Dihedral energy profiles for fragment 015 compared to DFT reference values. (b) MAE and RMSE for each model. DFT-level error threshold (red dashed line).

243 where ViSNet-SPICE2 and ViSNet-SPICE2-t1x lead the benchmark scoring  $\sim 1.0$ , followed closely  
 244 by NequIP-SPICE2 and MACE-SPICE2, MACE-SPICE2-t1x. Notably, all models completed the  
 245 full set of 500 fragments, demonstrating not only accuracy but robustness and generalisability across  
 246 a diverse chemical space.

247 The error bars shown on the right panel of Figure 4 underscore how consistent the models are,  
 248 with MAE values under 0.12 kcal/mol for all methods—well within chemical accuracy. MLIPs  
 249 outperform classical parameters like GAFF2 [55]. These results validate the capability of current  
 250 MLIPs to accurately model intramolecular potential energy surfaces, a critical requirement for reliable  
 251 conformational sampling, molecular docking, or pharmacophore prediction.

252 Taken together, this benchmark provides a clear example of how MLIPs can match DFT accuracy at a  
 253 fraction of the computational cost, making them practical for high-throughput screening or molecular  
 254 simulations involving flexible, drug-like molecules.

### 255 3.3.4 Small molecules benchmarks: conformer ranking

256 Figure 6 presents model performance on the conformer benchmark, showing MAE values by molecule  
 257 for three general-purpose MLIPs: NequIP-SPICE2, ViSNet-SPICE2, and MACE-SPICE2. All models  
 258 were trained on datasets that included adenosine (ADO) and efavirenz (EFA), while benzylpenicillin  
 259 (BPN) was excluded from training and thus acts as a stronger generalisation test.

260 Despite having seen ADO and EFA during training, none of the models reach the DFT-level MAE  
 261 threshold of 0.5 kcal/mol, pointing to persistent difficulty in accurately ranking conformers. ADO is  
 262 best predicted, while EFA shows higher errors due to its flexibility. BPN, which was unseen during  
 263 training, is the most challenging, though MACE-SPICE2 shows slightly better generalisation. All  
 264 models outperform GAFF2 [55], especially on EFA. Still, as seen in Appendix C, Figure 7, predicted  
 265 vs. DFT energy plots show strong agreement and near-perfect Spearman correlations across all  
 266 molecules.

267 This consistency suggests that while the models may struggle to reproduce exact conformer energy  
 268 magnitudes (as seen in the MAE analysis), they are highly effective at preserving the correct energetic  
 269 ordering. In practical applications like conformer selection or ranking, such ordinal accuracy can  
 270 be more important than precise energetic reproduction, particularly when used in combination with  
 271 scoring functions or downstream screening.

272 Interestingly, the performance gap between in-training-set molecules (ADO, EFA) and the out-of-  
 273 distribution case (BPN) is far less pronounced here than in absolute MAE terms—highlighting that

274 model generalisation, at least in terms of correlation, is relatively robust. These findings reinforce  
275 the importance of using multiple complementary metrics (e.g., MAE and rank correlation) when  
276 evaluating MLIP performance for conformational energetics.

### 277 3.3.5 Biomolecules benchmarks

278 The biomolecules benchmark (Appendix C, Table 6) provides a fitting conclusion to our compre-  
279 hensive assessment, highlighting the potential for MLIP models to operate effectively in complex,  
280 biologically relevant regimes. The biomolecules benchmark is the most computationally intensive  
281 one, as it involves solvated systems with 1000 to 4000 atoms in total (Appendix C, Table 4).

282 All top models successfully completed the protein folding stability benchmark (3/3 test cases, see  
283 Appendix C), all models achieve similar scores  $\sim 0.525$ , but there is room for improvement. This level  
284 of agreement underscores the growing maturity of MLIPs for macromolecular tasks. The Protein  
285 Sampling benchmark across different MLIP models shows that models trained on the SPICE2 dataset  
286 (e.g., ViSNet-SPICE2, NequIP-SPICE2, MACE-SPICE2) significantly outperform their t1x-trained  
287 counterparts, with ViSNet-SPICE2 achieving the highest score (0.928) and full coverage (12/12  
288 systems). Taken together, the results from this and all previous benchmarks reinforce a central  
289 conclusion: while task-specific training offers advantages in specialised domains, the leading models  
290 demonstrate strong, transferable performance across molecular scales and properties, setting the stage  
291 for robust deployment in real-world chemistry and biology applications.

## 292 3.4 Conclusions and future outlook

293 The MLIPAudit suite provides a comprehensive and diverse evaluation framework for MLIPs,  
294 spanning small-molecule geometrical and conformational energetics, reactivity, molecular liquids,  
295 and biomolecular stability and sampling. Results show that while specialised models trained on the  
296 t1x dataset excel in targeted tasks such as reaction barrier prediction, general-purpose architectures  
297 like ViSNet-SPICE2, NequIP-SPICE2, and MACE-SPICE2 exhibit strong and transferable accuracy  
298 across a wide range of benchmarks, often surpassing classical force fields and closely matching DFT  
299 reference data in others. Notably, the ViSNet model trained on SPICE2 and t1x from the OMOL  
300 dataset leads the small-molecule benchmarks, highlighting the promise of hybrid training strategies  
301 and possibly reflecting the importance of the underlying level of theory used in data generation.

302 Despite this progress, performance gaps persist, especially in condensed-phase systems and energeti-  
303 cally subtle regimes, indicating that further improvements are needed. While MLIPAudit establishes  
304 a unified and reproducible evaluation suite, it also has limitations. The current set of models is  
305 biased toward graph neural network architectures, and the benchmarks rely primarily on DFT data  
306 of varying origin, which may introduce systematic bias. Efficiency and robustness-oriented metrics  
307 (e.g., uncertainty calibration and scalability) are not yet fully assessed, and several critical chemical  
308 regimes, such as transition-metal systems, enzyme catalysis, and extreme thermodynamic conditions,  
309 remain under-represented due to limited reference data.

310 A further challenge lies in maintaining truly blind test sets. As the community continually expands  
311 training datasets, ensuring that future benchmark systems remain unseen becomes increasingly  
312 difficult. In future iterations, we will explore generating dedicated blind datasets and curated QM  
313 reference sets, though this task will remain increasingly complex.

314 Future releases will introduce more demanding simulation tasks, such as free-energy estimation,  
315 reactive condensed-phase processes, and protein–ligand systems. By evolving alongside the MLIP  
316 community and enabling continuous contribution, MLIPAudit aims not only to benchmark progress  
317 but to support rigorous, open, and scalable development of next-generation ML interatomic potentials.  
318 By continually broadening the scope and complexity of MLIPAudit, we hope to accelerate the  
319 development of MLIPs that are not only accurate but also general, scalable, and ready for real-world  
320 deployment across the chemical sciences.

## 321 References

- 322 [1] Matthias Rupp, Alexandre Tkatchenko, Klaus-Robert Müller, and O. Anatole von Lilienfeld.  
323 Fast and accurate modeling of molecular atomization energies with machine learning. *Physical*  
324 *Review Letters*, 108(5):058301, 2012.
- 325 [2] Keith T Butler, Daniel W Davies, Hugh Cartwright, Olexandr Isayev, and Aron Walsh. Machine  
326 learning for molecular and materials science. *Nature*, 559(7715):547–555, 2018.
- 327 [3] Gerbrand Ceder and Kristin Persson. The stuff of dreams. *Scientific American*, 309(3):36–41,  
328 2013.
- 329 [4] William D Cornell, Piotr Cieplak, Christopher I Bayly, Ian R Gould, Kenneth M Merz, David M  
330 Ferguson, et al. A second generation force field for the simulation of proteins, nucleic acids,  
331 and organic molecules. *Journal of the American Chemical Society*, 117(19):5179–5197, 1995.
- 332 [5] Alexander D MacKerell Jr, Donald Bashford, Michael Bellott, Roland L Dunbrack Jr, John D  
333 Evanseck, Michael J Field, et al. All-atom empirical potential for molecular modeling and  
334 dynamics studies of proteins. *The Journal of Physical Chemistry B*, 102(18):3586–3616, 1998.
- 335 [6] Walter Kohn and Lu Jeu Sham. Self-consistent equations including exchange and correlation  
336 effects. *Physical Review*, 140(4A):A1133, 1965.
- 337 [7] Robert G Parr and Weitao Yang. *Density-functional theory of atoms and molecules*. Oxford  
338 University Press, 1989.
- 339 [8] Yury Lysovorskiy, Chris Van Den Oord, Alexey Bochkarev, Shyue Ping Menon, Matteo Rinaldi,  
340 Tobias Hammerschmidt, Michael Mrovec, Alexander Thompson, Gábor Csányi, Christoph  
341 Ortner, et al. Performant implementation of the atomic cluster expansion (pace) and application  
342 to copper and silicon. *npj Computational Materials*, 7:97, 2021.
- 343 [9] Ilyes Batatia, David P Kovacs, Gregor Simm, Christoph Ortner, and Gábor Csányi. Mace:  
344 Higher-order equivariant message passing neural networks for fast and accurate force fields.  
345 *Advances in Neural Information Processing Systems*, 35, 2022.
- 346 [10] Dávid Péter Kovács, Ilyes Batatia, Eszter Sara Arany, and Gabor Csanyi. Evaluation of the  
347 mace force field architecture: From medicinal chemistry to materials science. *The Journal of*  
348 *Chemical Physics*, 159(4):044118, 2023.
- 349 [11] Sebastian Batzner, Alexander Musaelian, Linfeng Sun, Michael Geiger, Jonathan P Mailoa,  
350 Marc Kornbluth, Nicola Molinari, Tyle Smidt, and Boris Kozinsky. E(3)-equivariant graph  
351 neural networks for data-efficient and accurate interatomic potentials. *Nature Communications*,  
352 13:2453, 2022.
- 353 [12] Vitalii Zaverkin, Daniel Holzmüller, Luca Bonferraro, and Johannes Kästner. Transfer learning  
354 for chemically accurate interatomic neural network potentials. *Physical Chemistry Chemical*  
355 *Physics*, 25:5383, 2023.
- 356 [13] Mehdi Haghighatlari, Jia Li, Xiangyu Guan, Oliver Zhang, Abhishek Das, Christoph J Stein,  
357 Fatemeh Heidar-Zadeh, Meng Liu, Martin Head-Gordon, Lucas Bertels, et al. Newtonnet: A  
358 newtonian message passing network for deep learning of interatomic potentials and forces.  
359 *Digital Discovery*, 1:333, 2022.
- 360 [14] Alexander V Shapeev. Moment tensor potentials: A class of systematically improvable inter-  
361 atomic potentials. *Multiscale Modeling & Simulation*, 14(3):1153–1173, 2016.
- 362 [15] David Anstine, Roman Zubatyuk, and Olexandr Isayev. Aimnet2: A neural network potential to  
363 meet your neutral, charged, organic, and elemental-organic needs. *Chemical Science*, 2025.

