# Descriptor-based Foundation Models for Molecular Property Prediction

Jackson W. Burns[1][†], Akshat S. Zalte[1][†], William H. Green[1][*]

[1]Department of Chemical Engineering, MIT, Cambridge, MA.

*Corresponding author(s). E-mail(s): whgreen@mit.edu;
[†]These authors contributed equally to this work.

## Abstract

Fast and accurate prediction of molecular properties with machine learning is pivotal to scientific advancements across myriad domains. Foundation models in particular have proven especially effective, enabling accurate training on small, real-world datasets. This study introduces **CheMeleon**, a novel molecular foundation model pre-trained on deterministic molecular descriptors from the Mordred package, leveraging a Directed Message-Passing Neural Network to predict these descriptors in a noise-free setting. Unlike conventional approaches relying on noisy experimental data or biased quantum mechanical simulations, **CheMeleon** uses low-noise molecular descriptors to learn rich molecular representations. Evaluated on 58 benchmark datasets from Polaris and MoleculeACE, **CheMeleon** achieves a win rate of 79% on Polaris tasks, outperforming baselines like Random Forest (46%), **fastprop** (39%), and **Chemprop** (36%), and a 97% win rate on MoleculeACE assays, surpassing Random Forest (63%) and other foundation models. However, it struggles to distinguish activity cliffs like many of the tested models. The t-SNE projection of **CheMeleon**'s learned representations demonstrates effective separation of chemical series, highlighting its ability to capture structural nuances. These results underscore the potential of descriptor-based pre-training for scalable and effective molecular property prediction, opening avenues for further exploration of descriptor sets and unlabeled datasets.

**Keywords:** foundation models, descriptors, chemistry

1

# 1 Main

Fast and accurate prediction of molecular properties is pivotal to scientific advancements in drug discovery, material design, and catalysis. In recent years, machine learning has emerged as a transformative tool for molecular property prediction, offering rapid and scalable alternatives to resource-intensive experiments and quantum chemistry calculations. These data-driven models have demonstrated remarkable success in predicting critical properties, such as Absorption, Distribution, Metabolism, and Excretion (ADME) properties [1] and physicochemical properties like lipophilicity and solubility [2, 3], leveraging diverse architectures and molecular representations to accelerate discovery across myriad domains.

The development of machine learning models relies on a diverse set of architectures, each tailored to balance predictive power and data efficiency. Widely accessible through user-friendly Python packages, representative models include the ubiquitous Random Forest[4], **fastprop** [5], and **Chemprop** [2, 6]. All these methods employ direct supervised fitting, where randomly initialized parameters are trained solely on the labeled target dataset. Random Forest (RF), an ensemble of decision trees trained on molecular fingerprints (e.g., Morgan Fingerprints [7]) or descriptors, serves as a robust baseline for its effectiveness on small datasets [4]. Similarly, **fastprop**, a feedforward neural network using fixed molecular descriptors like topological or physicochemical features, achieves strong performance on small datasets but has a lower performance ceiling due to its reliance on predefined representations [5]. In contrast, **Chemprop**, a graph neural network based on a Directed Message-Passing Neural Network (D-MPNN), excels by learning optimal feature representations from molecular graphs during training, offering generality across a wide range of tasks like solubility [8], radical thermochemistry [9], and infrared spectra [10] prediction, among others.

Studies have shown that models leveraging learned representation (LR) tend to outperform those using fixed molecular representations [2, 11, 12]. However, in the limit of small datasets ($\lesssim \mathcal{O}(1{,}000)$ samples), **Chemprop** and other LR models struggle, often outperformed by classical methods like Random Forest [13]. Learning molecular representations from scratch using deep neural networks, such as in **Chemprop**, requires large datasets because the initial input features are sparse and uninformative [14]. This forces the model to simultaneously learn both a suitable representation and the target property mapping, which is challenging in low-data regimes and leads to poor generalization and overfitting.

To unlock the performance potential of LR models without suffering these limitations, recent work has focused on learning meaningful, general representations using molecular foundation models. These models are typically pre-trained on large unlabeled datasets via self-supervised learning and aim to provide rich, transferable features that can be fine-tuned for various downstream tasks. We discuss several foundation models for chemistry in the following section and refer readers to the comprehensive review by Choi et al. [15] for a more in-depth treatment of the literature.

A wide range of approaches has been developed to create molecular foundation models, differing in both molecular representations and training strategies. Inspired by advances in natural language processing (NLP), many efforts rely on text-based representations such as SMILES (Simplified Molecular Input Line Entry System) [16, 17]

and SELFIES [18], which naturally lend themselves to transformer-based architectures. Notable examples include ChemBERTa-2 [19] and MolFormer [20], both of which are pre-trained on $\mathcal{O}(10M)$ or more SMILES strings to learn general-purpose molecular representations. More recent work has shifted toward graph-based representations, which enable more parameter-efficient models by naturally capturing molecular structure that is absent in SMILES. [21, 22] These methods have evolved through a range of training strategies, broadly categorized as supervised or unsupervised based on standard machine learning conventions. Among unsupervised approaches, **MolCLR** [23] employs contrastive learning with graph neural networks, including graph convolution and graph isomorphism networks. Another model, GROVER [24], combines graph neural networks with Transformer-style attention and is trained using contextual property and motif prediction tasks.

The best-performing models in the current literature have extended graph-based learning approaches using supervised pre-training. All supervised pre-training methods broadly follow the workflow illustrated in Figure 4. A molecular representation method first embeds the input molecule into a numerical vector, which is then passed through a regressor to predict a set of target labels. The choice of these labels distinguishes different foundation model strategies.

State-of-the-art pre-training labels, such as experimental results and computed quantum mechanics properties, dominate this space. Models such as **MolE** [25] and the work of Beaini et al. adopt this strategy, reporting excellent results across many common benchmarks. While these labels are typically of high quality, their availability is limited, especially compared to the large unlabeled datasets leveraged in unsupervised learning. To attain the scale required for generalizable models, experimental datasets must be combined. This necessarily results in sparsity but also likely large stochastic and even systematic error from inter-laboratory variation. For instance, combining identical assays from different labs often results in low correlation coefficients, as shown by Landrum and Riniker. This poses a challenge for modeling as most methods struggle to regress in the presence of systematic error or implicitly rely on error being stochastic.

An alternative strategy is to use only labels generated from quantum mechanical (QM) simulations. Prominent models like **MolGPS** [28] and **GraphQPT** [29] are trained on such data. Although simulated datasets can be significantly larger, QM methods are often parameterized for narrow chemical subspaces, limiting their transferability. Moreover, simulation outputs can carry non-random, method-specific biases that are difficult to quantify or correct.

Both of these methods suffer from the aforementioned presence of systematic error and sparsity, hampering scaling and downstream model performance. In this work, we seek to address these shortcomings with **CheMeleon**, a D-MPNN foundation model pre-trained to predict Mordred precomputed descriptors for $\mathcal{O}(1M)$ molecules from PubChem [30]. These descriptors are calculated directly from the molecular graph using counting, aggregation, and complexity algorithms, occasionally using external data such as measured atomic volumes. As a result, **CheMeleon** exploits the low-noise signal provided by Mordred to learn a general and highly effective representation.

We demonstrate that fine-tuning **CheMeleon** across 58 different benchmark sets targeting solubility, lipophilicity, and activity yields consistent, statistically significant improvements over baseline **Chemprop** models. Further, **CheMeleon** outperforms baselines such as Random Forest and other foundation models while providing chemically meaningful embeddings. **CheMeleon** is open source, permissively licensed, and seamlessly integrated into the **Chemprop** software package for easy reuse.

## 2 Results

Herein we present summary statistics and selected results following the aforementioned experimental procedure. Complete results are available in Appendix B and in the corresponding source code (see Section 3.1). Although the selection of results shown here and the accompanying commentary are representative of the entire set.

### 2.1 Pre-training

Descriptors were calculated with Mordred and saved to disk using the Zarr format [31] to allow fast reloading and decrease training time. This format is especially appropriate for the dense target data generated by Mordred, characteristic of using molecular descriptors. **CheMeleon** was then constructed using **Chemprop**, having 6 hidden layers of dimension 2048 in the D-MPNN and a three-layer FNN of the same size, and then trained using PyTorch Lightning [32]. After training, **CheMeleon** achieved a Root Mean Squared Error of 0.14 averaged across all of the rescaled ($\mu$ of zero with $\sigma$ of one) and Winsorized ($\sigma$ of six) descriptors. This number is mentioned for completeness rather than its importance; the model is never actually deployed to calculate these descriptors, especially given that they can be calculated exactly using Mordred.

