



PAPER • OPEN ACCESS

Multi-fidelity transfer learning for quantum chemical data using a robust density functional tight binding baseline

To cite this article: Mengnan Cui *et al* 2025 *Mach. Learn.: Sci. Technol.* **6** 015071

View the [article online](#) for updates and enhancements.

You may also like

- [Improving the generative performance of chemical autoencoders through transfer learning](#)

Nicolae C Iovanac and Brett M Savoie

- [The good, the bad, and the ugly of atomistic learning for 'clusters-to-bulk' generalization](#)

Mikoaj J Gawkowski, Mingjia Li, Benjamin X Shi et al.

- [Bridging the gap between high-level quantum chemical methods and deep learning models](#)

Viki Kumar Prasad, Alberto Otero-de-la-Roza and Gino A DiLabio



PAPER

OPEN ACCESS

RECEIVED

21 November 2024

REVISED

3 March 2025

ACCEPTED FOR PUBLICATION

18 March 2025

PUBLISHED

28 March 2025

Original Content from this work may be used under the terms of the [Creative Commons Attribution 4.0 licence](#).

Any further distribution of this work must maintain attribution to the author(s) and the title of the work, journal citation and DOI.



Multi-fidelity transfer learning for quantum chemical data using a robust density functional tight binding baseline

Mengnan Cui^{1,2} , Karsten Reuter¹ and Johannes T Margraf^{1,2,*} ¹ Fritz Haber Institute of the Max Planck Society, Berlin, Germany² University of Bayreuth, Bavarian Center for Battery Technology (BayBatt), Bayreuth, Germany

* Author to whom any correspondence should be addressed.

E-mail: johannes.margraf@uni-bayreuth.de**Keywords:** multi-fidelity, transfer, learning, quantum, chemical, robustSupplementary material for this article is available [online](#)

Abstract

Machine learning has revolutionized the development of interatomic potentials over the past decade, offering unparalleled computational speed without compromising accuracy. However, the performance of these models is highly dependent on the quality and amount of training data. Consequently, the current scarcity of high-fidelity datasets (i.e. beyond semilocal density functional theory) represents a significant challenge for further improvement. To address this, this study investigates the performance of transfer learning (TL) across multiple fidelities for both molecules and materials. Crucially, we disentangle the effects of multiple fidelities and different configuration/chemical spaces for pre-training and fine-tuning, in order to gain a deeper understanding of TL for chemical applications. This reveals that negative transfer, driven by noise from low-fidelity methods such as a density functional tight binding baseline, can significantly impact fine-tuned models. Despite this, the multi-fidelity approach demonstrates superior performance compared to single-fidelity learning. Interestingly, it even outperforms TL based on foundation models in some cases, by leveraging an optimal overlap of pre-training and fine-tuning chemical spaces.

1. Introduction

Spurred by the high computational costs of first-principles electronic structure methods, the development of machine learning (ML) interatomic potentials has enabled accurate atomistic simulations for previously inaccessible systems [1, 2]. Early efforts are exemplified by the pioneering works of Behler and Parrinello [3], Popelier [4], as well as the Gaussian Approximation Potentials of Csányi and co-workers [5], among others. These involve the construction of rotationally invariant representations of atomic environments, combined with shallow neural networks or kernel regression methods. Subsequently, graph neural networks were developed, which expand local environment representations through message-passing mechanisms, such as in the SchNet [6], PhysNet [7], DimeNet [8], DTNN [9], and GemNet [10] models. Most recently, equivariant networks such as NequIP [11], PaiNN [12], SpookyNet [13], NewtonNet [12], and MACE [14] have emerged that currently represent the state-of-the-art in atomistic ML. Among these, the MACE model employed in this work builds on the advantages of atomic cluster expansion by constructing high-body ordered equivariant features and mapping them to structural potential energy through a message-passing neural network. Its predictive accuracy and generalization capability have been demonstrated for a wide range of applications [15], e.g. for bond dissociation energies [16] or molecular vibrations [17].

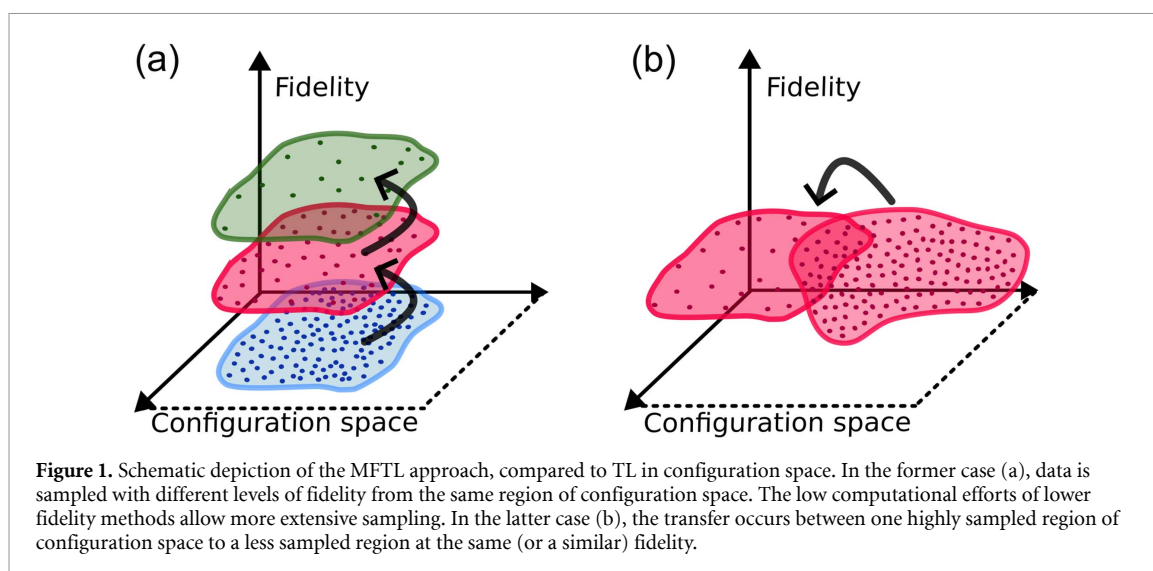
These methodological developments have led to remarkable improvements in accuracy, but also increased the resource demand of training and inference [18, 19]. In addition to increasing model capacity by using deeper or more complex architectures, accuracy can usually also be improved by increasing the amount of training data. This has the advantage that it does not affect the inference cost of the model, only the training cost [20]. It should be noted, that more data implies improved coverage of the configuration space

of interest with higher structural diversity in this context. Clearly, adding repetitive or irrelevant data is not helpful. For a given level of accuracy, the lower computational cost at inference time can thus be offset by increased cost for training data generation. Unfortunately, this can itself be prohibitively expensive, e.g. when highly accurate reference methods such as coupled cluster (CC) theory or quantum Monte Carlo are used. In such contexts, transfer learning (TL) is commonly used [21–23], meaning that a model that is pre-trained on one large dataset (for example a general and/or low-fidelity one) is fine-tuned on another (for example a specialized or high-fidelity one). In the best case, beneficial features of the pre-trained model can be maintained throughout the fine-tuning, leading to more accurate and robust models for a given training set size.