- 364 [16] A Kabylda, V Vassilev-Galindo, S Chmiela, I Poltavsky, and Alexandre Tkatchenko. Efficient  
365 interatomic descriptors for accurate machine learning force fields of extended molecules. *Nature*  
366 *Communications*, 14:3562, 2023.
- 367 [17] Jörg Behler. Four generations of high-dimensional neural network potentials. *Chemical Reviews*,  
368 121(16):10037–10072, 2021.
- 369 [18] Federico Musil, Andrea Grisafi, Albert P Bartók, Christoph Ortner, Gábor Csányi, and Michele  
370 Ceriotti. Physics-inspired structural representations for molecules and materials. *Chemical*  
371 *Reviews*, 121(16):9759–9815, 2021.
- 372 [19] Volker L Deringer, Albert P Bartók, Noam Bernstein, Daniel M Wilkins, Michele Ceriotti, and  
373 Gábor Csányi. Gaussian process regression for materials and molecules. *Chemical Reviews*,  
374 121(16):10073–10141, 2021.
- 375 [20] Bing Huang and O Anatole Von Lilienfeld. Ab initio machine learning in chemical compound  
376 space. *Chemical Reviews*, 121(16):10001–10036, 2021.
- 377 [21] Murray S Daw and MI Baskes. Embedded-atom method: Derivation and application to  
378 impurities, surfaces, and other defects in metals. *Physical Review B*, 29(12):6443, 1984.
- 379 [22] Volker L Deringer, Noam Bernstein, Gábor Csányi, C Ben Mahmoud, Michele Ceriotti, Mark  
380 Wilson, David A Drabold, and Steven R Elliott. Origins of structural and electronic transitions  
381 in disordered silicon. *Nature*, 589:59–64, 2021.
- 382 [23] William J Baldwin, Xiaoxuan Liang, Johan Klarbring, Marija Dubajic, Diego Dell’Angelo,  
383 Charles Sutton, Chiara Caddeo, Samuel D Stranks, Alessandro Mattoni, Aron Walsh, et al.  
384 Dynamic local structure in caesium lead iodide: Spatial correlation and transient domains.  
385 *Small*, 20(2303565), 2023.
- 386 [24] Christopher W Rosenbrock, Konstantin Gubaev, Alexander V Shapeev, László B Pártay, Noam  
387 Bernstein, Gábor Csányi, and Gus L W Hart. Machine-learned interatomic potentials for alloys  
388 and alloy phase diagrams. *npj Computational Materials*, 7:24, 2021.
- 389 [25] Christoph Brunken, Olivier Peltre, Heloise Chomet, Lucien Walewski, Manus McAuliffe,  
390 Valentin Heyraud, Solal Attias, Martin Maarand, Yessine Khanfir, Edan Toledo, Fabio Falcioni,  
391 Marie Bluntzer, Silvia Acosta-Gutiérrez, and Jules Tilly. Machine learning interatomic poten-  
392 tials: library for efficient training, model development and simulation of molecular systems.  
393 *arXiv preprint*, 2025.
- 394 [26] Yuan Chiang, Tobias Kreiman, Elizabeth Weaver, Ishan Amin, Matthew Kuner, Christine Zhang,  
395 Aaron Kaplan, Daryl Chrzan, Samuel Blau, Aditi Krishnapriyan, and Mark Asta. MLIP Arena:  
396 Fair and transparent benchmark of machine-learned interatomic potentials. *AI4Mat ICLR*, 2025.
- 397 [27] Xiang Fu, Zhenghao Wu, Wujie Wang, Tian Xie, Sinan Keten, Rafael Gomez-Bombarelli,  
398 and Tommi Jaakkola. Forces are not Enough: Benchmark and Critical Evaluation for Ma-  
399 chine Learning Force Fields with Molecular Simulations. *Transactions on Machine Learning*  
400 *Research*, 4, 2023.
- 401 [28] Tristan Maxson, Ademola Soyemi, Benjamin W. J. Chen, and Tibor Szilvási. Enhancing  
402 the quality and reliability of machine learning interatomic potentials through better reporting  
403 practices. *The Journal of Physical Chemistry C*, 2024.
- 404 [29] Christoph Ortner and Yangshuai Wang. A framework for a generalisation analysis of  
405 machine-learned interatomic potentials. *arXiv preprint*, 2022.
- 406 [30] Michael J. Waters and James M. Rondinelli. Benchmarking structural evolution methods for  
407 training of machine learned interatomic potentials. *arXiv preprint*, 2022.

- 408 [31] Yunxing Zuo, Chi Chen, Xiangguo Li, Zhi Deng, Yiming Chen, Jörg Behler, Gábor Csányi,  
409 Alexander V. Shapeev, Aidan P. Thompson, Mitchell A. Wood, and Shyue Ping Ong. A  
410 performance and cost assessment of machine learning interatomic potentials. *arXiv preprint*,  
411 2019.
- 412 [32] Dávid Péter Kovács, J. Harry Moore, Nicholas J. Browning, Ilyes Batatia, Joshua T. Horton,  
413 Yixuan Pu, Venkat Kapil, William C. Witt, Ioan-Bogdan Magdău, Daniel J. Cole, and Gábor  
414 Csányi. MACE-OFF: Transferable Short Range Machine Learning Force Fields for Organic  
415 Molecules. *Journal of the American Chemical Society*, 2025.
- 416 [33] Dylan M. Anstine, Qiyuan Zhao, Roman Zubatiuk, Shuhao Zhang, Veerupaksh Singla, Philipp  
417 Nikitin, Brett M. Savoie, and Olexandr Isayev. AIMNet2-rxn: A Machine Learned Potential for  
418 Generalized Reaction Modeling on a Millions-of-Pathways Scale. *ChemRxiv preprint*, 2025.
- 419 [34] Brandon M. Wood, Misko Dzamba, Xiang Fu, Meng Gao, Muhammed Shuaibi, Luis Barroso-  
420 Luque, Kareem Abdelmaqsoud, Vahe Gharakhanyan, John R. Kitchin, Daniel S. Levine, Kyle  
421 Michel, Anuroop Sriram, Taco Cohen, Abhishek Das, Ammar Rizvi, Sushree Jagriti Sahoo,  
422 Zachary W. Ulissi, and C. Lawrence Zitnick. Uma: A family of universal models for atoms.  
423 *arXiv preprint*, 2025.
- 424 [35] Ask Hjorth Larsen, Jens Jørgen Mortensen, Jakob Blomqvist, Ivano E Castelli, Rune Christensen,  
425 Marcin Dułak, Jesper Friis, Michael N Groves, Bjørk Hammer, Cory Hargus, Eric D Hermes,  
426 Paul C Jennings, Peter Bjerre Jensen, James Kermode, John R Kitchin, Esben Leonhard  
427 Kolsbjerg, Joseph Kubal, Kristen Kaasbjerg, Steen Lysgaard, Jón Bergmann Maronsson, Tristan  
428 Maxson, Thomas Olsen, Lars Pastewka, Andrew Peterson, Carsten Rostgaard, Jakob Schiøtz,  
429 Ole Schütt, Mikkel Strange, Kristian S Thygesen, Tejs Vegge, Lasse Vilhelmsen, Michael  
430 Walter, Zhenhua Zeng, and Karsten W Jacobsen. The atomic simulation environment—a python  
431 library for working with atoms. *Journal of Physics: Condensed Matter*, 29(27):273002, 2017.  
432 URL <http://stacks.iop.org/0953-8984/29/i=27/a=273002>.
- 433 [36] S. R. Bahn and K. W. Jacobsen. An object-oriented scripting interface to a legacy electronic  
434 structure code. *Comput. Sci. Eng.*, 4(3):56–66, MAY-JUN 2002. ISSN 1521-9615. doi:  
435 10.1109/5992.998641.
- 436 [37] Joe D. Morrow, John L. A. Gardner, and Volker L. Deringer. How to validate machine-learned  
437 interatomic potentials. *The Journal of Chemical Physics*, 158:121501, 2023.
- 438 [38] Carmelo Gonzales, Eric Fuemmeler, Ellad Tadmor, Stefano Martiniani, and Santiago Miret.  
439 Benchmarking of Universal Machine Learning Interatomic Potentials for Structural Relaxation.  
440 *AI4Mat NeurIPS*, 2024.
- 441 [39] Anders S. Christensen, Sai Krishna Sirumalla, Zhuoran Qiao, Michael B. O’Connor, Daniel  
442 G. A. Smith, Feizhi Ding, Peter J. Bygrave, Animashree Anandkumar, Matthew Welborn,  
443 Frederick R. Manby, and Thomas F. Miller. OrbNet Denali: A machine learning potential for  
444 biological and organic chemistry with semi-empirical cost and DFT accuracy. *The Journal of*  
445 *Chemical Physics*, 155:204103, 2021.
- 446 [40] John L. Weber, Rishabh D. Guha, Garvit Agarwal, Yujing Wei, Aidan A. Fike, Xiaowei Xie,  
447 James Stevenson, Karl Leswing, Mathew D. Halls, Robert Abel, and Leif D. Jacobson. Efficient  
448 Long-Range Machine Learning Force Fields for Liquid and Materials Properties. *arXiv preprint*,  
449 2025.
- 450 [41] Yusong Wang, Tong Wang, Shaoning Li, Xinheng He, Mingyu Li, Zun Wang, Nanning Zheng,  
451 Bin Shao, and Tie-Yan Liu. Enhancing geometric representations for molecules with equivariant  
452 vector-scalar interactive message passing. *Nature Communications*, 15(313), 2024.
- 453 [42] Daniel S. Levine, Muhammed Shuaibi, Evan Walter Clark Spotte-Smith, Michael G. Taylor,  
454 Muhammad R. Hasyim, Kyle Michel, Ilyes Batatia, Gábor Csányi, Misko Dzamba, Peter