### 2.2 Polaris Benchmarks

Figure 1 shows the results for a select set of benchmarks covering especially relevant tasks. As shown, in many cases there are multiple winning models or *all* models achieve statistically indistinguishable performance, generally including **CheMeleon**. See Appendix B1 for a complete diagram showing every single benchmark run, the results of which are summarized in 1. **CheMeleon** narrowly outperforms the next best model **minimol** on this subset of benchmarks, achieving a Win Rate of 79%. This demonstrates the effectiveness of molecular descriptors for pretraining. Without the need for combining noisy experimental assays or running expensive QM simulations, **CheMeleon** is able to achieve state-of-the-art performance. Particularly noteworthy here is the performance of **Chemprop**, falling below all tested methods, including simple baselines like Random Forest. **CheMeleon**, which is derived from **Chemprop**, is able to dramatically improve this performance.

### 2.3 MoleculeACE Benchmarks

Figure 2 shows a selected subset of the tested models' RMSE between activity cliff subgroups across all assays in the MoleculeACE study. This figure is inspired by that in the original study but modified to show both the per-model statistical tests for
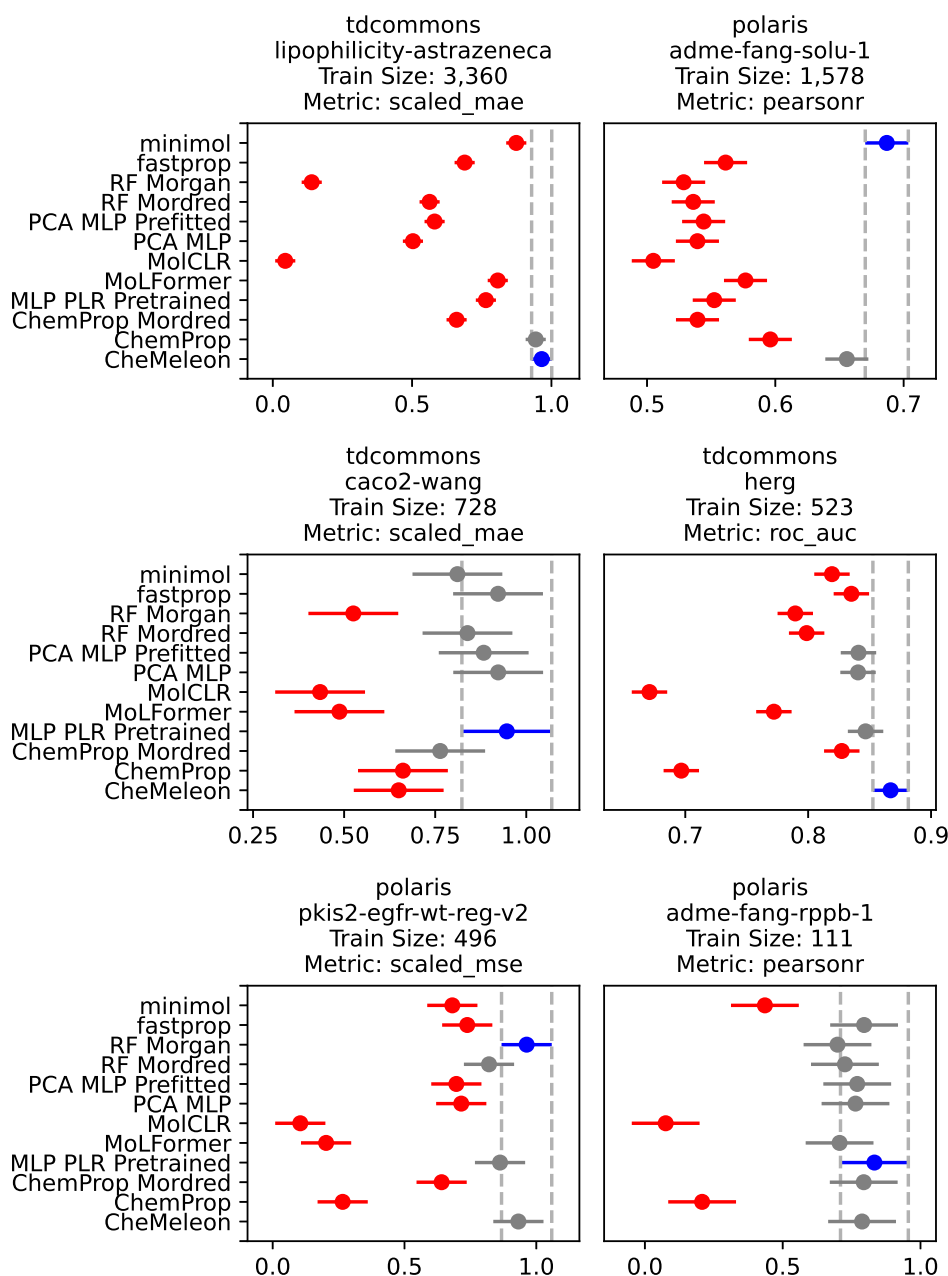
**Fig. 1**: Performance of each of the tested models across a set of different common molecular machine learning tasks. The origin of each benchmark set is shown as the first line of each subplot title, followed by the name of the dataset (which indicates the task), the size of the training data, and the metric used to evaluate model performance. Benchmarks are sorted in order of training size, decreasing. Models shown in blue are the absolute highest performers on the given benchmark, while models shown in gray are not practically different from the best performer according to the Tukey Honestly Significant Difference test based on the variance in test set performance across five repetitions, as laid out in Section 3.1. Models show in red *are* practically worse performers and are considered to have "lost" on the indicated benchmark.

**Table 1**: Model performance comparison on Polaris benchmarks.

| Model | Win Count | Win Rate (%) |
|---|---|---|
| **CheMeleon** | 22 | 79 |
| **minimol** | 20 | 71 |
| **MLP-PLR** Pre-fitted | 14 | 50 |
| RF Mordred | 13 | 46 |
| RF Morgan | 12 | 43 |
| **PCA MLP** | 11 | 39 |
| **fastprop** | 11 | 39 |
| **PCA MLP** Pre-fitted | 11 | 39 |
| **Chemprop** | 10 | 36 |
| **Chemprop**-Mordred | 9 | 32 |
| **MoLFormer** | 9 | 32 |
| **MolCLR** | 6 | 21 |

consistency and the intra-model test for absolute performance. The results for all models are shown in Appendix B and summary statistics for all models are presented in Table 2.

All of the tested models consistently fail to achieve practically identical performance between the cliff and noncliff compounds while also performing the best (filled blue markers). Across the tested benchmarks, **CheMeleon** does so only four times, which does at least outperform the next best model (Random Forest) with only two such cases. **CheMeleon** is statistically the best on 29 of the 30 benchmarks, making it the best performer in absolute terms. Notable for its consistency but inaccuracy in this benchmark is **Chemprop**. Although many of its differences are practically indistinguishable, the model achieves the highest performance in *zero* of the test benchmarks. Its apparent high consistency is an artifact of underfitting, which means that the model is equally bad at predicting cliff and noncliff compounds' activity.

The apparent difficulty of this benchmark for all tested methods suggests that foundation modeling, while able to improve the absolute performance, cannot discern the presence of such cliffs. As shown in Section 2.4, **CheMeleon** learns a representation that maps structurally similar molecules into the same feature space, hampering efforts to distinguish them during fine-tuning. Given the poor performance of the other tested foundation models **MoLFormer** and **minimol**, this challenge seems to be endemic to the field.

