

1 Supporting Information: Deep Learning  
2 Foundation Models from Classical Molecular  
3 Descriptors

4 Jackson W. Burns<sup>1†</sup>, Akshat Shirish Zalte<sup>1†</sup>,  
5 Charles R. A. Abreu<sup>1</sup>, Jochen Sieg<sup>2</sup>, Christian Feldmann<sup>2</sup>,  
6 Miriam Mathea<sup>2</sup>, William H. Green<sup>1\*</sup>

7 <sup>1</sup>Department of Chemical Engineering, MIT, Cambridge, Massachusetts,  
8 USA.

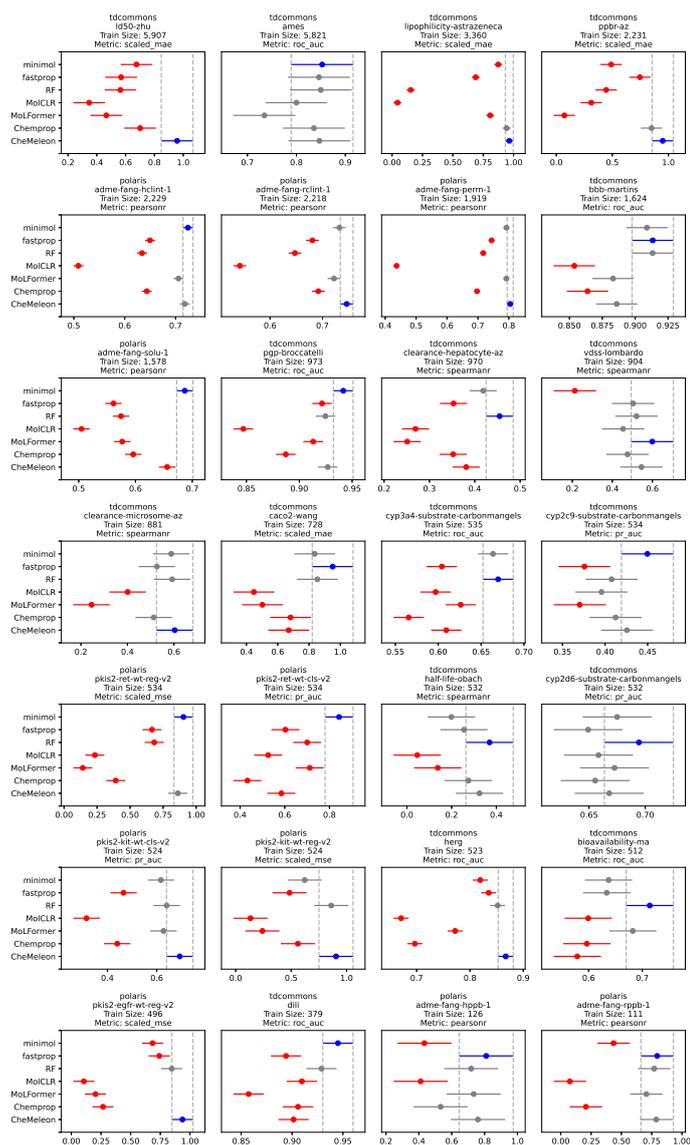
9 <sup>2</sup>BASF SE, Ludwigshafen, Germany.

10 \*Corresponding author(s). E-mail(s): [whgreen@mit.edu](mailto:whgreen@mit.edu);

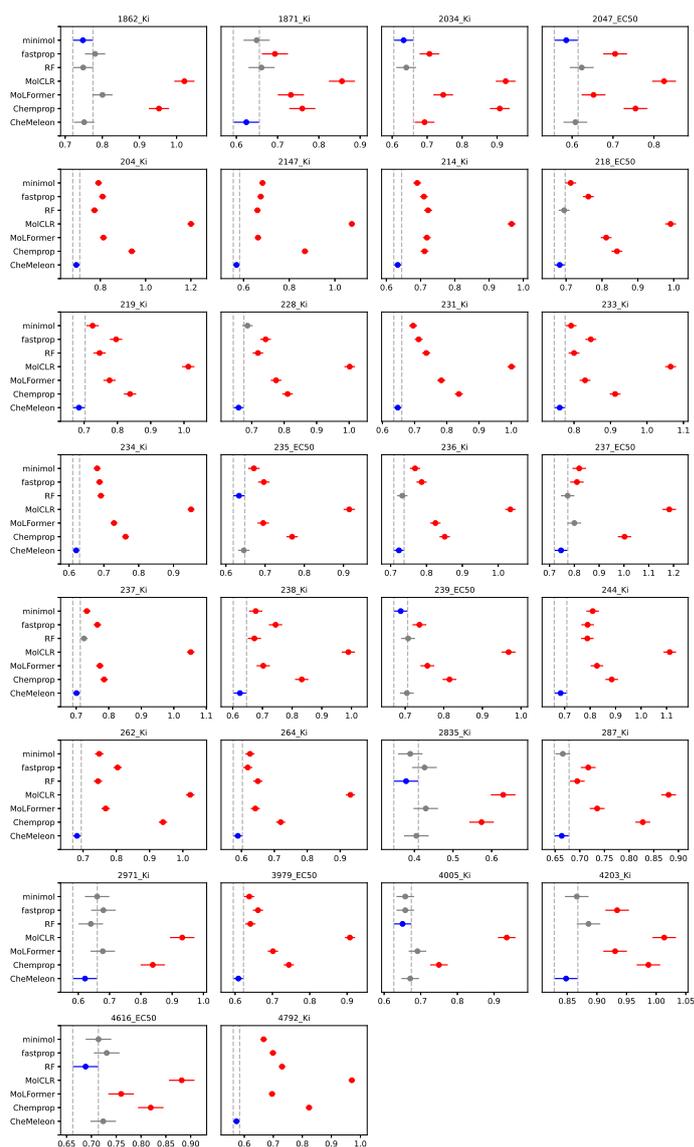
11 <sup>†</sup>These authors contributed equally to this work.