- 455 Eastman, Nathan C. Frey, Xiang Fu, Vahe Gharakhanyan, Aditi S. Krishnapriyan, Joshua A.  
456 Rackers, Sanjeev Raja, Ammar Rizvi, Andrew S. Rosen, Zachary Ulissi, Santiago Vargas,  
457 C. Lawrence Zitnick, Samuel M. Blau, and Brandon M. Wood. The Open Molecules 2025  
458 (OMol25) Dataset, Evaluations, and Models. *arXiv preprint*, 2025.
- 459 [43] Yunsheng Liu, Xingfeng He, and Yifei Mo. Discrepancies and error evaluation metrics for  
460 machine learning interatomic potentials. *npj Computational Materials*, 9(174), 2023.
- 461 [44] Janosh Riebesell, Rhys E. A. Goodall, Philipp Benner, Yuan Chiang, Bowen Deng, Gerbrand  
462 Ceder, Mark Asta, Alpha A. Lee, Anubhav Jain, and Kristin A. Persson. Matbench Discovery –  
463 A framework to evaluate machine learning crystal stability predictions. *arXiv preprint*, 2024.
- 464 [45] Hendrik Kraß, Ju Huang, and Seyed Mohamad Moosavi. MOFSimBench: Evaluating Universal  
465 Machine Learning Interatomic Potentials In Metal–Organic Framework Molecular Modeling.  
466 *arXiv preprint*, 2025.
- 467 [46] Fabian Zills, Sheena Agarwal, Tiago Goncalves, Srishti Gupta, Edvin Fako, Shuang Han, Imke  
468 Mueller, Christian Holm, and Sandip De. MLIPX: Machine Learned Interatomic Potential  
469 eXploration. *ChemRxiv preprint*, 2025.
- 470 [47] Peter Eastman, Pavan Kumar Behara, David L Dotson, Raimondas Galvelis, John E Herr, Josh T  
471 Horton, Yuezhi Mao, John D Chodera, Benjamin P Pritchard, Yuanqing Wang, et al. Spice, a  
472 dataset of drug-like molecules and peptides for training machine learning potentials. *Scientific  
473 Data*, 10(11), 2023.
- 474 [48] Mathias Schreiner, Arghya Bhowmik, Tejs Vegge, Jonas Busk, and Ole Winther. Transition1x  
475 - a dataset for building generalizable reactive machine learning potentials. *Scientific Data*, 9  
476 (779), 2022.
- 477 [49] Rebecca Brew, Ian Nelson, Meruyert Binayeva, Amlan Nayak, Wyatt Simmons, Joseph Gair,  
478 and Corin Wagen. Wiggle150: Benchmarking density functionals and neural network potentials  
479 on highly strained conformers. *ChemRxiv preprint*, 2025.
- 480 [50] Bowen Deng, Peichen Zhong, KyuJung Jun, Janosh Riebesell, Kevin Han, Christopher J. Bartel,  
481 and Gerbrand Ceder. Chgnet as a pretrained universal neural network potential for charge-  
482 informed atomistic modelling. *Nature Machine Intelligence*, 5(9):1031–1041, September 2023.  
483 ISSN 2522-5839. doi: 10.1038/s42256-023-00716-3. URL [http://dx.doi.org/10.1038/  
484 s42256-023-00716-3](http://dx.doi.org/10.1038/s42256-023-00716-3).
- 485 [51] L. Pattanaik Grambow and W. H. Green. Reactants, products, and transition states of elementary  
486 chemical reactions based on quantum chemistry. *Scientific Data*, 7(137), 2020.
- 487 [52] W. L. Jorgensen, J. Chandrasekhar, J. D. Madura, R. W. Impey, and M. L. Klein. Comparison  
488 of simple potential functions for simulating liquid water. *Journal of Chemical Physics*, 79:  
489 926–935, 1983.
- 490 [53] Gaia Camisasca, Harshad Pathak, Kjartan Thor Wikfeldt, and Lars G. M. Pettersson. Radial  
491 distribution functions of water: Models vs experiments. *The Journal of Chemical Physics*, 151:  
492 044502, 2019.
- 493 [54] R. E. Skyner, J. B. O. Mitchell, and C. R. Groom. Probing the average distribution of water in  
494 organic hydrate crystal structures with radial distribution functions (rdfs). *CrystEngComm*, 19:  
495 641–652, 2017. doi: 10.1039/C6CE02119K.
- 496 [55] Xibing He, Viet H. Man, Wei Yang, Tai-Sung Lee, and Junmei Wang. A fast and high-quality  
497 charge model for the next generation general amber force field. *The Journal of Chemical  
498 Physics*, 153:114502, 2020.

- 499 [56] Robert T. McGibbon, Kyle A. Beauchamp, Matthew P. Harrigan, Christoph Klein, Jason M.  
500 Swails, Carlos X. Hernández, Christian R. Schwantes, Lee-Ping Wang, Thomas J. Lane, and  
501 Vijay S. Pande. Mdtraj: A modern open library for the analysis of molecular dynamics  
502 trajectories. *Biophysical Journal*, 109(8):1528 – 1532, 2015. doi: 10.1016/j.bpj.2015.08.015.
- 503 [57] Skolnick J. Zhang Y. Scoring function for automated assessment of protein structure template  
504 quality. *Proteins.*, 57(4):702–710, 2004.
- 505 [58] Sander C. Kabsch W. Dictionary of protein secondary structure: pattern recognition of hydrogen-  
506 bonded networks in three-dimensional structures. *Biopolymers*, 22(12):2577–637, 1983.
- 507 [59] S.C. Lovell, J.M. Word, J.S. Richardson, and D.C. Richardson. The penultimate rotamer library.  
508 *Proteins*, 40:389–408, 2000.
- 509 [60] David M. Mount Songrit Maneewongvatana. Analysis of approximate nearest neighbor search-  
510 ing with clustered point sets. *ArXiv*, 1999. doi: <https://doi.org/10.48550/arXiv.cs/9901013>.
- 511 [61] Pauli Virtanen, Ralf Gommers, Travis E. Oliphant, Matt Haberland, Tyler Reddy, David  
512 Cournapeau, Evgeni Burovski, Pearu Peterson, Warren Weckesser, Jonathan Bright, Stéfan J.  
513 van der Walt, Matthew Brett, Joshua Wilson, K. Jarrod Millman, Nikolay Mayorov, Andrew  
514 R. J. Nelson, Eric Jones, Robert Kern, Eric Larson, C J Carey, İlhan Polat, Yu Feng, Eric W.  
515 Moore, Jake VanderPlas, Denis Laxalde, Josef Perktold, Robert Cimrman, Ian Henriksen, E. A.  
516 Quintero, Charles R. Harris, Anne M. Archibald, Antônio H. Ribeiro, Fabian Pedregosa, Paul  
517 van Mulbregt, and SciPy 1.0 Contributors. SciPy 1.0: Fundamental Algorithms for Scientific  
518 Computing in Python. *Nature Methods*, 17:261–272, 2020. doi: 10.1038/s41592-019-0686-2.
- 519 [62] Narbe Mardirossian and Martin Head-Gordon. wb97m-v: A combinatorially optimized, range-  
520 separated hybrid, meta-gga density functional with vv10 nonlocal correlation. *The Journal of*  
521 *Chemical Physics*, 144(21):214110, 06 2016. ISSN 0021-9606. doi: 10.1063/1.4952647. URL  
522 <https://doi.org/10.1063/1.4952647>.
- 523 [63] Simon Boothroyd, Pavan Kumar Behara, Owen C. Madin, David F. Hahn, Hyesu Jang, Vytautas  
524 Gapsys, Jeffrey R. Wagner, Joshua T. Horton, David L. Dotson, Matthew W. Thompson,  
525 Jessica Maat, Trevor Gokey, Lee-Ping Wang, Daniel J. Cole, Michael K. Gilson, John D.  
526 Chodera, Christopher I. Bayly, Michael R. Shirts, and David L. Mobley. Development and  
527 benchmarking of open force field 2.0.0: The sage small molecule force field. *Journal of*  
528 *Chemical Theory and Computation*, 19(11):3251–3275, 2023. doi: 10.1021/acs.jctc.3c00039.  
529 URL <https://doi.org/10.1021/acs.jctc.3c00039>.
- 530 [64] J. S. Smith, O. Isayev, and A. E. Roitberg. Ani-1: an extensible neural network potential  
531 with dft accuracy at force field computational cost. *Chem. Sci.*, 8:3192–3203, 2017. doi:  
532 10.1039/C6SC05720A. URL <http://dx.doi.org/10.1039/C6SC05720A>.
- 533 [65] Philip R. Evans. An introduction to stereochemical restraints. *Acta Crystallographica Section D:*  
534 *Biological Crystallography*, 63(Pt 1):58–61, January 2007. doi: 10.1107/S090744490604604X.  
535 URL <https://doi.org/10.1107/S090744490604604X>. Epub 2006 Dec 13.
- 536 [66] Markus Bursch, Jan-Michael Mewes, Andreas Hansen, and Stefan Grimme. Best-practice dft  
537 protocols for basic molecular computational chemistry. *Angewandte Chemie International*  
538 *Edition*, 61(42):e202205735, 2022. doi: <https://doi.org/10.1002/anie.202205735>. URL <https://onlinelibrary.wiley.com/doi/abs/10.1002/anie.202205735>.
- 540 [67] Frank Neese et al. *Section 4.5: Nudged Elastic Band Method*. Max-Planck-  
541 Institut für Kohlenforschung, Mülheim an der Ruhr, Germany, 2024. URL  
542 [https://orca-manual.mpi-muelheim.mpg.de/contents/structureactivity/](https://orca-manual.mpi-muelheim.mpg.de/contents/structureactivity/neb.html)  
543 [neb.html](https://orca-manual.mpi-muelheim.mpg.de/contents/structureactivity/neb.html). Accessed: 2025-10-31; available from [https://orca-manual.mpi-](https://orca-manual.mpi-muelheim.mpg.de/contents/structureactivity/neb.html)  
544 [muelheim.mpg.de/contents/structureactivity/neb.html](https://orca-manual.mpi-muelheim.mpg.de/contents/structureactivity/neb.html).