## 2.4 Foundation Fingerprint

Adapting the provided source code from the O'Boyle and Sayle study, we generated the LR for each of the molecules provided in the Single Assay benchmark. We calculated the cosine distance between the embeddings within each series to arrive at a **CheMeleon**-based sorting order. This was compared with the sort order derived from Atom Pair Fingerprints [35], Topological-torsion fingerprints [36], Morgan [7], RDKit [37] fingerprint, Avalon [38] fingerprint, and MACCS Keys [39] as they collectively make up the standard set of molecular fingerprints used for molecular representation, such as for the Random Forest model used here. The same **CheMeleon** fingerprints

**Table 2**: Model performance comparison on MoleculeACE benchmarks.

| Model | Win Count | Win Rate (%) |
|---|---|---|
| **CheMeleon** | 29 | 97 |
| RF Morgan | 19 | 63 |
| **minimol** | 13 | 43 |
| RF Mordred | 8 | 27 |
| **fastprop** | 5 | 17 |
| **Chemprop**-Mordred | 4 | 13 |
| **MLP-PLR** Pre-trained | 4 | 13 |
| **PCA MLP** Pre-fitted | 3 | 10 |
| **MoLFormer** | 3 | 10 |
| **Chemprop** | 0 | 0 |

are then subject to the methods of Orlov et al. to generate a two-dimensional projection and a few of the series present in the benchmark are highlighted.

The t-SNE projection shows that **CheMeleon** can separate large and small molecules in its feature space and that the three highlighted assays are far apart relative to each other. Each assay also provides further chemical insight. On top, all five species are small modifications of the original structure at the same site and are thus all projected into similar locations in chemical space. The left assay shows that modifications removing the methylenedioxy group cause two subsets of the series to be projected into *different* regions of chemical space. The final assay (shown at the bottom) follows a similar trend, in this case with all four compounds following the lead compound being projected into a different region of chemical space than the lead compound due to the initial significant structural modification.

We believe that perhaps the most profound conclusion of this study is that simple molecular descriptors can be used for effective pre-training of deep message-passing neural networks. Systematic error-ridden experimental datasets, combined from those scraped from the literature, are not required, nor are expensive QM simulations. One can achieve effective pre-training with readily calculated and understood quantities instead.

We hope that this will open an entirely new line of investigation in terms of *which* set of molecular descriptors is the most effective for this task. For the present study, we focus on the Mordred descriptors in particular due to their ease of interoperation with code, but any set of descriptors such as DRAGON [42] or PaDEL [43] can be used. The choice of unlabeled data interacts with the selection of descriptors as well. Other unlabeled datasets besides PubChem, which are more relevant to other domains, could be used, such as ChEMBL for drug-like molecules [44], the Collection of Open Natural Products database [45] for natural products. The appropriate selection of both datasets and their corresponding descriptors likely requires significant domain expertise but may yield substantial rewards.

# 3 Methods

Figure 4 visualizes the workflow for the present study. Each descriptor represents some algorithm operating on the molecular graph. These can be as simple as counting the presence of common functional groups and as complex as performing repeated shortest-path graph navigation, as in the Wiener index [46]. There are even 'parameterized' descriptors based on externally estimated sub-properties, such as atomic volume, as used in the molecular McGowan Volume descriptor [47]. Thus, during pre-training **Chemprop** 'learns' to calculate descriptors in the FNN, and we reuse the LR during fine-tuning. See Appendix A for more information about the implementation of model training.

Given that the descriptors being predicted as part of this workflow can themselves be used as molecular embeddings, such as with **fastprop** [5], we also sought to train a model that is trained to ingest this embedding and reproduce it. This is analogous to other SMILES-based foundation models such as ChemBERTa [19], which ingest SMILES strings and are pre-trained to reproduce them using Masked Language Modeling. To that end, we use the **MLP-PLR** as described by Gorishniy et al. as an autoencoder during pre-training. This is because it shows state-of-the-art performance relative to MLPs and is competitive with the more complex transformer methods. Also important is that the initial PLR embedding step of the input increases the dimension of the latent representation without allowing for a direct copy operation, preventing the model from learning the identity function during pre-training. See Appendix A for more information.

Finally, as another baseline, we use `scikit-learn`'s Principal Component Analysis (PCA) [41] in combination with a Multi Layer Perceptron (MLP), dubbed **PCA MLP**, to first project the many features into a smaller dimension and then regress the projected hyper-descriptors. The key idea with this method is that, given the high dimensionality of the Mordred feature set $(1,613)$ the use of dimensionality reduction will decrease the required number of features in the downstream model and thus improve performance. This method can be used both on the pre-calculated large set of descriptors (**PCA MLP** Pre-fitted) or on the fine-tuning dataset only (**PCA MLP**); both are done here.

For details about the comparison models, see the associated publications and Appendix C.

## 3.1 Benchmarks

We run the present models and all models referenced in C on a collection of 58 unique benchmark datasets, to be described in the following paragraph. Across all benchmarks, we only perform single-task regression or classification. Although multi-task has been shown to improve performance on chemical datasets [50], we seek to demonstrate the performance on individual tasks to facilitate comparisons better; all of the present models *can* be used in a multitask fitting approach. The present procedure enables us to quantify the performance of the models subject to the most significant source of uncontrollable, stochastic error. Towards that same end, we do not perform

hyperparameter optimization for any model on any benchmark dataset. Such a procedure is known to lead to overfitting [51], particularly on the small datasets used in this study.

For all benchmarks in this paper, we adopt a single test set as defined by the original creators of each dataset. We do this partly in deference to previous studies, facilitating easy comparison to other models, but also to enable rigorous statistical comparisons between our tested models. To generate the said statistics, we adopt the following procedure to obtain replicate measurements for each benchmark.

1. Select a random seed for the current repetition
2. Based on the model architecture, use this random seed to:
   (a) (for deep models) select the validation set, used for early stopping to detect overfitting
   (b) (for random forest) set the initial feature embeddings
3. Train the model using the training data, using validation for early stopping where applicable
4. Apply the trained model to the testing data Described the statistical tests in much greater detail.

The Polaris benchmarking suite is a collection of expert-curated datasets covering a broad range of molecular property prediction targets [52]. Users have also made available the datasets from the Therapeutic Data Commons (TDC) [53]. For this study, we selected a set of 28 unique benchmarks covering relevant targets such as solubility, physiology, and biophysics from both Polaris and TDC. A subset of particularly relevant benchmarks whose performance across models is representative of the entire set has been visualized in 1. Each benchmark uses a different metric, automatically calculated by Polaris, all of which are either naturally scaled or rescaled (as done here) to range from zero to one, with one being the best.

The remaining 29 benchmarks are a set curated in the MoleculeACE study [33], each of which is an assay measuring the activity of small molecules against a biologically relevant target. The metric of interest in each benchmark is the difference in the Root Mean Squared Error (RMSE) of the predictions between "cliff" and "noncliff" molecules. van Tilborg et al. provide a thorough definition of what this precisely entails. In summary, the molecules are partitioned into training and testing such that a series of structurally similar molecules ("noncliff molecules"), where some members have dramatically different activity ("cliff molecules"), are placed in full into the two sets. This allows calculation of the RMSE of test set predictions for the two subgroups, used for absolute performance comparisons, and then the difference in RMSE. For the present study, we perform a one-sided, one-sample t-test on this quantity to determine if the model performance is consistent between the two subgroups. This is then summarized across all assays as the "consistency rate", a percentage reflecting the frequency with which the model maintains the same performance on cliff compounds as noncliff compounds.

Although not central to the present work, we have provided performance results on the MoleculeNet benchmark collection [54] to facilitate comparisons with historical models (see Appendix B.3). As discussed at length by Walters, many of the

benchmarks included in this collection cover irrelevant or poorly defined tasks and contain data curation errors. The physical chemistry, biophysics, and physiology categories from MoleculeNet are still represented in our selected benchmarks, but with the introduction of Polaris, they are now more curated, standardized, and easier to access.

For each benchmark, we perform the Tukey Honestly Significant Difference (HSD) test as implemented in `SciPy` [56] to compare all of the tested models simultaneously as suggested by Polaris authors [52]. For each benchmark, the best-performing model and all models with performance that is statistically indistinguishable from the best performer at 95% confidence are considered winners. This is aggregated across all benchmarks to arrive at a win count and win rate.

Finally, we sought a non-modeling-based evaluation method to provide foundational insights into the learned representation of **CheMeleon**. To that end, we reproduce the "Single Assay" benchmark from O'Boyle and Sayle, which enables quantitative comparisons of arbitrary molecular fingerprints' capacity to reproduce human-defined sorting order across thousands of series of small molecules. To further visualize the result of this benchmark, we adopt the methods of Orlov et al. to visualize how the **CheMeleon** LR projects chemical series into its feature space. This provides insights into how the model has learned to separate common chemical moieties.