TL is an appealing idea and widely used in chemistry [24, 25]. Within the TL literature, the similarity between datasets for pre-training and fine-tuning is considered crucial to the success of this method. When the overlap is limited, traditional models often fail to generalize effectively due to insufficient shared features or representations [21, 26]. In this context, the overlap between quantum chemical datasets should be considered both in terms of the configuration space covered by the structures and the fidelities of the quantum chemical methods. For instance, Hutchinson *et al* implemented TL to enhance the accuracy of experimental band gap predictions, where comparatively cheaper density functional theory (DFT) band gaps are used for model pre-training (transfer from low- to high-fidelity) [27]. Likewise, Frey *et al* demonstrated that the MEGNet model, pre-trained on tens of thousands of 3D bulk crystals, could be fine-tuned to efficiently predict the properties of 2D materials (transfer from one chemical space to another) [28, 29]. With the recent advent of broadly applicable foundation models for materials and molecules, TL is becoming even more relevant [15, 28, 30]. However, there are two common issues that need to be avoided in this context. On one hand, negative transfer can occur, meaning that features learned during pre-training can in some cases be detrimental to the task in the fine-tuning step. On the other hand, catastrophic forgetting can occur, meaning that the fine-tuning essentially overwrites all pre-trained information, rendering the pre-training step irrelevant [25, 31]. Both of these are especially pertinent when overlap between the pre-training and fine-tuning datasets is insufficient. In chemical applications, this idea of dataset overlap relates both to the types of structures included in each set (i.e. how similar are bulk crystals and 2D materials) and the fidelity of the reference data (i.e. how good is the agreement between low and high fidelity labels).

When foundation models are fine-tuned, the main focus is on structural overlap. For example, models like MACE-MP-0 [15] (for materials) and MACE-OFF23 [32] (for molecules) are able to extrapolate remarkably well to diverse systems, including liquids, amorphous phases, and extreme conditions (e.g. high temperatures and pressures). This success is notable given that these models were trained exclusively on near-ground state structures of inorganic crystals and isolated molecules and clusters, respectively. This remarkable extrapolation capability on highly diverse systems is due to the high-body order equivariant graph neural network architecture of MACE, as well as efforts to prepare diverse large training sets. Nevertheless, the further the configurations of the intended application are from those in the pre-training dataset, the more additional data will be required to obtain an accurate fine-tuned model. This explains the appeal of TL in a multi-fidelity setting [33, 34]. Here, additional data can be generated cheaply with low costs method beforehand, for exactly the kinds of structures that are of interest for a given application (i.e. with perfect structural overlap between pre-training and fine-tuning sets). For this reason, multi-fidelity approaches have been widely used in chemical applications such as molecular crystal structure prediction, hybrid models with long-range electrostatics, materials screening, and implicit solvation. These applications use both TL and other multi-fidelity settings, including Δ -ML [35, 36], multi-task learning [37–39], or meta-learning [40]. In principle, there is thus a trade-off between structural overlap and fidelity overlap. In practice, both aspects are usually confounded however, since pre-training is often performed with pre-existing databases (e.g. the Materials Project [41] or SPICE datasets [42]) with a predefined chemical space and level of theory. Meanwhile, the chemical space and level of theory for fine-tuning are defined by the target application.

The goal of this paper is to systematically explore the impact of structural and fidelity overlap in TL. To achieve this, we use the recently reported periodic table baseline parameters (PTBPs) for density functional tight binding (DFTB) calculations. These parameters allow us to generate data for molecular and material datasets at a low computational expense. With this, we generate customized low-fidelity datasets for arbitrary configuration spaces, with perfect structural overlap. The corresponding multi-fidelity TL (MFTL) models are compared with TL models based on pre-trained foundation models, which by definition feature a lower degree of structural overlap, but are trained with higher fidelity reference data.