- 545 [68] Tobias Morawietz, Andreas Singraber, Christoph Dellago, and Jörg Behler. How van der  
546 waals interactions determine the unique properties of water. *Proceedings of the National*  
547 *Academy of Sciences*, 113(30):8368–8373, 2016. doi: 10.1073/pnas.1602375113. URL <https://www.pnas.org/doi/abs/10.1073/pnas.1602375113>.
- 549 [69] Lorenzo D’Amore, David F. Hahn, David L. Dotson, Joshua T. Horton, Jamshed Anwar, Ian  
550 Craig, Thomas Fox, Alberto Gobbi, Sirish Kaushik Lakkaraju, Xavier Lucas, Katharina Meier,  
551 David L. Mobley, Arjun Narayanan, Christina E. M. Schindler, William C. Swope, Pieter J. in ’t  
552 Veld, Jeffrey Wagner, Bai Xue, and Gary Tresadern. Collaborative assessment of molecular  
553 geometries and energies from the open force field. *Journal of Chemical Information and*  
554 *Modeling*, 62(23):6094—6104, 2022.
- 555 [70] Brajesh K. Rai, Vishnu Sresht, Qingyi Yang, Ray Unwalla, Meihua Tu, Alan M. Mathiowetz,  
556 and Gregory A. Bakken. Torsionnet: A deep neural network to rapidly predict small-molecule  
557 torsional energy profiles with the accuracy of quantum mechanics. *Journal of Chemical*  
558 *Information and Modeling*, 62(4):785–800, 2022.
- 559 [71] Oya Wahl and Thomas Sander. Tautobase: An open tautomer database. *Journal of Chemical*  
560 *Information and Modeling*, 60(3):1085–1089, 2020.
- 561 [72] Raghunathan Ramakrishnan, Pavlo O. Dral, Matthias Rupp, and O. Anatole von Lilienfeld.  
562 Quantum chemistry structures and properties of 134 kilo molecules. *Scientific Data*, 1:140022,  
563 2014.
- 564 [73] Lawrie B. Skinner; Congcong Huang; Daniel Schlesinger; Lars G. M. Pettersson; Anders  
565 Nilsson; Chris J. Benmore. Benchmark oxygen-oxygen pair-distribution function of ambient  
566 water from x-ray diffraction measurements with a wide q-range. *The Journal of Chemical*  
567 *Physics*, 138:074506, 2013.
- 568 [74] Yoshitada Murata Keiko Nishikawa. Liquid structure of carbon tetrachloride and long-range  
569 correlation. *Bulletin of the Chemical Society of Japan*, 52:293–298, 1979.
- 570 [75] Evert Jan Meijer Jan-Willem Handgraaf, Titus S van Erp. Ab initio molecular dynamics study  
571 of liquid methanol. *Chemical Physics Letters*, 367:617–624, 2003.
- 572 [76] László Pusztai Szilvia Pothoczki. Intermolecular orientations in liquid acetonitrile: New insights  
573 based on diffraction measurements and all-atom simulations. *Journal of Molecular Liquids*,  
574 225:160–166, 2017.
- 575 [77] Ilyes Batatia, Philipp Benner, Yuan Chiang, Alin M. Elena, Dávid P. Kovács, Janosh Riebesell,  
576 Xavier R. Advincula, Mark Asta, Matthew Avaylon, William J. Baldwin, Fabian Berger, Noam  
577 Bernstein, Arghya Bhowmik, Filippo Bigi, Samuel M. Blau, Vlad Cărare, Michele Ceriotti,  
578 Sanggyu Chong, James P. Darby, Sandip De, Flaviano Della Pia, Volker L. Deringer, Rokas  
579 Elijošius, Zakariya El-Machachi, Fabio Falcioni, Edvin Fako, Andrea C. Ferrari, John L. A.  
580 Gardner, Mikolaj J. Gawkowski, Annalena Genreith-Schriever, Janine George, Rhys E. A.  
581 Goodall, Jonas Grandel, Clare P. Grey, Petr Grigorev, Shuang Han, Will Handley, Hendrik H.  
582 Heenen, Kersti Hermansson, Christian Holm, Cheuk Hin Ho, Stephan Hofmann, Jad Jaafar,  
583 Konstantin S. Jakob, Hyunwook Jung, Venkat Kapil, Aaron D. Kaplan, Nima Karimitari,  
584 James R. Kermode, Panagiotis Kourtis, Namu Kroupa, Jolla Kullgren, Matthew C. Kuner,  
585 Domantas Kuryla, Guoda Liepuoniute, Chen Lin, Johannes T. Margraf, Ioan-Bogdan Magdău,  
586 Angelos Michaelides, J. Harry Moore, Aakash A. Naik, Samuel P. Niblett, Sam Walton Norwood,  
587 Niamh O’Neill, Christoph Ortner, Kristin A. Persson, Karsten Reuter, Andrew S. Rosen, Louise  
588 A. M. Rosset, Lars L. Schaaf, Christoph Schran, Benjamin X. Shi, Eric Sivonxay, Tamás K.  
589 Stenczel, Viktor Svahn, Christopher Sutton, Thomas D. Swinburne, Jules Tilly, Cas van der  
590 Oord, Santiago Vargas, Eszter Varga-Umbrich, Tejs Vegge, Martin Vondrák, Yangshuai Wang,  
591 William C. Witt, Thomas Wolf, Fabian Zills, and Gábor Csányi. A foundation model for  
592 atomistic materials chemistry, 2025. URL <https://arxiv.org/abs/2401.00096>.

## 593 A Benchmarks overview

594 Each benchmark in MLIP-Audit includes a brief introduction that outlines its purpose, helping  
595 users understand the relevance of the task and how it reflects molecular challenges. A link to the  
596 documentation is provided for users who want a deeper explanation of the benchmark’s design,  
597 scientific context, datasets and implementation details. A description of each benchmark’s dataset can  
598 be found in Appendix C-Table 4. This is followed by key performance metrics for the best-performing  
599 model, along with a summary of results across all analysed MLIP models. Depending on the nature  
600 of the benchmark, additional visualisations may be included, such as radial distribution functions for  
601 molecular liquids or torsion energy profiles for small molecules, which users can explore interactively  
or download for further analysis (Figure 5).