## Software and Data

The source code associated with this study, including pre-training, fine-tuning, numerical results, and visualization of the same, is available at GitHub.com/JacksonBurns/CheMeleon.

All datasets used in this work are available from their respective authors via Polaris or GitHub.

## Acknowledgments

## Appendix A    Training Details

All models were trained using PyTorch Lightning [32] and associated open source machine learning and cheminformatics packages. We used eight Nvidia 2080 Ti's for foundation model training and a Nvidia Quadro RTX 4000 laptop GPU for fine-tuning. The architecture for the **MLP-PLR** was determined via modest hyperparameter optimization using conventional Python toolkits. The architecture for **CheMeleon** was

as described in the main text and was arrived at without any hyperparameter tuning. Every model tested during fine-tuning had the following architecture, which was decided in deference to the suggested fine-tuning architecture from each of the original authors, where such an indication was given and held consistent otherwise:

- Random Forest: default settings as set by `scikit-learn`
- **PCA MLP**: number of hyper-descriptors (components in PCA) set to account for 95% of variance in the original data, fine-tuning with MLP used two layers of size 1,800.
- **MoLFormer**: tuned entire network with a single additional readout layer
- **MLP-PLR**: optimized model used an embedding dimension of 8 with a hidden size of 4096 with 3 layers, fine-tuning head used two layers of size 256
- **fastprop**: default settings as set by Burns and Green
- **Chemprop**: default setting as set by Heid et al.
- **minimol**: multi-layer skip connection MLP with regularization, hidden size of 512 as described by Kläser et al.

Note that we use mordred-community, a community-maintained fork of the original Mordred descriptor calculator, which is no longer maintained.

# Appendix B    Detailed Experimental Results

## B.1    Polaris Benchmarks

The complete set of Honestly Significant Difference diagrams following the conventions described in the main text (Figure 1) is shown in Figure B1. A high-quality version of this image is also available in the GitHub repository described in Section 3.1.

## B.2    MoleculeACE Benchmarks

The complete set of Honestly Significant Difference diagrams following the conventions described in the main text (Figure 1) is shown in Figure B2. A high-quality version of this image is also available in the GitHub repository described in Section 3.1.

The complete set of MoleculeACE-style diagrams following the conventions described in the main text (Figure 2) is shown in Figure B3. A high-quality version of this image is also available in the GitHub repository described in Section 3.1.

## B.3    Chemprop Benchmarks

This section extends our evaluation of **CheMeleon** to include all molecular property prediction benchmarks originally used to assess **Chemprop** [6]. We follow the same protocols, datasets, and splits to ensure direct comparability. These benchmarks include a mix of regression tasks (e.g., UV/Vis, SAMPL, QM9, PCQM4Mv2) and classification tasks (e.g., HIV, PCBA). It is important to note that **CheMeleon** is used out-of-the-box, while the **Chemprop** results are based on tuned hyperparameters.

**Table B1**: Comparison of model performance on the test set, trained with **CheMeleon** and **Chemprop**, for regression tasks, including UV/Vis peak absorption wavelength, logP for the SAMPL challenge (SAMPL 6, 7, and 9), and HOMO-LUMO gap for PCQM4Mv2 dataset. It is important to note that the hyperparameters were tuned for the **Chemprop** model.

| | CheMeleon | | | Chemprop | | |
|---|---|---|---|---|---|---|
| Task | MAE | RMSE | $R^2$ | MAE | RMSE | $R^2$ |
| UV/Vis | 16.86 | 31.36 | 0.911 | 16.9 | 31.1 | 0.913 |
| SAMPL6 | 0.27 | 0.29 | 0.812 | 0.34 | 0.46 | 0.525 |
| SAMPL7 | 0.43 | 0.62 | 0.127 | 0.33 | 0.49 | 0.449 |
| SAMPL9 | 0.81 | 1.01 | 0.799 | 0.92 | 1.10 | 0.758 |
| PCQM4Mv2 | 0.10 | 0.15 | 0.982 | 0.10 | 0.15 | 0.983 |

**Table B2**: Comparison of model performance on the test set, trained with **CheMeleon** and **Chemprop**, for different targets of QM9. The top group of tasks was trained together in a single multitask model. The bottom results are for single-task models.

| Model | Target | CheMeleon | | Chemprop | |
|---|---|---|---|---|---|
| | | MAE | RMSE | MAE | RMSE |
| multitask | mu | 0.347 | 0.607 | 0.339 | 0.595 |
| | alpha | 0.294 | 0.824 | 0.227 | 0.573 |
| | HOMO | 0.00249 | 0.00426 | 0.00245 | 0.00422 |
| | LUMO | 0.00247 | 0.00417 | 0.00238 | 0.00411 |
| | gap | 0.00338 | 0.00600 | 0.00333 | 0.00589 |
| | r2 | 18.1 | 36.0 | 17.2 | 33.4 |
| | ZPVE | 0.000463 | 0.000153 | 0.000235 | 0.000338 |
| | Cv | 0.134 | 0.331 | 0.108 | 0.223 |
| | U0 | 3.73 | 15.64 | 1.89 | 3.25 |
| | U298 | 3.75 | 15.75 | 1.90 | 3.26 |
| | H298 | 3.77 | 15.84 | 1.90 | 3.26 |
| | G298 | 3.53 | 14.51 | 1.84 | 3.21 |
| individual | gap | 0.00322 | 0.00588 | 0.00312 | 0.00591 |
| | U0 | 2.17 | 14.13 | 1.02 | 2.39 |

# Appendix C   Comparison Models

Unfortunately, many of the best models from the literature, such as those referenced in the 1, are not available for us to run and facilitate direct, fair, statistically rigorous comparisons between models. Grover [24], one of the first foundation models in this space, no longer provides model weights and has not been maintained since publication.

**Table B3**: Comparison of model performance on the test set, trained with **CheMeleon** and **Chemprop**, for HIV and PCBA classification tasks.

| | Split Type | CheMeleon | | | Chemprop | | |
|---|---|---|---|---|---|---|---|
| | | ROC-AUC | PRC-AUC | AP | ROC-AUC | PRC-AUC | AP |
| HIV | Random | 0.7681 | 0.3318 | 0.3369 | 0.7713 | 0.3048 | 0.3067 |
| PCBA | Random | 0.9098 | 0.2104 | 0.2151 | 0.9085 | 0.2146 | 0.2200 |
| PCBA (None) | Random | 0.9055 | 0.3693 | 0.3755 | 0.9075 | 0.3811 | 0.3852 |
| PCBA | Scaffold | 0.8895 | 0.2875 | 0.2913 | 0.8826 | 0.2886 | 0.2926 |

**MolE** [25] provides weights for a toy model used to demonstrate reproducibility, but not the full model. MolGPS [28] and Beaini et al. pre-trained foundation models following the concatenated datasets approach discussed previously, and released only the training datasets. GraphQPT [29] leveraged QM descriptors as targets during pre-training but have not yet made model checkpoints available.

Among the available models, we ran the following baseline, foundation, and classical models:

- `scikit-learn`'s Random Forest [41] with Morgan Fingerprint [7] or Mordred molecular descriptors, accessed via scikit-mol [57]: despite advances in deep methods, the work of Xia et al. has demonstrated that baseline methods can and do outperform deep methods.
- **fastprop** [5]: uses the Mordred descriptor set with a simple FNN but without any dimensionality reduction.
- **Chemprop**-Mordred: Combines the GNN-based molecular representation from **Chemprop** with Mordred descriptors as additional input features, following a similar strategy to the Chemprop-RDKit model used by Swanson et al..
- **minimol** pre-trained foundation model [22]: Graph Neural Network using an architecture similar to **Chemprop** pre-trained on a sparse labeled dataset from combined QM and biological assays, similar number of parameters to the present models allows for easy fine-tuning on commercial hardware.
- **MoLFormer** pre-trained foundation model [20]: SMILES-based transformer pre-trained using a massive corpus of unlabeled molecular structures.
- **MolCLR** pre-trained foundation model [23]: GNN-based model trained using contrastive learning on augmented 2D molecular graphs to learn general-purpose representations from unlabeled data.