2. Results and discussions

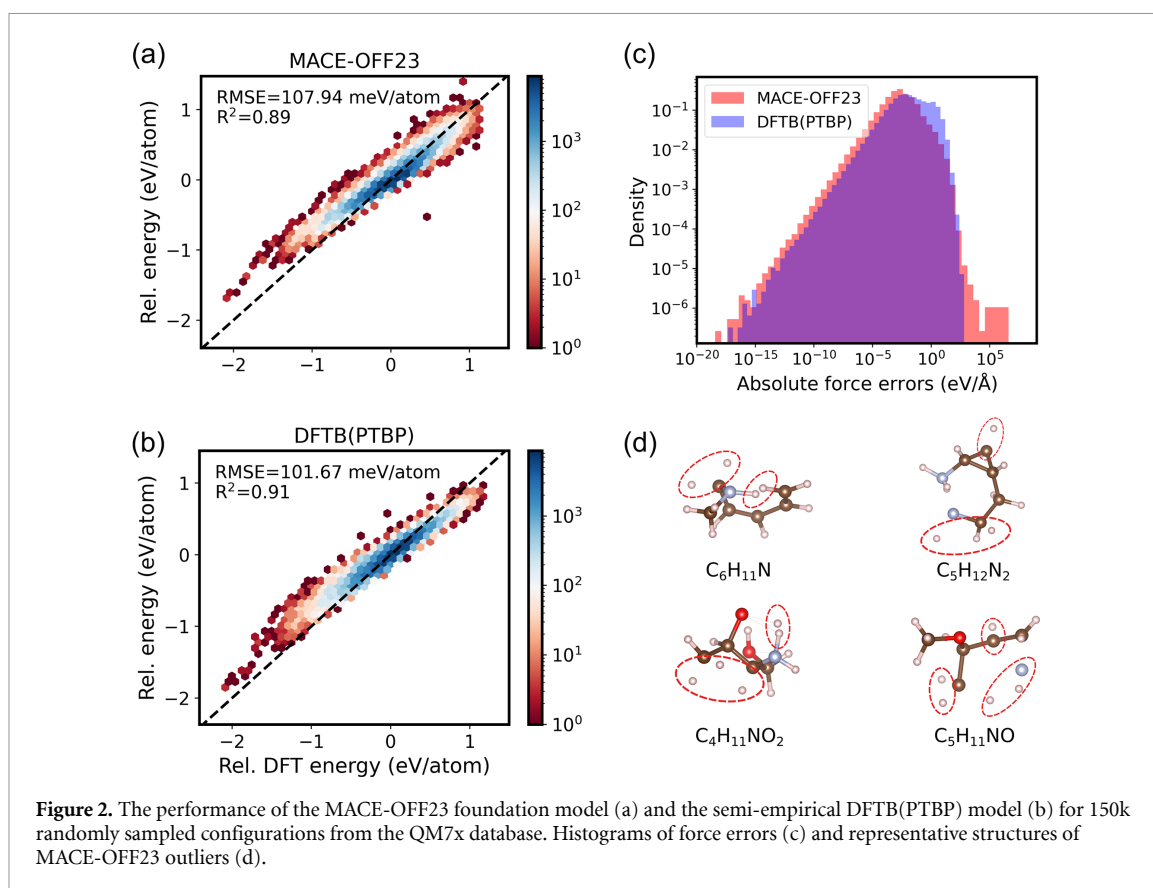
2.1. Model definition

Figure 1(a) illustrates the general MFTL workflow used herein. In the initial stage, DFTB is employed as an efficient method for generating training labels for a large sample from the target configuration space. The interatomic potential (e.g. MACE in this case) is pre-trained using this low-fidelity data. As discussed in the introduction, the key advantage of the resulting models (compared to existing foundational models) is that they are trained on data with perfect structural overlap for the target application. The pre-trained low-fidelity model is subsequently fine-tuned using a smaller but more accurate high-fidelity dataset. This fine-tuning step enhances the model's predictive performance. In principle, this fine-tuning process can be iteratively extended across multiple fidelities until the desired level of accuracy for the final target is achieved, as demonstrated further below [43]. This MFTL approach is in contrast to configuration space TL illustrated in figure 1(b).

2.2. Molecular systems

To demonstrate the benefit of MFTL, we first consider the QM7x dataset [44], which represents an extensive configuration space of small organic molecules in equilibrium and non-equilibrium configurations. Specifically, QM7x comprises approximately 4.2 million configurations based on the enumeration of small organic molecules containing up to 7 heavy atoms (i.e. C, N, O, S, Cl). The molecular sizes span 4–23 atoms in total. For each configuration, total energies and forces were calculated at the hybrid PBE0 level [45] with the many-body dispersion correction [46], hereafter referred to as PBE0+MBD. To evaluate the impact of training set size, we randomly sampled training and validation sets with sizes of 0.5k, 1k, 3k, 10k, and 50k configurations, respectively. An additional 50k configurations were reserved as an independent test set for final evaluations.

The performance of the low-fidelity DFTB method (using the PTBP parameters) compared to the target PBE0+MBD method is shown in figure 2. For comparison, we also show the performance of the recent MACE-OFF23 foundation model. Perhaps surprisingly, the foundation model displays similar error statistics as the PTBP model (root mean squared errors (RMSEs) of 107.94 and 101.67 meV atom⁻¹ for relative energies, respectively), despite it being trained on similar organic molecules, whereas PTBP is a simple DFTB model fitted on inorganic solids. These deviations can partially be attributed to the different levels of theory used for training MACE-OFF23 (ω B97M-D3) and for generating the QM7x data (PBE0+MBD). However, this does not explain the magnitude of the observed errors, since both methods are dispersion-corrected hybrid DFT functionals, which should perform similarly on this data. A more detailed investigation of the errors per molecule reveals that large deviations are exclusively observed for configurations with close interatomic contacts and/or broken covalent bonds. These occur in QM7x, because the non-equilibrium geometries are generated by normal mode sampling in rectilinear coordinates. In contrast, the SPICE set on which MACE-OFF23 is trained uses molecular dynamics (MD) to generate non-equilibrium structures, where close interatomic contacts or broken bonds are highly unlikely. As a consequence, the MACE-OFF23 RMSE is strongly impacted by a small number of outlier structures with unphysical bonding configurations. This becomes apparent when considering the histogram of force errors (figure 2(c)), which reveals that there



is a lower density of errors in the intermediate range (around $1 \text{ eV } \text{Å}^{-1}$) for MACE-OFF23, but a tail of very large errors with low density. In contrast, the distribution of PTBP force errors lacks this tail, highlighting the robustness of this simple physics-based model. Representative configurations, for which MACE-OFF23 displays large errors are shown in figure 2(d). These feature extremely short bonds such as 0.75 Å for a non-covalent H–H pair and 0.67 Å for an N–H bond, respectively. In general, although both models predict reasonable energies and forces for most configurations in QM7x, the electronic structure-based PTBP model demonstrates superior robustness, particularly for non-equilibrium systems. This highlights its suitability as a pre-training low-fidelity method for the MFTL approach.

For initial MFTL tests on QM7x, new MACE models (using the MACE-OFF23 model architecture) were pre-trained on 10k and 50k DFTB(PTBP) datapoints, and subsequently fine-tuned on 0.5k PBE0+MBD datapoints. For robust statistics, each training was repeated three times with randomly initialized weights. In figure 3, the performance of these MFTL models on the PBE0+MBD test set is shown, as a function of the number of epochs used for pre-training. In all cases, errors initially decrease but quickly stagnate and even increase for longer training times. This trend is particularly clear for the larger pre-training set, and when training on forces. Overall, this figure shows that the MLTF concept is sound, since using more low-fidelity DFTB(PTBP) datapoints improves the performance for the high-fidelity test set, even though the size of the high-fidelity dataset is constant. On the other hand, this analysis also provides clear evidence of negative transfer, since training for more epochs increases the error on the high-fidelity data (even when it still decreases the error on the low-fidelity data, see figure S1).

These results may appear somewhat counter-intuitive at first glance, since MFTL appears to benefit from more data but not from more pre-training. However, they can be understood from the perspective of the widely used early-stopping approach for model regularization [47]. Neural networks tend to learn more general (and thus more transferable) concepts and features during early training epochs and more specific details in later epochs. In other words, fully converging the pre-training teaches the model irrelevant (and indeed detrimental) details about the PTBP potential energy surface, which cannot be corrected by the small fine-tuning dataset. Note that we use training epochs as a convenient measure for the length of the training of a given model here. However, this metric is not meaningful when comparing different training set sizes, since the number of weight updates per epoch increases with training set size.