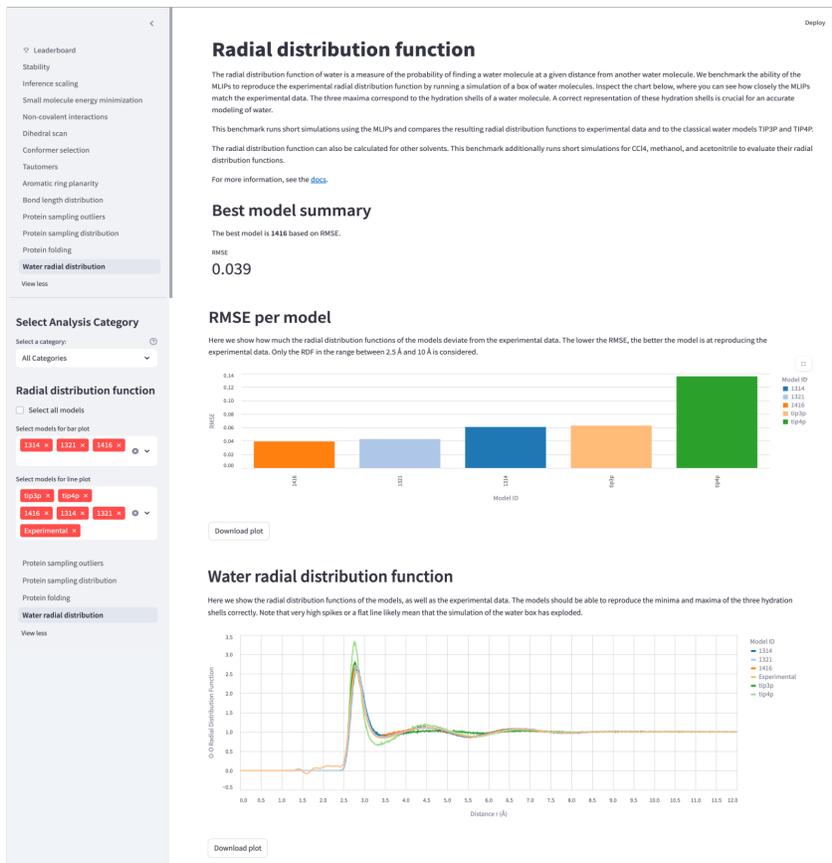


Figure 5: MLIPAudit interface

602

603 In the following subsections, we describe the composition, rationale, and evaluation criteria for each  
604 benchmark category: (i) general systems designed for molecular dynamics stability and scaling, (ii)  
605 small molecules relevant to pharmaceutical and materials chemistry, and (iii) biomolecules, which  
606 pose unique challenges due to their size, flexibility, and hierarchical structure.

### 607 A.1 General benchmarks

608 The general benchmarks implemented in MLIP Audit are system-agnostic and focus on fundamental  
609 molecular dynamics (MD) stability and performance metrics that are applicable across molecular  
610 systems. Two benchmarks are included in this category:

- 611 • **Stability:** assesses the dynamical stability of an MLIP during an MD simulation for a  
612 diverse set of large biomolecular systems. For each system, the benchmark performs an MD

613 simulation using the MLIP model in the NVT ensemble at 300 K for 100,000 steps (100 ps),  
614 leveraging the jax-md engine, as integrated via the mlip library[25]. The test monitors the  
615 system for signs of instability by detecting abrupt temperature spikes (“explosions”) and  
616 hydrogen atom drift. These indicators help determine whether the MLIP maintains stable  
617 and physically consistent dynamics over extended simulation times.

618 • **Inference Scaling:** evaluates how the computational cost of an MLIP scales with the system  
619 size. By running single, long MD episodes on a series of molecular systems of increasing  
620 size, we systematically assess the relationship between molecular complexity and inference  
621 performance. This benchmark is not used for scoring, but it aims at helping the user to pick  
622 the best model in terms of time-to-solution for the application task.

## 623 A.2 Small Molecules

624 MLIPAudit small-molecule benchmarks focus on the ability of MLIPs to reproduce the properties  
625 and dynamics of small organic molecules, including their conformational sampling and interactions  
626 with other molecules. In order of task complexity:

627 • **Bond Length:** evaluates the ability of MLIPs to accurately model the equilibrium bond  
628 lengths of small organic molecules during MD simulations. This is an important test to  
629 understand whether the MLIP respects basic chemistry throughout simulations. Accurate  
630 prediction of bond length is crucial for capturing the structural and electronic properties  
631 of any chemically relevant compounds. For each molecule in the dataset, the benchmark  
632 performs an MD simulation with the same configuration described in the stability benchmark.  
633 Throughout the trajectory, the positions of the bond atoms are tracked, and their deviation  
634 from a reference bond length of the QM-optimised starting structure is calculated. The  
635 average deviation over the trajectory provides a direct measure of the MLIP’s ability to  
636 maintain bond lengths under thermal fluctuations, enabling quantitative comparison to  
637 reference data or other models.

638 • **Ring Planarity:** evaluates the ability of MLIPs to preserve the planarity of aromatic  
639 and conjugated rings in small organic molecules during molecular dynamics simulations.  
640 Aromatic rings (e.g., benzene) are inherently planar due to delocalised  $\pi$  electrons. Ring  
641 planarity enforcement is crucial in molecular dynamics simulations because it preserves  
642 the correct geometry, electronic structure, and interactions of aromatic and conjugated  
643 systems. Without proper planarity (e.g., via improper torsions), simulations can produce  
644 chemically unrealistic distortions that compromise accuracy in energy, flexibility, and  
645 binding predictions. This is especially important in molecules like benzene, tyrosine side  
646 chains, nucleobases, and drug scaffolds, where planarity governs stacking, hydrogen bonding,  
647 and overall stability. For each molecule in the dataset, the benchmark performs an MD  
648 simulation with the same configuration described in the stability benchmark. Throughout  
649 the trajectory, the positions of the ring atoms are tracked, and their deviation from a perfect  
650 plane is quantified using the root mean square deviation (RMSD) from planarity. The ideal  
651 plane of the ring is computed using a principal component analysis of the ring’s atoms.  
652 The average deviation over the trajectory provides a direct measure of the MLIP’s ability  
653 to maintain ring planarity under thermal fluctuations, enabling quantitative comparison to  
654 reference data or other models.

655 • **Dihedral Scan:** evaluates the MLIP’s ability to reproduce torsional energy profiles of  
656 rotatable bonds in small molecules, aiming to approach the quantum-mechanical QM  
657 reference quality. Dihedral scans are essential for mapping how a molecule’s energy changes  
658 as bonds rotate, revealing preferred conformations and energy barriers. Beyond force field  
659 development, they are also used in studying reaction mechanisms, analysing conformational  
660 dynamics in drug discovery, validating quantum chemistry methods, and guiding the design  
661 of flexible or constrained molecules. For each molecule, the benchmark leverages the mlip  
662 library for model inference, comparing the predicted energies along a dihedral scan to QM  
663 reference energy profiles. The reference profile is shifted so that its global minimum is zero,

664 and the MLIP profile is aligned to the same conformer. Performance is quantified using  
665 the following metrics: MAE and RMSE. The Pearson correlation coefficient between the  
666 MLIP-predicted and reference datapoints and the mean barrier height error.

- 667 • **Non-covalent Interactions:** tests if the MLIP can reproduce interaction energies of molec-  
668 ular complexes driven by non-covalent interactions. Non-covalent interactions are of the  
669 highest importance for the structure and function of every biological molecule. This bench-  
670 mark assesses a broad range of interaction types: London dispersion, hydrogen bonds, ionic  
671 hydrogen bonds, repulsive contacts and sigma hole interactions. Assessing the accuracy of  
672 non-covalent interactions is crucial for evaluating how well computational models capture  
673 key forces like hydrogen bonding,  $\pi$ - $\pi$  stacking, and van der Waals interactions that govern  
674 molecular recognition, binding, and assembly. This is essential not only for force field  
675 development, but also for validating quantum methods, guiding molecular design, modelling  
676 biomolecular interfaces, and studying condensed-phase behaviour such as solvation and  
677 aggregation. The benchmark runs energy inference on all structures of the distance scans  
678 of bi-molecular complexes in the dataset. The key metric is the RMSE of the interaction  
679 energy, which is the minimum of the energy well in the distance scan, relative to the energy  
680 of the dissociated complex, compared to the reference data. For repulsive contacts, the  
681 maximum of the energy profile is used instead. Some of the molecular complexes in the  
682 benchmark dataset contain exotic elements (see dataset section). In case the MLIP has never  
683 seen an element of a molecular complex, this complex will be skipped in the benchmark.
- 684 • **Reference Geometry Stability:** assesses the MLIP’s capability to preserve the ground-state  
685 geometry of organic small molecules during energy minimisation, ensuring that initial DFT-  
686 optimised structures remain accurate and physically consistent. Each system is minimised  
687 using the Broyden–Fletcher–Goldfarb–Shanno (BFGS) algorithm (ASE default parameters).  
688 After minimisation, structural fidelity is assessed by computing the RMSD of all heavy  
689 atoms relative to the initial geometry, using the RMSD implementation provided by mdtraj  
690 [56].
- 691 • **Conformer Selection:** evaluates the MLIP’s ability to identify the most stable conformers  
692 within an ensemble of flexible organic molecules and accurately predict their relative energy  
693 differences. It focuses on capturing subtle intramolecular interactions and strain effects that  
694 influence conformational energies. These metrics assess both numerical accuracy and the  
695 MLIP’s ability to preserve relative conformer energetics, which is crucial for downstream  
696 applications such as conformational sampling and compound ranking.
- 697 • **Tautomers:** assesses the ability of MLIP to accurately predict the relative energies and  
698 stabilities of tautomeric forms of small molecules in vacuum. Tautomers are structural  
699 isomers that interconvert via proton transfer and/or double bond rearrangement, and ac-  
700 curately estimating the energy gap between them is an important measure of chemical  
701 accuracy in the MLIP framework. Tautomer ranking assesses a model’s ability to predict the  
702 relative stability of different tautomeric forms of a molecule, which is critical for accurately  
703 modelling protonation states, reactivity, and binding affinities. It is especially important in  
704 drug discovery, quantum method benchmarking, and cheminformatics, where tautomers  
705 can dramatically affect molecular properties and biological activity. For each molecule, the  
706 benchmark compares MLIP-predicted energies against QM reference data. Performance  
707 is quantified by comparing the absolute deviation of the energy difference between the  
708 tautomeric forms from the DFT data.
- 709 • **Reactivity:** assesses the MLIP’s capability to model chemical reactivity. The reactivity-tst  
710 benchmark tests the ability to predict the energy of transition states relative to the reaction’s  
711 reactants and products and thereby the activation energy and enthalpy of a reaction. This  
712 benchmark calculates the energy of reactants, products and transition states of a large dataset  
713 of reactions. From the difference between these states, the activation energy and enthalpy of  
714 formation can be calculated. The performance is quantified using the MAE and RMSE in  
715 activation energy and enthalpy of formation. The reactivity-neb benchmark evaluates the

716 capability to converge a set of nudged elastic band calculations with a known transition state.  
717 The performance is quantified by the percentage of converged calculations.