# References

[1] Swanson, K., Walther, P., Leitz, J., Mukherjee, S., Wu, J.C., Shivnaraine, R.V., Zou, J.: Admet-ai: a machine learning admet platform for evaluation of large-scale chemical libraries. Bioinformatics **40**(7), 416 (2024)

[2] Yang, K., Swanson, K., Jin, W., Coley, C., Eiden, P., Gao, H., Guzman-Perez, A., Hopper, T., Kelley, B., Mathea, M., Palmer, A., Settels, V., Jaakkola, T., Jensen, K., Barzilay, R.: Analyzing learned molecular representations for property prediction. Journal of Chemical Information and Modeling **59**(8), 3370–3388 (2019) https://doi.org/10.1021/acs.jcim.9b00237 https://doi.org/10.1021/acs.jcim.9b00237. PMID: 31361484

[3] Attia, L., Burns, J.W., Doyle, P.S., Green, W.H.: Organic solubility prediction at the limit of aleatoric uncertainty (2024) https://doi.org/10.26434/chemrxiv-2024-93qp3

[4] Svetnik, V., Liaw, A., Tong, C., Culberson, J.C., Sheridan, R.P., Feuston, B.P.: Random forest: a classification and regression tool for compound classification and qsar modeling. Journal of chemical information and computer sciences **43**(6), 1947–1958 (2003)

[5] Burns, J.W., Green, W.H.: Generalizable, fast, and accurate deepqspr with fastprop. Journal of Cheminformatics **17**(1), 73 (2025) https://doi.org/10.1186/s13321-025-01013-4

[6] Heid, E., Greenman, K.P., Chung, Y., Li, S.-C., Graff, D.E., Vermeire, F.H., Wu, H., Green, W.H., McGill, C.J.: Chemprop: A machine learning package for chemical property prediction. Journal of Chemical Information and Modeling **64**(1), 9–17 (2024) https://doi.org/10.1021/acs.jcim.3c01250 https://doi.org/10.1021/acs.jcim.3c01250. PMID: 38147829

[7] Morgan, H.L.: The generation of a unique machine description for chemical structures-a technique developed at chemical abstracts service. Journal of Chemical Documentation **5**(2), 107–113 (1965) https://doi.org/10.1021/c160017a018

[8] Chung, Y., Vermeire, F.H., Wu, H., Walker, P.J., Abraham, M.H., Green, W.H.: Group contribution and machine learning approaches to predict abraham solute parameters, solvation free energy, and solvation enthalpy. Journal of Chemical Information and Modeling **62**(3), 433–446 (2022)

[9] Zalte, A.S., Pang, H.-W., Doner, A.C., Green, W.H.: RIGR: Resonance invariant graph representation for molecular property prediction (2025) https://doi.org/10.26434/chemrxiv-2025-qgfxp

[10] McGill, C., Forsuelo, M., Guan, Y., Green, W.H.: Predicting infrared spectra with message passing neural networks. Journal of Chemical Information and Modeling

**61**(6), 2594–2609 (2021) https://doi.org/10.1021/acs.jcim.1c00055

[11] Magar, R., Wang, Y., Lorsung, C., Liang, C., Ramasubramanian, H., Li, P., Farimani, A.B.: Auglichem: data augmentation library of chemical structures for machine learning. Machine Learning: Science and Technology **3**(4), 045015 (2022)

[12] Aldeghi, M., Coley, C.W.: A graph representation of molecular ensembles for polymer property prediction. Chemical Science **13**(35), 10486–10498 (2022)

[13] Xia, J., Zhang, L., Zhu, X., Li, S.Z.: Why Deep Models Often cannot Beat Non-deep Counterparts on Molecular Property Prediction? (2023). https://arxiv.org/abs/2306.17702

[14] Li, S.-C., Wu, H., Menon, A., Spiekermann, K.A., Li, Y.-P., Green, W.H.: When do quantum mechanical descriptors help graph neural networks to predict chemical properties? Journal of the American Chemical Society **146**(33), 23103–23120 (2024)

[15] Choi, J., Nam, G., Choi, J., Jung, Y.: A perspective on foundation models in chemistry. JACS Au (2025)

[16] Weininger, D.: Smiles, a chemical language and information system. 1. introduction to methodology and encoding rules. Journal of Chemical Information and Computer Sciences **28**(1), 31–36 (1988) https://doi.org/10.1021/ci00057a005

[17] Weininger, D., Weininger, A., Weininger, J.L.: Smiles. 2. algorithm for generation of unique smiles notation. Journal of Chemical Information and Computer Sciences **29**(2), 97–101 (1989) https://doi.org/10.1021/ci00062a008

[18] Krenn, M., Häse, F., Nigam, A., Friederich, P., Aspuru-Guzik, A.: Self-referencing embedded strings (selfies): A 100% robust molecular string representation. Machine Learning: Science and Technology **1**(4), 045024 (2020)

[19] Chithrananda, S., Grand, G., Ramsundar, B.: ChemBERTa: Large-Scale Self-Supervised Pretraining for Molecular Property Prediction. arXiv (2020). https://doi.org/10.48550/ARXIV.2010.09885 . https://arxiv.org/abs/2010.09885

[20] Ross, J., Belgodere, B., Chenthamarakshan, V., Padhi, I., Mroueh, Y., Das, P.: Large-scale chemical language representations capture molecular structure and properties. Nature Machine Intelligence **4**(12), 1256–1264 (2022) https://doi.org/10.1038/s42256-022-00580-7

[21] Galkin, M., Yuan, X., Mostafa, H., Tang, J., Zhu, Z.: Towards foundation models for knowledge graph reasoning. arXiv preprint arXiv:2310.04562 (2023)

[22] Kläser, K., Banaszewski, B., Maddrell-Mander, S., McLean, C., Müller, L., Parviz, A., Huang, S., Fitzgibbon, A.: `MiniMol`: A Parameter-Efficient Foundation Model

for Molecular Learning (2024). https://arxiv.org/abs/2404.14986

[23] Wang, Y., Wang, J., Cao, Z., Barati Farimani, A.: Molecular contrastive learning of representations via graph neural networks. Nature Machine Intelligence **4**(3), 279–287 (2022) https://doi.org/10.1038/s42256-022-00447-x

[24] Rong, Y., Bian, Y., Xu, T., Xie, W., Wei, Y., Huang, W., Huang, J.: Self-supervised graph transformer on large-scale molecular data. In: Proceedings of the 34th International Conference on Neural Information Processing Systems. NIPS '20. Curran Associates Inc., Red Hook, NY, USA (2020)

[25] Méndez-Lucio, O., Nicolaou, C.A., Earnshaw, B.: Mole: a foundation model for molecular graphs using disentangled attention. Nature Communications **15**(1) (2024) https://doi.org/10.1038/s41467-024-53751-y

[26] Beaini, D., Huang, S., Cunha, J.A., Li, Z., Moisescu-Pareja, G., Dymov, O., Maddrell-Mander, S., McLean, C., Wenkel, F., Müller, L., Mohamud, J.H., Parviz, A., Craig, M., Koziarski, M., Lu, J., Zhu, Z., Gabellini, C., Klaser, K., Dean, J., Wognum, C., Sypetkowski, M., Rabusseau, G., Rabbany, R., Tang, J., Morris, C., Ravanelli, M., Wolf, G., Tossou, P., Mary, H., Bois, T., Fitzgibbon, A.W., Banaszewski, B., Martin, C., Masters, D.: Towards foundational models for molecular learning on large-scale multi-task datasets. In: The Twelfth International Conference on Learning Representations (2024). https://openreview.net/forum?id=Zc2aIcucwc

[27] Landrum, G.A., Riniker, S.: Combining ic50 or ki values from different sources is a source of significant noise. Journal of Chemical Information and Modeling **64**(5), 1560–1567 (2024) https://doi.org/10.1021/acs.jcim.4c00049

[28] Sypetkowski, M., Wenkel, F., Poursafaei, F., Dickson, N., Suri, K., Fradkin, P., Beaini, D.: On the scalability of GNNs for molecular graphs. In: The Thirty-eighth Annual Conference on Neural Information Processing Systems (2024). https://openreview.net/forum?id=klqhrq7fvB

[29] Fallani, A., Nugmanov, R., Arjona-Medina, J., Wegner, J.K., Tkatchenko, A., Chernichenko, K.: Pretraining graph transformers with atom-in-a-molecule quantum properties for improved admet modeling. Journal of Cheminformatics **17**(1) (2025) https://doi.org/10.1186/s13321-025-00970-0