To investigate the influence of the size of the fine-tuning set, MFTL models trained on 50k DFTB(PTBP) samples over 200 epochs were fine-tuned on varying amounts of PBE0+MBD data (see figure 4).

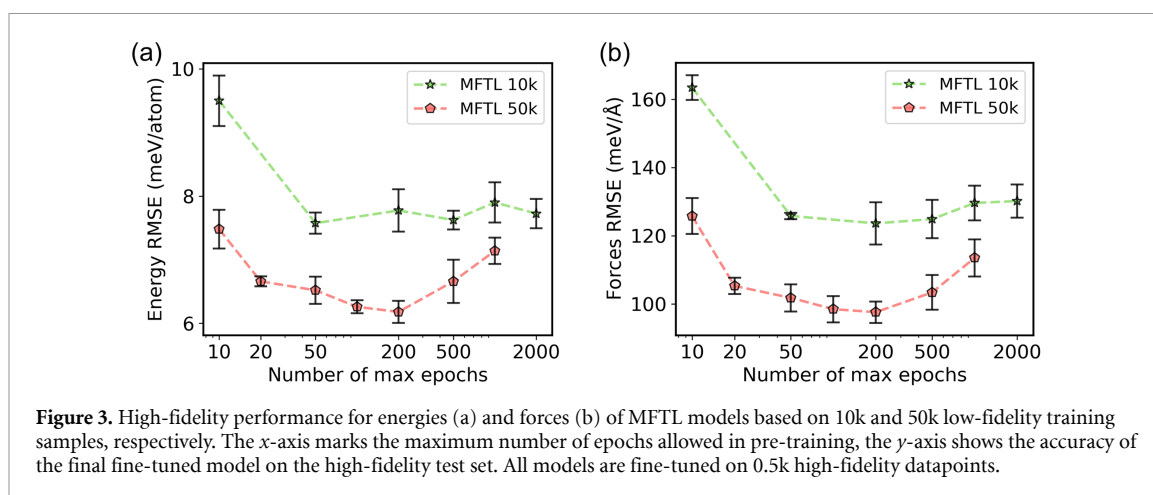


Figure 3. High-fidelity performance for energies (a) and forces (b) of MFTL models based on 10k and 50k low-fidelity training samples, respectively. The x-axis marks the maximum number of epochs allowed in pre-training, the y-axis shows the accuracy of the final fine-tuned model on the high-fidelity test set. All models are fine-tuned on 0.5k high-fidelity datapoints.

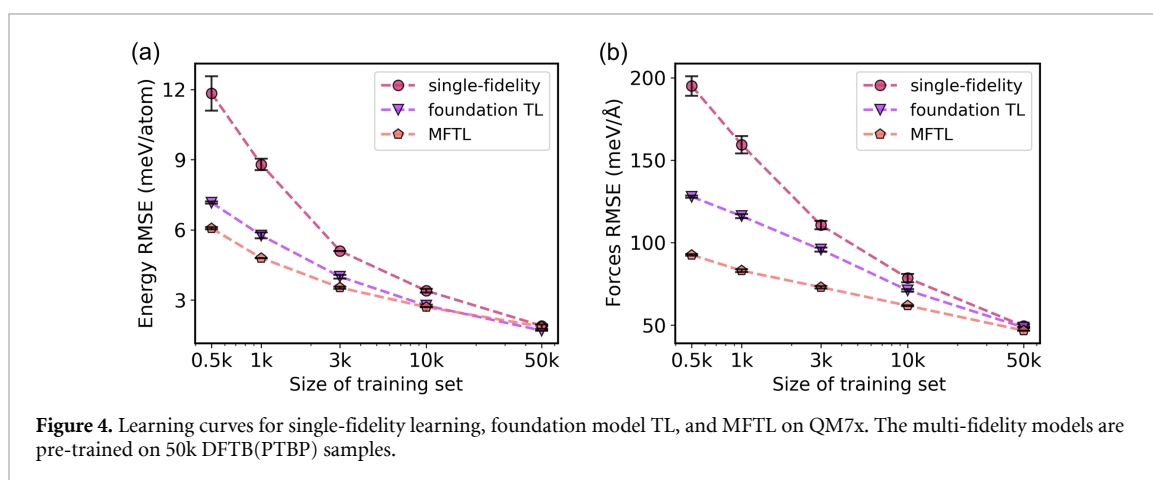


Figure 4. Learning curves for single-fidelity learning, foundation model TL, and MFTL on QM7x. The multi-fidelity models are pre-trained on 50k DFTB(PTBP) samples.