### 718 A.3 Molecular Liquids

719 The MLIP Audit molecular liquids benchmark focuses on assessing long-range interactions by  
720 computing the radial distribution function for different molecular liquids.

721 • **Radial Distribution Function:** assesses the ability of MLIP to accurately reproduce  
722 the radial distribution function (RDF) of liquids. The RDF characterises the local and  
723 intermediate-range structure of a liquid by describing how particle density varies as a  
724 function of distance from a reference particle. Accurate modelling of the RDF is essential  
725 for capturing both short-range ordering and long-range interactions, which are critical for  
726 understanding the microscopic structure and emergent properties of liquid systems. The  
727 benchmark performs an MD simulation using the MLIP model in the NVT ensemble at  
728 300 K for 500,000 steps, leveraging the jax-md engine from the mlip library. The starting  
729 configuration is already equilibrated. For every specific atom pair (e.g., oxygen-oxygen in  
730 water), the radial distribution function (RDF or  $g(r)$ ) is calculated from the simulation, as:

$$g(r) = \frac{1}{4\pi r^2 \rho N} \left\langle \sum_{i=1}^N \sum_{j \neq i}^N \delta(r - r_{ij}) \right\rangle \quad (1)$$

731 where:  $r$  is the distance from a reference particle,  $\rho$  is the average number density,  $N$  is the  
732 number of particles,  $r_{ij}$  is the distance between particles and  $\delta$  is the Dirac delta function.  
733 Angle brackets denote an ensemble average. For each test case, the benchmark computes  
734  $r_{\text{peak}} = \arg \max_r g(r)$  and compares it with the experimental value for the first solvation  
735 shell.

736 • **Tetrahedral Order Parameter:** evaluates the ability of an MLIP to reproduce the tetrahedral  
737 structure of liquid water by computing the tetrahedrality ( $q$ -number) around each water  
738 molecule. This descriptor quantifies how closely the local arrangement of neighbouring  
739 molecules matches an ideal tetrahedral geometry, a defining feature of hydrogen-bonded  
740 water networks and a key determinant of liquid water’s structural and thermodynamic  
741 properties. The benchmark performs an MD simulation in the NVT ensemble at 300 K for  
742 500,000 steps using the jax-md engine from the mlip library, starting from an equilibrated  
743 configuration. For each oxygen atom, the four nearest oxygen neighbours are identified, and  
744 the tetrahedral order parameter  $q$  is computed as:

$$q = 1 - \frac{3}{8} \sum_{j=1}^4 \sum_{k=j+1}^4 \left( \cos \psi_{jik} + \frac{1}{3} \right)^2 \quad (2)$$

745 where  $\psi_{jik}$  is the angle between vectors  $\mathbf{r}_{ij}$  and  $\mathbf{r}_{ik}$  connecting the central oxygen  $i$  to  
746 neighbours  $j$  and  $k$ . A value of  $q = 1$  corresponds to a perfect tetrahedral environment,  
747 while  $q = 0$  indicates a fully disordered one. For each test case, the benchmark reports the  
748 mean tetrahedrality  $\langle q \rangle$  and compares it against experimental and first-principles reference  
749 values, providing a stringent evaluation of a model’s ability to capture hydrogen-bond  
750 network structure in liquid water.

### 751 A.4 Biomolecules

752 MLIP Audit biomolecule benchmarks focus on assessing the properties and dynamics of proteins,  
753 including their folding behaviour, structural stability, and conformational sampling.

754 • **Protein Folding Stability:** evaluates the ability of an MLIP to preserve the structural  
755 integrity of experimentally determined protein conformations during MD simulations. It  
756 assesses the retention of secondary structure elements and overall compactness across a

757 set of known protein structures. This module analyses the folding trajectories of proteins  
 758 in MLIP simulations. For each molecule in the dataset, the benchmark performs an MD  
 759 simulation with the same configuration described in the stability benchmark. We track how  
 760 Root Mean Square Deviation (RMSD), TM Score [57], Dictionary of Secondary Structure  
 761 in Proteins (DSSP) [58] and Radius of Gyration change over time.

- 762 • **Sampling Outlier Detection:** Assesses the structural quality of sampled conformations  
 763 by computing backbone Ramachandran angles ( $\phi/\psi$ ) and side-chain rotamer angles ( $\chi$ ),  
 764 and identifying outliers through comparison with reference rotamer libraries [59]. For  
 765 each molecule in the dataset, the benchmark performs an MD simulation with the same  
 766 configuration described in the stability benchmark. The outlier detection identifies residues  
 767 whose dihedral angles fall outside expected ranges, relying on the fast KDtree [60] scipy  
 768 [61] implementation. The analysis provides a global percentage of outliers for backbone  
 769 and rotamers per structure, as well as a more detailed analysis per residue type.

## 770 B Benchmarks scoring

771 To enable consistent and fair comparison across models, we define a composite score that aggregates  
 772 performance over all compatible benchmarks. Each benchmark  $b \in \mathcal{B}$  may report one or more metrics  
 773  $x_{m,b}^{(i)}$ , where  $i = 1, \dots, N_b$  indexes the  $N_b$  metrics evaluated for the model  $m$ . For each metric, we  
 774 compute a normalised score using a soft thresholding function based on a DFT-derived reference  
 775 tolerance  $t_b^{(i)}$  (see 3):

$$s_{m,b}^{(i)} = \begin{cases} 1, & \text{if } x_{m,b}^{(i)} \leq t_b^{(i)} \\ \exp\left(-\alpha \cdot \frac{x_{m,b}^{(i)} - t_b^{(i)}}{t_b^{(i)}}\right), & \text{otherwise} \end{cases}$$

776 where  $\alpha$  is a tunable parameter controlling the steepness of the penalty (e.g.,  $\alpha = 3$ ). The per-  
 777 benchmark score is then computed as the average over all its metric scores:

$$s_{m,b} = \frac{1}{N_b} \sum_{i=1}^{N_b} s_{m,b}^{(i)}$$

778 Let  $\mathcal{B}_m \subseteq \mathcal{B}$  denote the subset of benchmarks for which the model  $m$  has valid data (i.e., benchmarks  
 779 compatible with its element set). The final model score is the mean over all benchmarks on which the  
 780 model could be evaluated:

$$S_m = \frac{1}{|\mathcal{B}_m|} \sum_{b \in \mathcal{B}_m} s_{m,b}$$

781 This scoring framework ensures that models are rewarded for meeting or exceeding DFT-level  
 782 accuracy. In the current version, full benchmarks are skipped if a model does not have all the  
 783 necessary chemical elements to run all the test cases. This is true for all benchmarks, but non-covalent  
 784 interactions, in which we do a per-test-case exception. Benchmarks with multiple metrics contribute  
 785 proportionally, and the result is a single interpretable score  $S_m \in [0, 1]$  that balances physical fidelity,  
 786 chemical coverage, and overall model robustness. The thresholds for the different benchmarks have  
 787 been chosen based on the literature. In the case of tautomers, energy differences are very small;  
 788 therefore, we’ve chosen a stricter threshold of 1-2 kcal/mol, which is not enough for classification.  
 789 Thresholds for biomolecules are borrowed from traditional literature in molecular modelling.

Table 3: Score thresholds across benchmarks.

<b>Benchmark</b>	<b>Metric</b>	<b>Threshold</b>
Reference Geometry Stability	RMSD (Å)	0.075 [62]
Non-covalent Interactions	Absolute deviation from reference interaction energy (kcal/mol)	1.0 [62]
Dihedral Scan	Mean barrier error (kcal/mol)	1.0 [63]
Conformer Selection	MAE (kcal/mol)	0.5
	RMSE (kcal/mol)	1.5 [64]
Tautomers	Absolute deviation ( $\Delta G$ )	0.05
Ring Planarity	Deviation from plane (Å)	0.05 [65]
Bond Length Distribution	Avg. fluctuation (Å)	0.05 [62]
Reactivity-TST	Activation Energy (kcal/mol)	3.0 [66]
	Enthalpy (kcal/mol)	2.0 [66]
Reactivity-NEB	Final force convergence (eV/Å)	0.05 [67]
Radial Distribution Function	RMSE (Å)	0.1 [68]
Protein Sampling Outliers	Ramachandran ratio	0.1
	Rotamers Ratio	0.03
Protein Folding Stability	min(RMSD) (Å)	2.0
	max(TM-Score)	0.5

790 **C Supporting Figures and Tables**

Table 4: Datasets used for the different benchmarks in MLIPAudit.