[30] Kim, S., Chen, J., Cheng, T., Gindulyte, A., He, J., He, S., Li, Q., Shoemaker, B.A., Thiessen, P.A., Yu, B., Zaslavsky, L., Zhang, J., Bolton, E.E.: Pubchem 2025 update. Nucleic Acids Research **53**(D1), 1516–1525 (2024) https://doi.org/10.1093/nar/gkae1059

[31] Miles, A.: zarr-developers/zarr-python: v3.0.7. Zenodo (2025). https://doi.org/10.5281/ZENODO.3773449 . https://zenodo.org/doi/10.5281/zenodo.3773449

[32] Falcon, W., The PyTorch Lightning team: PyTorch Lightning. https://doi.org/10.5281/zenodo.3828935 . https://github.com/Lightning-AI/lightning

[33] Tilborg, D., Alenicheva, A., Grisoni, F.: Exposing the limitations of molecular machine learning with activity cliffs. Journal of Chemical Information and Modeling **62**(23), 5938–5951 (2022) https://doi.org/10.1021/acs.jcim.2c01073 https://doi.org/10.1021/acs.jcim.2c01073. PMID: 36456532

[34] O'Boyle, N.M., Sayle, R.A.: Comparing structural fingerprints using a literature-based similarity benchmark. Journal of Cheminformatics **8**(1) (2016) https://doi.org/10.1186/s13321-016-0148-0

[35] Carhart, R.E., Smith, D.H., Venkataraghavan, R.: Atom pairs as molecular features in structure-activity studies: definition and applications. Journal of Chemical Information and Computer Sciences **25**(2), 64–73 (1985) https://doi.org/10.1021/ci00046a002

[36] Nilakantan, R., Bauman, N., Dixon, J.S., Venkataraghavan, R.: Topological torsion: a new molecular descriptor for sar applications. comparison with other descriptors. Journal of Chemical Information and Computer Sciences **27**(2), 82–85 (1987) https://doi.org/10.1021/ci00054a008

[37] Landrum, G.: rdkit. Zenodo (2025). https://doi.org/10.5281/ZENODO.591637 . https://zenodo.org/doi/10.5281/zenodo.591637

[38] Gedeck, P., Rohde, B., Bartels, C.: Qsar - how good is it in practice? comparison of descriptor sets on an unbiased cross section of corporate data sets. Journal of Chemical Information and Modeling **46**(5), 1924–1936 (2006) https://doi.org/10.1021/ci050413p

[39] Durant, J.L., Leland, B.A., Henry, D.R., Nourse, J.G.: Reoptimization of mdl keys for use in drug discovery. Journal of Chemical Information and Computer Sciences **42**(6), 1273–1280 (2002) https://doi.org/10.1021/ci010132r

[40] Orlov, A.A., Akhmetshin, T.N., Horvath, D., Marcou, G., Varnek, A.: From high dimensions to human insight: Exploring dimensionality reduction for chemical space visualization. Molecular Informatics **44**(1) (2024) https://doi.org/10.1002/minf.202400265

[41] Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., Duchesnay, E.: Scikit-learn: Machine learning in Python. Journal of Machine Learning Research **12**, 2825–2830 (2011)

[42] Tetko, I.V., Gasteiger, J., Todeschini, R., Mauri, A., Livingstone, D., Ertl, P., Palyulin, V.A., Radchenko, E.V., Zefirov, N.S., Makarenko, A.S., Tanchuk, V.Y., Prokopenko, V.V.: Virtual computational chemistry laboratory – design and

description. Journal of Computer-Aided Molecular Design **19**(6), 453–463 (2005) https://doi.org/10.1007/s10822-005-8694-y

[43] Yap, C.W.: Padel-descriptor: An open source software to calculate molecular descriptors and fingerprints. Journal of Computational Chemistry **32**(7), 1466–1474 (2010) https://doi.org/10.1002/jcc.21707

[44] Gaulton, A., Bellis, L.J., Bento, A.P., Chambers, J., Davies, M., Hersey, A., Light, Y., McGlinchey, S., Michalovich, D., Al-Lazikani, B., Overington, J.P.: Chembl: a large-scale bioactivity database for drug discovery. Nucleic Acids Research **40**(D1), 1100–1107 (2011) https://doi.org/10.1093/nar/gkr777

[45] Sorokina, M., Merseburger, P., Rajan, K., Yirik, M.A., Steinbeck, C.: Coconut online: Collection of open natural products database. Journal of Cheminformatics **13**(1) (2021) https://doi.org/10.1186/s13321-020-00478-9

[46] Wiener, H.: Structural determination of paraffin boiling points. Journal of the American Chemical Society **69**(1), 17–20 (1947) https://doi.org/10.1021/ja01193a005 https://doi.org/10.1021/ja01193a005. PMID: 20291038

[47] Abraham, M.H., McGowan, J.C.: The use of characteristic volumes to measure cavity terms in reversed phase liquid chromatography. Chromatographia **23**(4), 243–246 (1987) https://doi.org/10.1007/bf02311772

[48] Moriwaki, H., Tian, Y.-S., Kawashita, N., Takagi, T.: Mordred: a molecular descriptor calculator. Journal of Cheminformatics **10**(1) (2018) https://doi.org/10.1186/s13321-018-0258-y

[49] Gorishniy, Y., Rubachev, I., Babenko, A.: On Embeddings for Numerical Features in Tabular Deep Learning. arXiv (2022). https://doi.org/10.48550/ARXIV.2203.05556 . https://arxiv.org/abs/2203.05556

[50] Capela, F., Nouchi, V., Deursen, R.V., Tetko, I.V., Godin, G.: Multitask Learning On Graph Neural Networks Applied To Molecular Property Predictions (2019). https://arxiv.org/abs/1910.13124

[51] Tetko, I.V., Deursen, R., Godin, G.: Be aware of overfitting by hyperparameter optimization! Journal of Cheminformatics **16**(1) (2024) https://doi.org/10.1186/s13321-024-00934-w

[52] Ash, J.R., Wognum, C., Rodríguez-Pérez, R., Aldeghi, M., Cheng, A.C., Clevert, D.-A., Engkvist, O., Fang, C., Price, D.J., Hughes-Oliver, J.M., al.: Practically significant method comparison protocols for machine learning in small molecule drug discovery. ChemRxiv (2024) https://doi.org/10.26434/chemrxiv-2024-6dbwv-v2

[53] Huang, K., Fu, T., Gao, W., Zhao, Y., Roohani, Y., Leskovec, J., Coley, C.W.,

Xiao, C., Sun, J., Zitnik, M.: Therapeutics data commons: Machine learning datasets and tasks for drug discovery and development. Proceedings of Neural Information Processing Systems, NeurIPS Datasets and Benchmarks (2021)

[54] Wu, Z., Ramsundar, B., Feinberg, E.N., Gomes, J., Geniesse, C., Pappu, A.S., Leswing, K., Pande, V.: Moleculenet: a benchmark for molecular machine learning. Chemical Science **9**(2), 513–530 (2018) https://doi.org/10.1039/c7sc02664a

[55] Walters, P.: We Need Better Benchmarks for Machine Learning in Drug Discovery — practicalcheminformatics.blogspot.com. https://practicalcheminformatics. blogspot.com/2023/08/we-need-better-benchmarks-for-machine.html. [Accessed 13-05-2025]

[56] Virtanen, P., Gommers, R., Oliphant, T.E., Haberland, M., Reddy, T., Cournapeau, D., Burovski, E., Peterson, P., Weckesser, W., Bright, J., van der Walt, S.J., Brett, M., Wilson, J., Millman, K.J., Mayorov, N., Nelson, A.R.J., Jones, E., Kern, R., Larson, E., Carey, C.J., Polat, İ., Feng, Y., Moore, E.W., VanderPlas, J., Laxalde, D., Perktold, J., Cimrman, R., Henriksen, I., Quintero, E.A., Harris, C.R., Archibald, A.M., Ribeiro, A.H., Pedregosa, F., van Mulbregt, P., SciPy 1.0 Contributors: SciPy 1.0: Fundamental Algorithms for Scientific Computing in Python. Nature Methods **17**, 261–272 (2020) https://doi.org/10.1038/s41592-019-0686-2