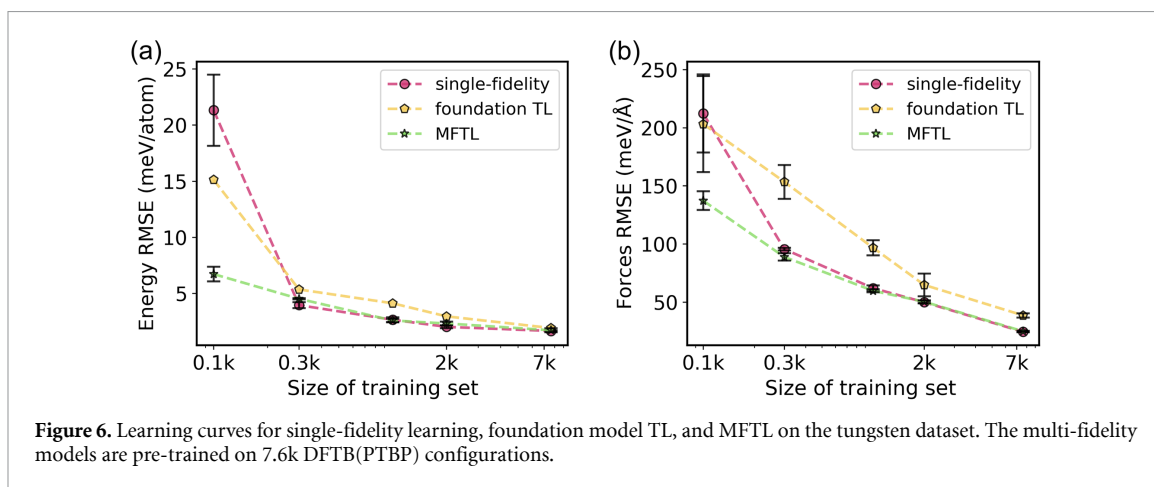
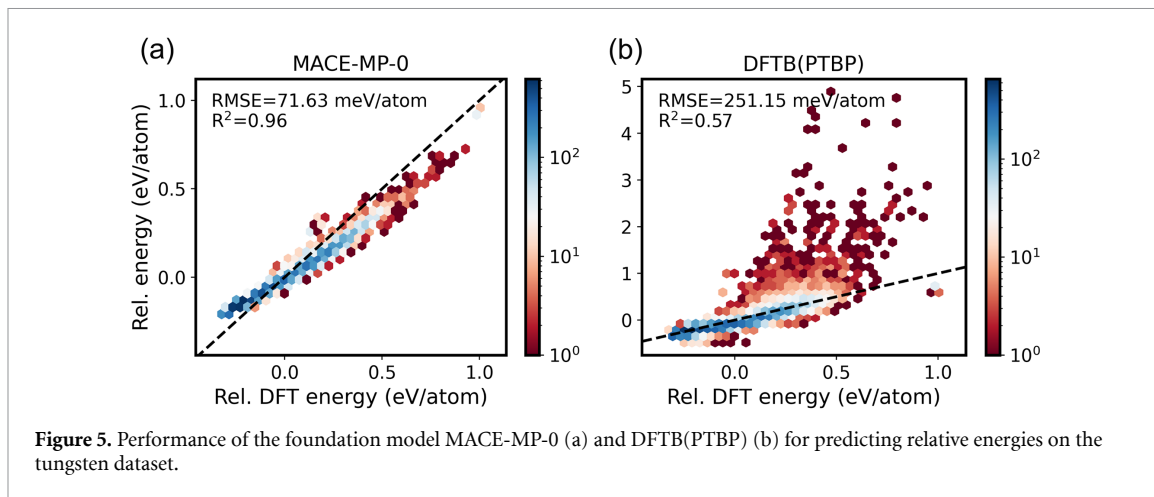
For comparison, we also trained single-fidelity models from scratch on the same data, as well as fine-tuning the MACE-OFF23 foundation model. We find that both TL approaches outperform the single-fidelity model, with the improvements being particularly pronounced for the smallest high-fidelity training set (0.5k configurations). Here, the MFTL model performs best with energy and force RMSEs of $6.1 \text{ meV atom}^{-1}$ and $92.5 \text{ meV \AA}^{-1}$, respectively, compared to $11.8 \text{ meV atom}^{-1}$ and $195.0 \text{ meV \AA}^{-1}$ for single-fidelity learning.

Overall, the differences between the MFTL and fine-tuned foundation models are small but still significant (RMSEs of $7.17 \text{ meV atom}^{-1}$ and $127.87 \text{ meV \AA}^{-1}$ with 0.5k configurations). This shows the benefit of pre-training on configurations directly sampled from the target dataset. While MACE-OFF23 is certainly a better model than PTBP for describing organic molecules in general, the MLFT model benefits from a better description of short interatomic distances. These are important in QM7x but not included in the training of MACE-OFF23. Because of the efficiency of DFTB(PTBP) (and semi-empirical models in general), generating such custom pre-training data only leads to a small computational overhead relative to the generation of high-fidelity data.

Given the good performance of MLFT, it is also worth comparing with Δ -ML [35], which is perhaps the most straightforward multi-fidelity ML approach. Here, instead of targeting the full high-fidelity reference data, the difference between high- and low-fidelity targets is learned. This difference typically displays lower variance and is thus easier to learn than the full high-fidelity label. For QM7x, MLFT and Δ -ML display similar performance (see figure S2). However, Δ -ML has the downside that the low-fidelity method needs to be evaluated for each prediction. While semi-empirical methods are computationally efficient for small molecules, they display less favorable scaling with system size than atomistic ML models. This makes MLFT a more versatile approach than Δ -ML overall.

2.3. Materials

Although QM7x is a highly diverse dataset, small organic molecules are generally a manageable task for ML potentials. This is because of the highly systematic nature of organic chemistry, which can ultimately be reduced to a limited number of atomic environments (functional groups). In contrast, interatomic potentials for materials involving various surfaces, defects, and crystal structures can be more challenging. In particular,



it was recently shown that transition metals display many-body interactions that are difficult to describe with interatomic potentials [48]. We therefore next investigate the performance of MFTL on a dataset containing 1.58k diverse configurations of elemental tungsten (i.e. vacancies, low-index surfaces, gamma-surfaces, and dislocation cores), previously reported in [49]. Energies and forces for this dataset were computed with the PBE functional, which serves as the high-fidelity target.

Since the PTBP model was only fitted to simple crystals, its performance for the tungsten set is rather poor, with some large outliers (see figures 5 and S3), leading to energy and force RMSEs of 251.15 meV and $5.28 \text{ eV } \text{\AA}^{-1}$, respectively. Nonetheless, it provides at least a qualitatively correct baseline in most cases. In contrast, the MACE-MP-0 foundation model performs much better. In terms of energies, non-equilibrium structures are systematically overstabilized, consistent with the previously reported mode-softening of MACE-MP-0 and other foundation models [50]. Nonetheless, the predicted energies and forces show excellent correlation with the reference DFT calculations and no significant outliers.

To develop MFTL models for this dataset, we isolated 1k structures each for validation and testing, and split the remaining dataset into training sets of 0.1k, 0.3k, 1k, 2k, and 7.6k configurations. As for the QM7x dataset, we observe negative transfer that can be mitigated by stopping the pre-training early. Indeed, we find that the best results are observed when stopping after just 6 epochs in this case (see figure S4). This is much earlier than for QM7x, likely due to the fact that the low-fidelity model is significantly noisier in this case. Learning curves for MFTL, single-fidelity models, and the fine-tuned foundation models can be found in figure 6. As above, we find that MFTL is highly beneficial for the smallest training set (100 configurations), with an almost four-fold improvement of the energy RMSE, compared to the single-fidelity model. For larger datasets, the performance of single- and multi-fidelity models is nearly indistinguishable, however.