Benchmark	Dataset Name	Link/Citation	Content Description
General Stability	In-house dataset	Released with MLIPAudit	3 small molecules in vacuum (1 HCNO-only, 1 with halogens, 1 with sulfur). 2 peptides in vacuum (Neurotensin PDBid 2LNF and Oxytocin PDBid 7OFG). 1 protein in vacuum (PDBid 1A7M). 1 peptide in pure water (Oxytocin). 1 peptide in water with Cl <sup>-</sup> counterions (Neurotensin).
Inference Scaling	In-house dataset	Released with MLIPAudit	Proteins in vacuum. PDBids: 1AY3, 1UAO, 1AB7, 1P79, 1BIP, 1A5E, 1A7M, 2BQV, 1J7H, 5KGZ, 1VSQ, 1JRS.
Reference Geometry Stability	OpenFF	[69]	200 molecules for the neutral dataset and 20 for the charged dataset. The subsets are constructed so that the chemical diversity, as represented by Morgan fingerprints, is maximised.
Non-covalent Interactions	NCI-ATLAS subsets: D442x10, HB375x10, HB300SPXx10, IHB100x10, R739x5, SH250x10	<a href="http://www.nciatlas.org/">http://www.nciatlas.org/</a>	QM optimised geometries of distance scans of bi-molecular complexes, where the two molecules interact via non-covalent interactions with associated energies.
Dihedral Scan	In-house recomputed TorsionNet 500 dataset at $\omega$ B97M-D3(BJ) DFT-level.	[70]	500 structures of drug-like molecules and their energy profiles around selected rotatable bonds at $\omega$ B97M-D3(BJ) DFT-level.
Conformer Selection	Wiggle 150	[49]	50 conformers each of three molecules: Adenosine, Benzylpenicillin, and Efavirenz.
Tautomers	In-house recomputed Tautobase dataset at $\omega$ B97M-D3(BJ) DFT-level.	[71]	2,792 tautomer pairs sourced from the Tautobase dataset. After generation of the structures and minimisation at xtb level, the QM energies were computed in-house using $\omega$ B97M-D3(BJ)/def2-TZVPPD level of theory.
Ring Planarity	QM9 subset	[72]	One representative molecule each, containing substructures for benzene, furan, imidazole, purine, pyridine and pyrrole.
Bond Length	QM9 subset	[72]	One representative molecule each, containing the bond types C-C, C=C, C#C, C-N, C-O, C=O and C-F.
Reactivity	Grambow dataset	[51]	Reactants, products and transition states of 11960 reactions.
Radial Distribution Function	In-house solvent boxes	Released with MLIPAudit. Reference data: [73–76]	Water, CCl <sub>4</sub> , Acetonitrile, Methanol.
Protein Folding Stability	In-house dataset	Released with MLIPAudit	3 solvated proteins: Chignolin, Orexin and Trp Cage. PDBids: 1UAO, 2JOF, 1CQ0.

## 791 D Model training details

792 Three of the models presented in this paper were released as part of the mlip library [25]: ViSNet-  
 793 SPICE2, MACE-SPICE2, and NequIP-SPICE2. Details on how these models were trained, alongside  
 794 training data details and hyperparameters can be found in the original reference.

Table 5: MLIPAudit test-cases overlap with models training dataset for internal models only

Benchmark Category	Benchmark	Overlap [%]
Small-Molecule	Reference Geometry Stability	0
Small-Molecule	Bond Length distribution	0
Small-Molecule	Ring Planarity	0
Small-Molecule	Conformer selection	66.7
Small-Molecule	Dihedral scan	1.4
Small-Molecule	Tautomers	8.4
Small-Molecule	Non-covalent interactions	-
Small-Molecule	Reactivity	-
Molecular liquids	RDF	0
Biomolecules	Folding stability	0
Biomolecules	Sampling	0

Table 6: Category-based rankings (aggregated scores by benchmark category)

Source	Rank	Category	Model Name	Score	Metrics
External	1	General	UMA-Small	1.00	1/1
External	1	General	MACE-SPICE2	1.00	1/1
External	1	General	MACE-MP	1.00	1/1
Internal	1	General	ViSNet-SPICE2	1.00	1/1
Internal	1	General	MACE-SPICE2	1.00	1/1
Internal	2	General	NequIP-SPICE2	0.90	1/1
Internal	3	General	ViSNet-SPICE2-t1x	0.75	1/1
Internal	3	General	ViSNet-t1x	0.00	1/1
Internal	3	General	NequIP-t1x	0.00	1/1
Internal	3	General	MACE-t1x	0.00	1/1
External	1	Small-molecules	UMA-Small	0.56	7/9
External	2	Small-molecules	MACE-OFF	0.50	8/9
External	3	Small-molecules	MACE-MP	0.36	7/9
Internal	1	Small-molecules	ViSNet-SPICE2-t1x	0.65	9/9
Internal	2	Small-molecules	ViSNet-SPICE2	0.52	9/9
Internal	2	Small-molecules	NequIP-SPICE2	0.52	9/9
Internal	3	Small-molecules	MACE-SPICE2	0.51	9/9
Internal	4	Small-molecules	NequIP-t1x	0.16	6/9
Internal	5	Small-molecules	MACE-t1x	0.14	6/9
Internal	6	Small-molecules	ViSNet-t1x	0.11	6/9
External	1	Molecular-liquids	UMA-Small	0.98	2/2
External	2	Molecular-liquids	MACE-OFF	0.73	2/2
External	-	Molecular-liquids	MACE-MP	0.0	2/2
Internal	1	Molecular-liquids	NequIP-SPICE2	0.97	2/2
Internal	1	Molecular-liquids	MACE-SPICE2	0.97	2/2
Internal	1	Molecular-liquids	MACE-SPICE2	0.97	2/2
Internal	2	Molecular-liquids	ViSNet-SPICE2-t1x	0.95	2/2
Internal	-	Molecular-liquids	ViSNet-t1x	0.0	2/2
Internal	-	Molecular-liquids	NequIP-t1x	0.0	2/2
Internal	-	Molecular-liquids	MACE-t1x	0.0	2/2
External	1	Biomolecules	UMA-Small	0.92	2/2
External	1	Biomolecules	MACE-OFF	0.92	2/2
External	-	Biomolecules	MACE-MP	0.79	2/2
Internal	1	Biomolecules	ViSNet-SPICE2	1.00	2/2
Internal	1	Biomolecules	NequIP-SPICE2	1.00	2/2
Internal	2	Biomolecules	ViSNet-SPICE2-t1x	0.43	2/2
Internal	-	Biomolecules	ViSNet-t1x	0.0	2/2
Internal	-	Biomolecules	NequIP-t1x	0.0	2/2
Internal	-	Biomolecules	MACE-t1x	0.0	2/2

Table 7: Single benchmarks rankings

Source	Rank	Benchmark	Model Name	Score	Test Cases
External	1	General Stability	UMA-Small	1.00	8/8
External	1	General Stability	MACE-OFF	1.00	8/8
External	1	General Stability	MACE-MP	1.00	8/8
Internal	1	General Stability	ViSNet-SPICE2	1.00	8/8
Internal	1	General Stability	MACE-SPICE2	1.00	8/8
Internal	2	General Stability	Nequip-SPICE2	0.90	8/8
Internal	3	General Stability	ViSNet-SPICE2-t1x	0.75	8/8
Internal	4	General Stability	MACE-t1x	0.44	8/8
External	1	Solvent RDF	UMA-Small	0.95	3/3
External	2	Solvent RDF	MACE-OFF	0.73	3/3
External	-	Solvent RDF	MACE-MP	0.00	0/3
Internal	1	Solvent RDF	Nequip-SPICE2	0.97	3/3
Internal	2	Solvent RDF	MACE-SPICE2	0.94	3/3
Internal	2	Solvent RDF	ViSNet-SPICE2	0.94	3/3
Internal	3	Solvent RDF	ViSNet-SPICE2-t1x	0.90	3/3
External	1	Water RDF	UMA-small	1.00	1/1
External	2	Water RDF	MACE-OFF	0.56	1/1
External	-	Water RDF	MACE-MP	0.00	1/1
Internal	1	Water RDF	ViSNet-SPICE2	1.00	1/1
Internal	1	Water RDF	MACE-SPICE2	1.00	1/1
Internal	1	Water RDF	Nequip-SPICE2	1.00	1/1
Internal	1	Water RDF	ViSNet-SPICE2-t1x	1.00	1/1
Internal	-	Water RDF	ViSNet-t1x	0.00	1/1
Internal	-	Water RDF	MACE-t1x	0.00	1/1
Internal	-	Water RDF	Nequip-t1x	0.00	1/1
External	1	Water Ang. Dist.	MACE-OFF	1.00	1/1
External	2	Water Ang. Dist.	UMA-Small	0.76	1/1
External	-	Water Ang. Dist.	MACE-MP	0.00	1/1
Internal	1	Water Ang. Dist.	Nequip-SPICE2	0.72	1/1
Internal	2	Water Ang. Dist.	ViSNet-SPICE2-t1x	0.60	1/1
Internal	3	Water Ang. Dist.	Visnet-SPICE2	0.59	1/1
Internal	4	Water Ang. Dist.	MACE-SPICE2	0.51	1/1
Internal	-	Water Ang. Dist.	ViSNet-t1x	0.00	1/1
Internal	-	Water Ang. Dist.	MACE-t1x	0.00	1/1
Internal	-	Water Ang. Dist.	Nequip-t1x	0.00	1/1
External	1	Protein Folding Stability	UMA-Small	1.00	3/3
External	1	Protein Folding Stability	MACE-OFF	1.00	3/3
External	1	Protein Folding Stability	MACE-MP	1.00	3/3
Internal	1	Protein Folding Stability	ViSNet-SPICE2	1.00	3/3
Internal	1	Protein Folding Stability	Nequip-SPICE2	1.00	3/3
Internal	2	Protein Folding Stability	MACE-SPICE2	0.33	3/3
Internal	-	Protein Folding Stability	ViSNet-SPICE2-t1x	0.00	3/3

795 We present in Table 9 below details about the other models presented as examples in the paper. Note  
796 that training details for UMA-Small, MACE-OFF and MACE-MP can be found in Ref. [34], [32],  
797 and [77] respectively.