[57] Bjerrum, E.J., Bachorz, R.A., Bitton, A., Choung, O.-h., Chen, Y., Esposito, C., Ha, S.V., Poehlmann, A.: Scikit-Mol brings cheminformatics to Scikit-Learn. ChemRxiv (2023) https://doi.org/10.26434/chemrxiv-2023-fzqwd . preprint. Accessed 2023-12-06
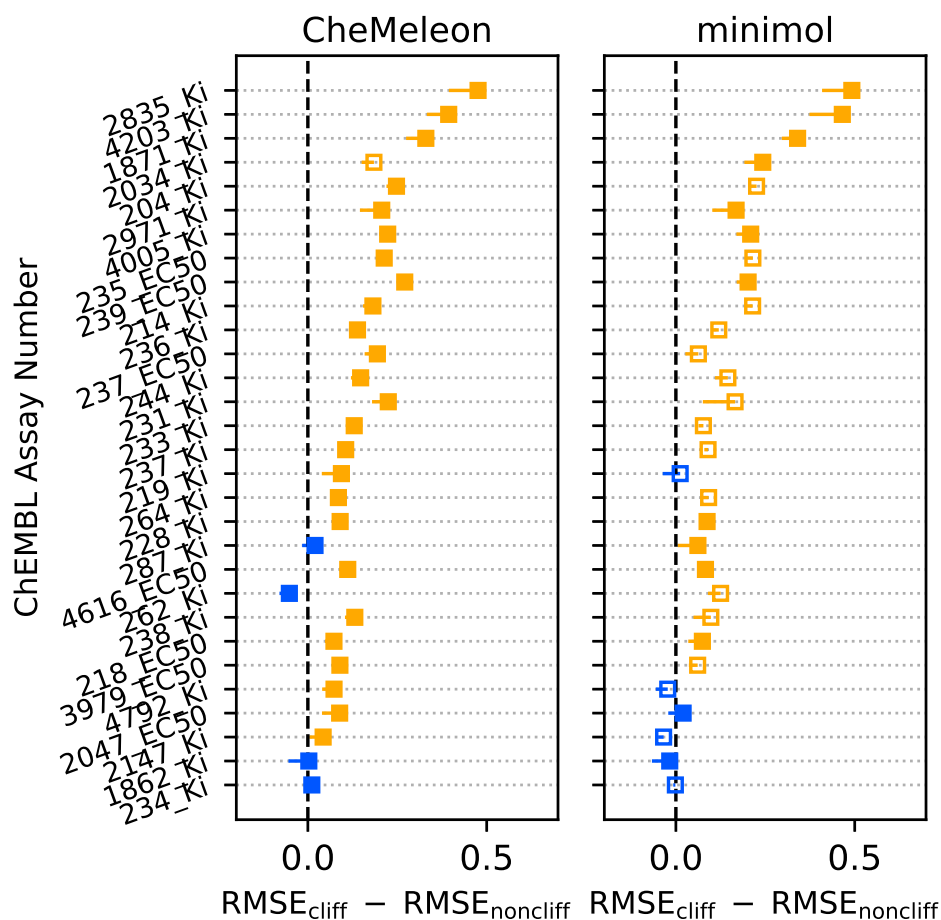
**Fig. 2**: Performance of models across the ChEMBL assays curated as part of the MoleculeACE study [33]. Each assay is one row along the horizontal axis and the plotted value on the x-axis is the difference in root mean squared error for predictions of molecules in the cliff set ("cliff") and those not in the cliff set ("noncliff"). Dots shown in blue are not practically different from zero, a positive result indicating that the model performance on the two sets is indistinguishable. The filling of the dots indicates the RMSE of the model across the entire test set relative to the other models, with filled dots indicating that the given model was statistically the best or indistinguishable from the best performer and hollow dots indicating that it was practically worse.
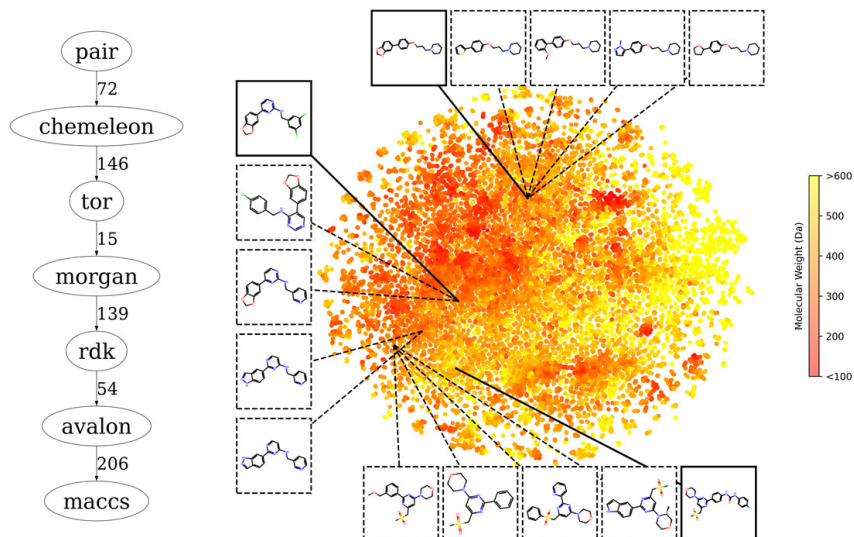
**Fig. 3**: (left) Hasse diagram showing the difference in performance on the O'Boyle and Sayle Single Assay benchmark when using **CheMeleon** as a fingerprint in comparison to the traditional molecular fingerprints Atom Pair Fingerprints [35] ("pair"), Topological-torsion fingerprints [36] ("tor"), Morgan [7] fingerprint ("morgan"), RDKit [37] fingerprint ("rdk"), Avalaon [38] fingerprint ("avalon"), and MACCS Keys [39] ("maccs"). Directed edge indicates that a practically significant difference in correctness of sorting order exists between the two nodes, with edge