Interestingly, the foundation model TL scheme is a much smaller improvement over single-fidelity learning for the smallest dataset. In fact, the fine-tuned foundation model performs somewhat worse than the single fidelity model for the larger training sets. This indicates that significant negative transfer is occurring here, despite the good performance of the foundation model. This can be attributed to the fact that the MPTraj dataset used to train MACE-MP-0 only contains very few pure tungsten configurations. Specifically,

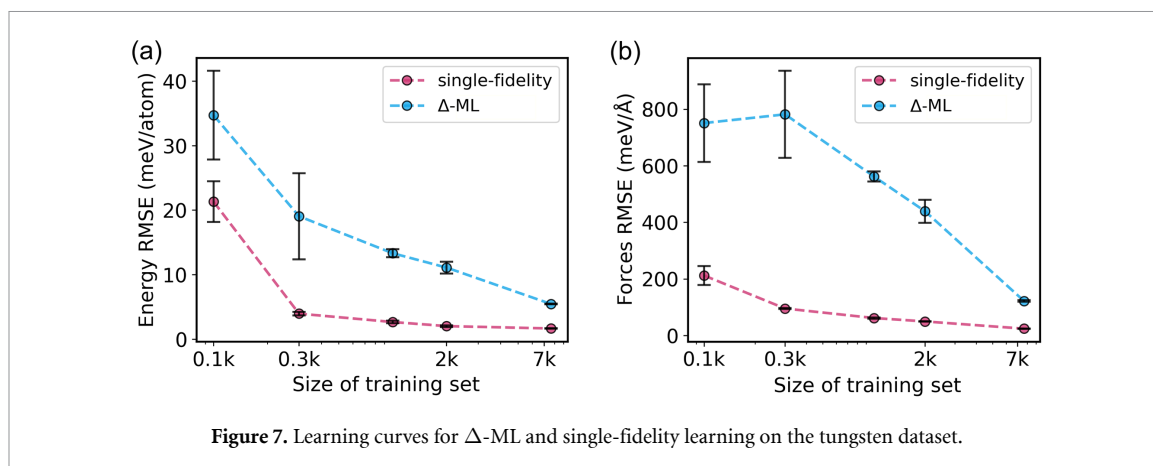


Figure 7. Learning curves for Δ -ML and single-fidelity learning on the tungsten dataset.

the Materials Project only contains eight different tungsten samples, all of which are simple crystals. In contrast, the MFTL model is pre-trained on the full range of atomic environments included in the dataset.

It should be emphasized that all examples discussed up to this point use the simplest TL strategy of retraining all pre-trained weights on the new dataset. For MACE-MP-0, a multi-head TL approach was recently developed, which uses separate read-out heads for the pre-training and fine-tuning data. Additionally, this approach retains a subsample of the pre-training data during the fine-tuning step. This can mitigate both catastrophic forgetting and negative transfer. We also applied this multi-head strategy to the tungsten set, finding much improved results, almost en par with MLFT (see figure S5). This indicates that negative transfer is indeed the likely cause of the discrepancy between MLFT and foundation model fine-tuning. For comparison, Δ -ML models were also developed for this dataset. As shown in figure 7, exceptionally high errors (larger than for single-fidelity models) were observed for these models, however. This can be attributed to the high level of noise in the DFTB(PTBP) baseline data. With early stopping during the pre-training phase, MFTL is nevertheless highly robust, even under these circumstances.

2.4. Multiple fidelities

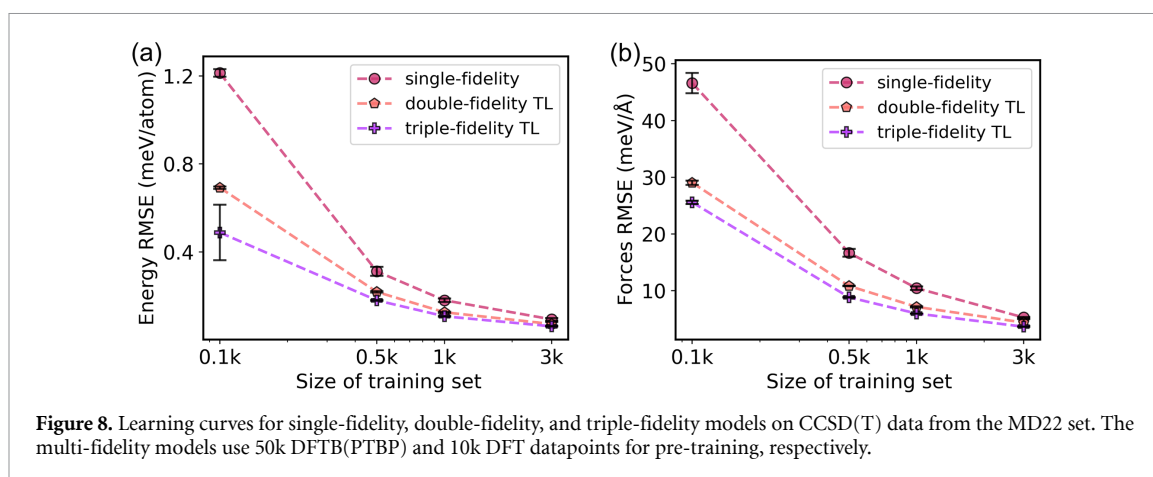
So far, all MFTL examples we discussed used a single low fidelity level (DFTB) and a target high fidelity level (DFT). As demonstrated by von Lilienfeld and co-workers in the Δ -ML context, quantum chemical data is also well suited for developing models with more levels of fidelity, e.g. including highly accurate wavefunction methods like CC theory [51]. To explore this idea further, we used the DFTB(PTBP) pre-trained and PBE0+MBD fine-tuned models developed on the QM7x dataset (see above) as baselines for further TL on CCSD(T) targets from the MD22 dataset [52].

Specifically, we used a set of 4500 non-equilibrium configurations of Benzene, Malonaldehyde, and Toluene (1500 structures for each molecule). These molecules were chosen because the MD22 set provides a set of consistent high fidelity CCSD(T)/cc-pVDZ energies and forces for them. From this combined dataset, we uniformly sampled training sets with 0.1k, 0.5k, 1k, and 3k configurations. Validation and test sets of 0.75k structures each were also generated. We then trained single-fidelity models (trained from scratch on CCSD(T) data), double-fidelity models (TL from DFT to CCSD(T)), and triple-fidelity models (TL from DFTB to DFT to CCSD(T)). Here, 50k DFTB and 10k DFT datapoints from the QM7x dataset were used as the pre-training samples. The corresponding results are shown in figure 8.

This shows that using multiple levels of fidelity is indeed beneficial, as the triple-fidelity model performs best across all training set sizes. Compared to the single-fidelity model, it achieves up to 59.75% and 45.09% improvement in energy and forces error, respectively, yielding RMSEs of $0.49 \text{ meV atom}^{-1}$ and $25.59 \text{ meV \AA}^{-1}$ with only 100 CCSD(T) datapoints. Importantly, it even outperforms the double-fidelity model pre-trained on 10k DFT configurations. This confirms that the extensive amount of DFTB data contains information relevant to the CCSD(T) learning task, beyond what is provided by the DFT pretraining.