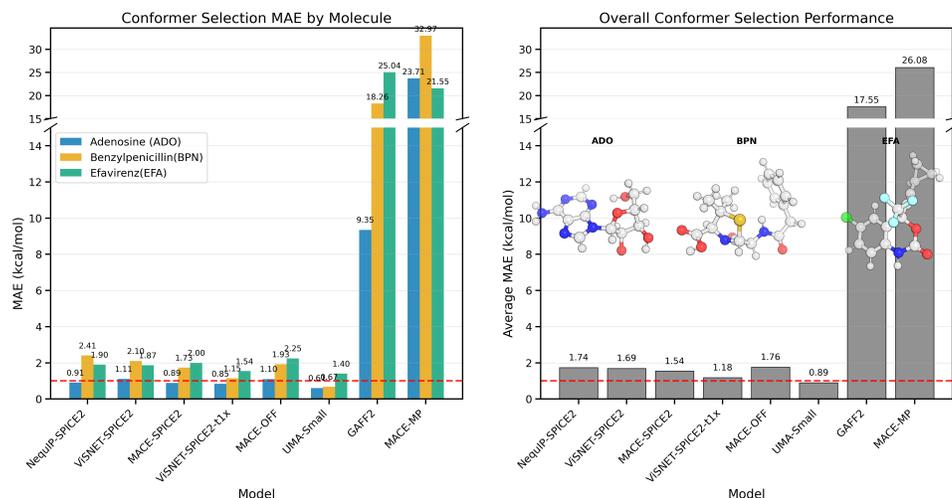


Figure 6: Conformer selection benchmark across three pharmaceutically relevant molecules: adenosine (ADO), benzylpenicillin (BPN), and efavirenz (EFA). MAE is computed with respect to DFT reference conformer energies. DFT threshold (red dashed line at 0.5 kcal/mol). Insets depict representative 3D conformers for each molecule.

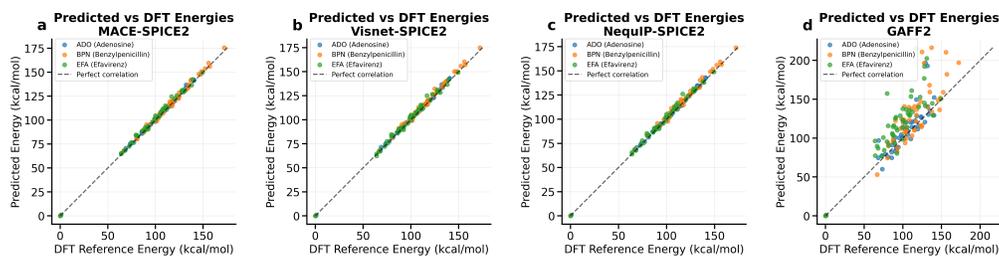


Figure 7: Predicted vs. DFT conformer energies for adenosine (ADO, blue), benzylpenicillin (BPN, orange), and efavirenz (EFA, green).

Table 8: Single benchmarks rankings (cont.)

Source	Rank	Benchmark	Model Name	Score	Test Cases
External	1	Reference Geometry Stability	UMA-Small	0.98	220/220
External	1	Reference Geometry Stability	MACE-OFF	0.93	220/220
External	1	Reference Geometry Stability	MACE-MP	0.50	220/220
Internal	1	Reference Geometry Stability	ViSNet-SPICE2-t1x	0.97	220/220
Internal	1	Reference Geometry Stability	ViSNet-SPICE2	0.97	220/220
Internal	2	Reference Geometry Stability	MACE-SPICE2	0.96	220/220
Internal	3	Reference Geometry Stability	Nequip-SPICE2	0.94	220/220
External	1	Conformer Selection	UMA-Small	0.29	3/3
External	-	Conformer Selection	MACE-OFF	0.00	3/3
External	-	Conformer Selection	MACE-MP	0.00	3/3
Internal	1	Conformer Selection	ViSNet-SPICE2-t1x	0.05	3/3
Internal	2	Conformer Selection	Nequip-SPICE2	0.03	3/3
Internal	2	Conformer Selection	MACE-SPICE2	0.03	3/3
Internal	-	Conformer Selection	Visnet-SPICE2	0.00	3/3
External	1	Dihedral Scan	UMA-Small	0.71	500/500
External	2	Dihedral Scan	MACE-OFF	0.66	500/500
External	2	Dihedral Scan	MACE-MP	0.40	500/500
Internal	1	Dihedral Scan	ViSNet-SPICE2-t1x	0.70	500/500
Internal	2	Dihedral Scan	ViSNet-SPICE2	0.69	500/500
Internal	2	Dihedral Scan	Nequip-SPICE2	0.66	500/500
Internal	3	Dihedral Scan	MACE-SPICE2	0.65	500/500
External	1	Non-covalent Interactions	UMA-Small	0.84	2192/2206
External	2	Non-covalent Interactions	MACE-OFF	0.70	1728/2206
External	3	Non-covalent Interactions	MACE-MP	0.44	2206/2206
Internal	1	Non-covalent Interactions	MACE-SPICE2	0.75	1807/2206
Internal	2	Non-covalent Interactions	Visnet-SPICE2	0.73	1807/2206
Internal	2	Non-covalent Interactions	Nequip-SPICE2	0.73	1807/2206
Internal	3	Non-covalent Interactions	Visnet-SPICE2-t1x	0.68	1807/2206
Internal	4	Non-covalent Interactions	Nequip-t1x	0.44	689/2206
Internal	5	Non-covalent Interactions	MACE-t1x	0.43	689/2206
Internal	6	Non-covalent Interactions	Visnet-t1x	0.21	689/2206
External	1	Reactivity	UMA-Small	0.86	11961/11961
External	1	Reactivity	MACE-OFF	0.12	11961/11961
External	1	Reactivity	MACE-MP	0.04	11961/11961
Internal	1	Reactivity	Visnet-SPICE2-t1x	0.77	11961/11961
Internal	2	Reactivity	Nequip-t1x	0.44	11961/11961
Internal	3	Reactivity	MACE-t1x	0.43	11961/11961
Internal	4	Reactivity	Visnet-t1x	0.22	11961/11961
Internal	5	Reactivity	MACE-SPICE2	0.10	11961/11961
Internal	6	Reactivity	Visnet-SPICE2	0.05	11961/11961
Internal	7	Reactivity	Nequip-SPICE2	0.04	11961/11961
Internal	1	Nudged Elastic Band	Visnet-SPICE2-t1x	0.58	100/100
Internal	1	Nudged Elastic Band	Nequip-t1x	0.58	100/100
Internal	2	Nudged Elastic Band	MACE-t1x	0.44	100/100
Internal	3	Nudged Elastic Band	Visnet-t1x	0.38	100/100
External	1	Tautomers	UMA-Small	0.23	1391/1391
External	2	Tautomers	MACE-OFF	0.07	1391/1391
External	-	Tautomers	MACE-MP	0.00	1391/1391
Internal	1	Tautomers	Nequip-SPICE2	0.11	1391/1391
Internal	2	Tautomers	Visnet-SPICE2	0.10	1391/1391
Internal	3	Tautomers	Visnet-SPICE2-t1x	0.09	1391/1391
Internal	3	Tautomers	MACE-SPICE2	0.05	1391/1391
External	1	Bond Length	UMA-Small	1.00	8/8
External	1	Bond Length	MACE-OFF	1.00	8/8
External	1	Bond Length	MACE-MP	1.00	8/8
Internal	1	Bond Length	Visnet-SPICE2-t1x	1.00	8/8
Internal	1	Bond Length	Visnet-SPICE2	1.00	8/8
Internal	1	Bond Length	MACE-SPICE2	1.00	8/8
Internal	1	Bond Length	Nequip-SPICE2	1.00	8/8
External	1	Ring Planarity	MACE-OFF	0.99	6/6
External	2	Ring Planarity	UMA-Small	0.98	6/6
External	1	Ring Planarity	MACE-MP	0.80	6/6
Internal	1	Ring Planarity	Visnet-SPICE2-t1x	1.00	6/6

Table 9: Example models training details

<b>Model</b>	<b>Dataset</b>	<b>Hyperparameters</b>
ViSNet-SPICE2	Original version of SPICE2 [47], as curated in [25] - includes only neutral systems	As described in [25]
MACE-SPICE2	Original version of SPICE2 [47], as curated in [25] - includes only neutral systems	As described in [25]
NequIP-SPICE2	Original version of SPICE2 [47], as curated in [25] - includes only neutral systems	As described in [25]
ViSNet-t1x	Original version of Transition-1X [48], trained on 1M samples, randomly sampled with 95/5 train/val split.	Same as ViSNet-SPICE2
MACE-t1x	Original version of Transition-1X [48], trained on 1M samples, randomly sampled with 95/5 train/val split.	Same as MACE-SPICE2
NequIP-t1x	Original version of Transition-1X [48], trained on 1M samples, randomly sampled with 95/5 train/val split.	Same as NequIP-SPICE2
ViSNet-SPICE2(charged)-t1x	SPICE2 and Transition-1X as recomputed in the OMOL dataset [42]. SPICE2 is curated as is described in [25]. T1X includes 50k samples, selected among transition states, reactants and products.	Same as ViSNet-SPICE2, except for the number of channels with is increased to 256.
ViSNet-SPICE2(neutral)-t1x	SPICE2 and Transition-1X as recomputed in the OMOL dataset [42]. SPICE2 is curated as is described in [25]. T1X includes 1M samples, selected among transition states, reactants and products.	Same as ViSNet-SPICE2, except for the number of channels with is increased to 256.