weight indicating the number of series in the entire benchmark set for which a practical difference was observed. See the original study for a more detailed discussion on the formulation and interpretation of this Hasse diagram [34]. (right) t-SNE project computed using `scikit-learn` [41] with perplexity hyperparameter set such that highly similar neighborhoods in the original feature space are preserved in the projected space following the procedure laid out by Orlov et al. Three of the series present in the benchmark data have been highlighted with the "lead" molecule shown in bold and the consecutively more dissimilar molecules shown after it.
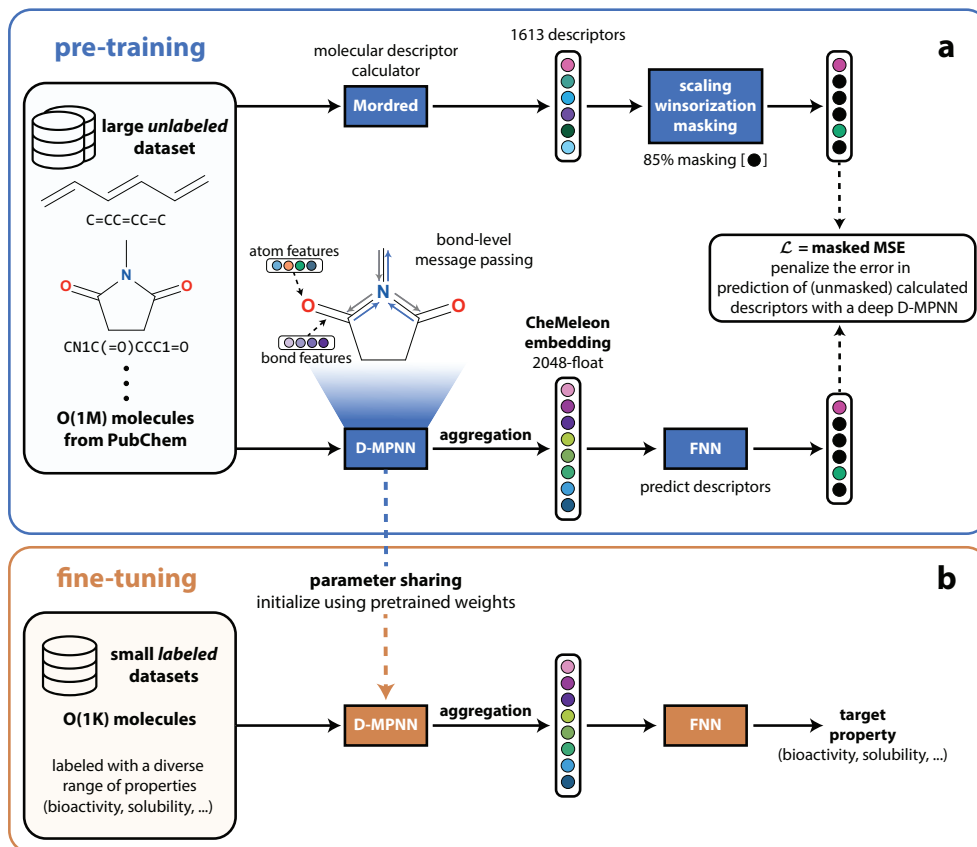
**Fig. 4**: Workflow for the present study. (a) A large corpus of unlabeled SMILES strings is randomly selected from PubChem [30] and featurized into a vector of molecular descriptors using Mordred [48]. **Chemprop** is used to train a Directed-Message Passing Neural Network [2, 6] (D-MPNN) to predict these descriptors using masked loss analogous to ChemBERTa as a form of regularization [19]. (b) The resulting D-MPNN is then reused for subsequent fine-tuning on smaller downstream datasets labeled with quantities of interest, such as bioactivity.
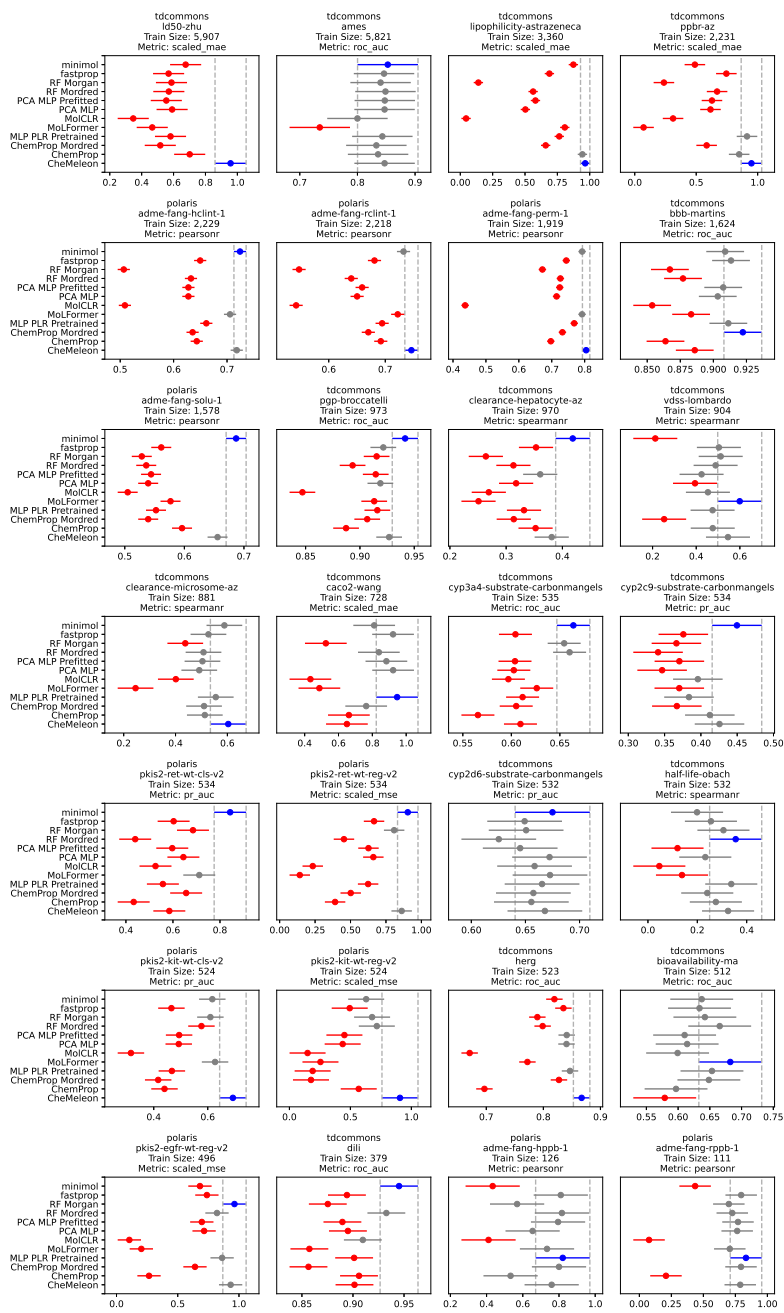
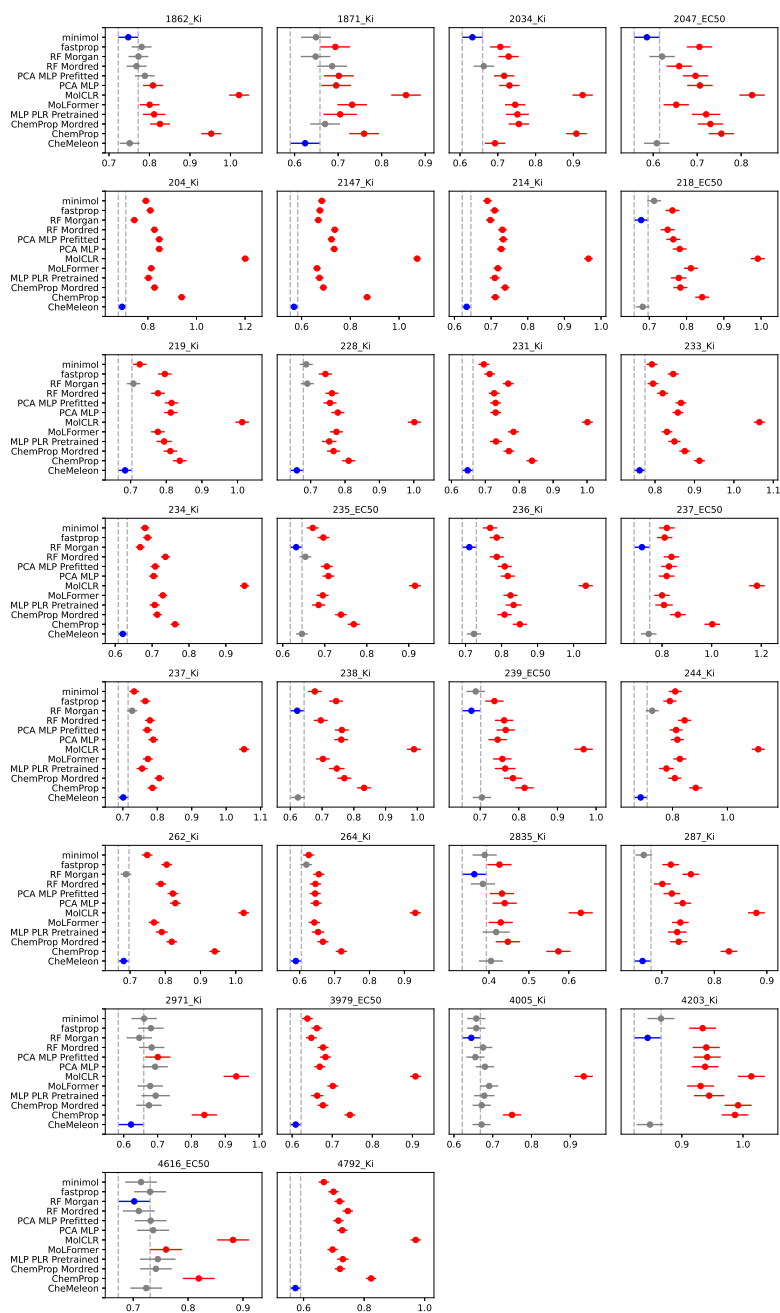**Fig. B1**: All results for Polaris benchmarks following the same conventions as main text figure 1.

23

**Fig. B2**: All results for MoleculeACE benchmarks following the same conventions as the Polaris figures.
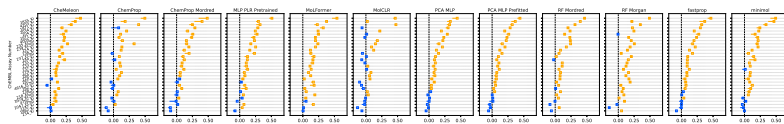
**Fig. B3**: All MoleculeACE consistency results following the same conventions as the main text Figure 2.