3. Conclusion

In this study, we have investigated the properties of MFTL models based on a robust DFTB(PTBP) baseline. By drawing low- and high-fidelity configurations from the same datasets, the effects of structural and fidelity overlap in quantum chemical TL could be disentangled. We find that noise in the low-fidelity labels can be detrimental in both TL and Δ -ML settings. Early-stopping of the pre-training proved to be an efficient way to mitigate this issue, however. With this approach, MFTL even outperforms the straightforward fine-tuning



of high quality foundation models, both for molecular and materials datasets. This is somewhat surprising, since the foundation models are generally more accurate than the DFTB(PTBP) baseline. However, the lower structural overlap between pre-training and fine-tuning datasets causes some negative transfer in the fine-tuning of the foundation models. For the challenging tungsten dataset, we found that this can be mitigated with a more sophisticated multi-head fine-tuning strategy for the foundation model.

More broadly, our results indicate that the current approach of training foundation models to achieve the highest possible accuracy on large single-fidelity databases is non-optimal from the perspective of fine-tuning for specific applications. In future work, multi-fidelity approaches and early-stopping should be investigated for foundation models as well. Inexpensive electronic structure models like DFTB(PTBP) or small basis DFT would allow a massive exploration of materials configuration space [53, 54]. This could increase the applicability and robustness of the next generation of foundation models.

Data availability statement

The data that support the findings of this study are openly available at the following URL/DOI: <https://gitlab.com/mncui/ptbplus.git>.

Acknowledgments

The authors gratefully acknowledge the Max Planck Computing and Data Facility (MPCDF) for providing computing time.

ORCID iDs

Mengnan Cui  <https://orcid.org/0000-0001-6910-2142>

Karsten Reuter  <https://orcid.org/0000-0001-8473-8659>

Johannes T Margraf  <https://orcid.org/0000-0002-0862-5289>

References

- [1] Margraf J, Jung H, Scheurer C and Reuter K 2023 Exploring catalytic reaction networks with machine learning *J. Chem. Theory Comput.* **6** 112–21
- [2] Margraf J 2023 Science-driven atomistic machine learning *Angew. Chem. Int. Ed.* **62** e202219170
- [3] Behler J and Parrinello M 2007 Generalized neural-network representation of high-dimensional potential-energy surfaces *Phys. Rev. Lett.* **98** 1–4
- [4] Handley C, Hawe G, Kell D and Popelier P 2009 Optimal construction of a fast and accurate polarisable water potential based on multipole moments trained by machine learning *Phys. Chem. Chem. Phys.* **11** 6365–76
- [5] Bartók A, Payne M, Kondor R and Csányi G 2010 Gaussian approximation potentials: the accuracy of quantum mechanics, without the electrons *Phys. Rev. Lett.* **104** 136403
- [6] Schütt K, Kindermans P-J, Sauceda H E, Chmiela S, Tkatchenko A and Müller K-R 2017 SchNet: a continuous-filter convolutional neural network for modeling quantum interactions *Advances in Neural Information Processing Systems* vol **30**
- [7] Unke O and Meuwly M 2019 PhysNet: a neural network for predicting energies, forces, dipole moments and partial charges *J. Chem. Theory Comput.* **15** 3678–93
- [8] Gastegger J, Groß J and Günnemann S 2020 Directional message passing for molecular graphs *Int. Conf. on Learning Representations*
- [9] Schütt K, Arbabzadah F, Chmiela S, Müller K and Tkatchenko A 2017 Quantum-chemical insights from deep tensor neural networks *Nat. Commun.* **8** 13890

- [10] Gasteiger J, Becker F and Günnemann S 2021 GemNet: universal directional graph neural networks for molecules *Advances Neural Information Processing Systems* vol 34 pp 6790–802
- [11] Batzner S, Musaelian A, Sun L, Geiger M, Mailoa J, Kornbluth M, Molinari N, Smidt T and Kozinsky B 2022 E(3)-equivariant graph neural networks for data-efficient and accurate interatomic potentials *Nat. Commun.* **13** 2453
- [12] Haghghatdari M et al 2022 NewtonNet: a Newtonian message passing network for deep learning of interatomic potentials and forces *Digit. Discov.* **1** 333–43
- [13] Unke O, Chmiela S, Gastegger M, Schütt K, Sauceda H and Müller K 2021 SpookyNet: learning force fields with electronic degrees of freedom and nonlocal effects *Nat. Commun.* **12** 7273
- [14] Batatia I, Kovacs D, Simm G, Ortner C and Csányi G 2022 MACE: higher order equivariant message passing neural networks for fast and accurate force fields *Advances Neural Information Processing Systems* vol 35 pp 11 423–11 436
- [15] Batatia I and Benner P, 2024 A foundation model for atomistic materials chemistry, (arXiv:2401.00096)
- [16] Gelžinytė E, Ören M, Segall M and Csányi G 2024 Transferable machine learning interatomic potential for bond dissociation energy prediction of drug-like molecules *J. Chem. Theory Comput.* **20** 164–77
- [17] Pracht P, Pillai Y, Kapil V, Csányi G, Gönnheimer N, Vondrák M, Margraf J and Wales D 2024 Efficient composite infrared spectroscopy: combining the double-harmonic approximation with machine learning potentials *J. Chem. Theory Comput.* **20** 10986–004
- [18] Zhang D, Bi H, Dai F, Jiang W, Liu X, Zhang L and Wang H 2024 Pretraining of attention-based deep learning potential model for molecular simulation *npj Comput. Mater.* **10** 1–8
- [19] Zhouyin Z, Gan Z, Pandey S K, Zhang L and Gu Q 2025 Learning local equivariant representations for quantum operators *The 13th Int. Conf. on Learning Representations*
- [20] Stocker S, Gasteiger J, Becker F, Günnemann S and Margraf J 2022 How robust are modern graph neural network potentials in long and hot molecular dynamics simulations? *Mach. Learn.: Sci. Technol.* **3** 045010
- [21] Pan S and Yang Q 2010 A survey on transfer learning *IEEE Trans. Knowl. Data Eng.* **22** 1345–59
- [22] Käser S, Itza Vazquez-Salazar L, Meuwly M and Töpfer K 2023 Neural network potentials for chemistry: concepts, applications and prospects *Digit. Discov.* **2** 28–58
- [23] Zhuang F, Qi Z, Duan K, Xi D, Zhu Y, Zhu H, Xiong H and He Q 2021 A comprehensive survey on transfer learning *Proc. IEEE* **109** 43–76
- [24] Rowe P, Deringer D, Gasparotto P, Csányi G and Michaelides A 2022 Erratum: an accurate and transferable machine learning potential for carbon *J. Chem. Phys.* **156** 2020–2
- [25] Chen C and Ong S 2021 AtomSets as a hierarchical transfer learning framework for small and large materials datasets *npj Comput. Mater.* **7** 1–9
- [26] Zhang W, Deng L, Zhang L and Wu D 2023 A survey on negative transfer *IEEE/CAA J. Autom. Sinica* **10** 305–29
- [27] Hutchinson M, Antono E, Gibbons B, Paradiso S, Ling J and Meredig B 2017 Overcoming data scarcity with transfer learning (arXiv:1711.05099)
- [28] Chen C, Ye W, Zuo Y, Zheng C and Ong S 2019 Graph networks as a universal machine learning framework for molecules and crystals *Chem. Mater.* **31** 3564–72
- [29] Frey N, Akinwande D, Jariwala D and Shenoy V 2020 Machine learning-enabled design of point defects in 2D materials for quantum and neuromorphic information processing *ACS Nano* **14** 13406–17
- [30] Devereux C, Smith J, Huddleston K, Barros K, Zubatyuk R, Isayev O and Roitberg A 2020 Extending the applicability of the ANI deep learning molecular potential to sulfur and halogens *J. Chem. Theory Comput.* **16** 4192–202
- [31] Fitzgerald T and Thomaz A L 2015 Skill demonstration transfer for learning from demonstration *10th Annual ACM/IEEE Int. Conf. on Human-Robot Interaction Extended Abstracts (2 March 2015)* pp 187–8
- [32] Kovács D, Moore J, Browning N, Batatia I, Horton J, Kapil V, Witt W, Magdău I, Cole D and Csányi G 2023 MACE-OFF23: transferable machine learning force fields for organic molecules (arXiv:2312.15211)
- [33] Batra R, Pilania G, Uberuaga B and Ramprasad R 2019 Multifidelity information fusion with machine learning: a case study of dopant formation energies in hafnia *ACS Appl. Mater. Interfaces* **11** 24 906–24 918
- [34] Goodlett S, Turney J, Schaefer H I 2023 Comparison of multifidelity machine learning models for potential energy surfaces *Int. J. Chem. Phys.* **159** 044111
- [35] Ramakrishnan R, Dral P, Rupp M and von Lilienfeld O 2015 Big data meets quantum chemistry approximations: the delta-machine learning approach *J. Chem. Theory Comput.* **11** 2087–96
- [36] Wengert S, Csányi G, Reuter K and Margraf J 2022 A hybrid machine learning approach for structure stability prediction in molecular co-crystal screenings *J. Chem. Theory Comput.* **18** 4586–93
- [37] Fare C, Fenner P, Benatan M, Varsi A and Pyzer-Knapp E 2022 A multi-fidelity machine learning approach to high throughput materials screening *npj Comput. Mater.* **8** 1–9
- [38] Buterez D, Janet J, Kiddle S, Oglic D and Lió P 2024 Transfer learning with graph neural networks for improved molecular property prediction in the multi-fidelity setting *Nat. Commun.* **15** 1517
- [39] Chen K, Kunkel C, Cheng B, Reuter K and Margraf J T 2023 Physics-inspired machine learning of localized intensive properties *Chem. Sci.* **14** 4913–22
- [40] Allen A E, Lubbers N, Matin S, Smith J, Messerly R, Tretiak S and Barros K 2024 Learning together: towards foundation models for machine learning interatomic potentials with meta-learning *npj Comput. Mater.* **10** 145301
- [41] Jain A et al 2013 Commentary: the materials project: a materials genome approach to accelerating materials innovation *APL Mater.* **1** 011002
- [42] Eastman P et al 2023 SPICE, a dataset of drug-like molecules and peptides for training machine learning potentials *Sci. Data* **10** 11
- [43] Zaspel P, Huang B, Harbrecht H and von Lilienfeld O 2019 Boosting quantum machine learning models with a multilevel combination technique: Pople diagrams revisited *J. Chem. Theory Comput.* **15** 1546–59
- [44] Hoja J, Medrano Sandonas L, Ernst B, Vazquez-Mayagoitia A, DiStasio R and Tkatchenko A 2021 QM7-x, a comprehensive dataset of quantum-mechanical properties spanning the chemical space of small organic molecules *Sci. Data* **8** 43
- [45] Adamo C and Barone V 1999 Toward reliable density functional methods without adjustable parameters: the PBE model *Int. J. Chem. Phys.* **110** 6158–70
- [46] Tkatchenko A, DiStasio R, Car R and Scheffler M 2012 Accurate and efficient method for many-body van der Waals interactions *Phys. Rev. Lett.* **108** 236402
- [47] Caruana R, Lawrence S and Giles C 2000 Overfitting in neural nets: backpropagation, conjugate gradient and early stopping *Advances in Neural Information Processing Systems* vol 13

- [48] Owen C, Torrisi S, Xie Y, Batzner S, Bystrom K, Coulter J, Musaelian A, Sun L and Kozinsky B 2024 Complexity of many-body interactions in transition metals via machine-learned force fields from the TM23 data set *npj Comput. Mater.* **10** 1–16
- [49] Szlachta W, Bartók A and Csányi G 2014 Accuracy and transferability of Gaussian approximation potential models for tungsten *Phys. Rev. B* **90** 104108
- [50] Deng B, Choi Y, Zhong P, Riebesell J, Anand S, Li Z, Jun K, Persson K A and Ceder G 2025 Systematic softening in universal machine learning interatomic potentials *npj Comput. Mater.* **11** 47
- [51] Zaspel P, Huang B, Harbrecht H and von Lilienfeld O A 2019 Boosting quantum machine learning models with a multilevel combination technique: Pople diagrams revisited *J. Chem. Theory Comput.* **15** 1546–59
- [52] Chmiela S, Vassilev-Galindo V, Unke O T, Kabylda A, Sauceda H E, Tkatchenko A and Müller K-R 2023 Accurate global machine learning force fields for molecules with hundreds of atoms *Sci. Adv.* **9** 1875
- [53] Cui M, Reuter K and Margraf J 2024 Obtaining robust density functional tight-binding parameters for solids across the periodic table *J. Chem. Theory Comput.* **20** 5276–90
- [54] Keller E, Morgenstein J, Reuter K and Margraf J 2024 Small basis set density functional theory method for cost-efficient, large-scale condensed matter simulations *Int. J. Chem. Phys.* **161** 074104