12 **S1 Detailed Experimental Results**

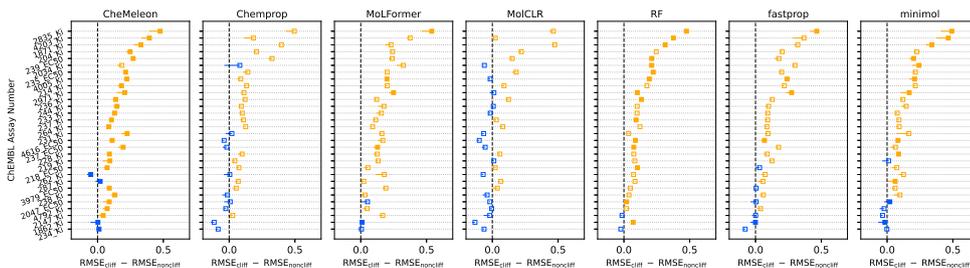
13 The complete set of MoleculeACE-style diagrams and Honestly Significant Difference  
14 diagrams following the conventions described in the main text (Figures 2 and 3, respec-  
15 tively) are provided in Figures S1, S2, and S3. A high-quality version of these images  
16 is also available in the GitHub repository described in Section 3.5.



**Supplementary Figure S1:** Performance of all of the tested models across all of the tested Polaris benchmark 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 by training size in decreasing order. 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 3.2. Models shown in red *are* practically worse performers and are considered to have “lost” on the indicated benchmark.



**Supplementary Figure S2:** Performance of all of the tested models on the MoleculeACE study in terms of Root Mean Square Error (RMSE). The ChEMBL assay and target corresponding to each benchmark is shown as the title of each subplot. 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 3.2. Models shown in red *are* practically worse performers and are considered to have “lost” on the indicated benchmark.



**Supplementary Figure S3:** Performance of all tested models across the ChEMBL assays [1] curated as part of the MoleculeACE study [2]. 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. Filled dots indicating that the given model was statistically the best or indistinguishable from the best performer and hollow dots indicating that it was significantly worse than the best model for that dataset.

### 17 S1.1 MoleculeNet Benchmarks

18 This section extends our evaluation of **CheMeleon** to include all molecular property  
 19 prediction benchmarks originally used to assess **Chemprop** [3], in particular  
 20 the MoleculeNet benchmarks [4]. These were excluded from the main text because  
 21 more relevant benchmarks have since been devised [5]. We used the same datasets  
 22 and data splits as in the original **Chemprop** paper. These benchmarks include a mix  
 23 of regression tasks (UV/Vis, SAMPL, QM9, PCQM4Mv2) and classification tasks  
 24 (HIV, PCBA). No hyperparameter tuning was performed for either **CheMeleon** or  
 25 **Chemprop** models.

**Table S1:** 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 the PCQM4Mv2 dataset.

Task	CheMeleon			Chemprop		
	MAE	RMSE	R <sup>2</sup>	MAE	RMSE	R <sup>2</sup>
UV/Vis	16.86	31.36	0.911	20.00	34.2	0.894
SAMPL6	0.28	0.31	0.784	0.30	0.39	0.662
SAMPL7	0.48	0.70	-0.102	0.36	0.55	0.301
SAMPL9	0.85	1.04	0.786	0.94	1.09	0.763
PCQM4Mv2	0.09	0.15	0.982	0.10	0.16	0.982

**Table S2:** Comparison of model performance on the test set, trained with **CheMeleon** and **Chemprop**, for different targets of QM9. The tasks in the top group were trained together in a single multitask model. The bottom results are for single-task models.

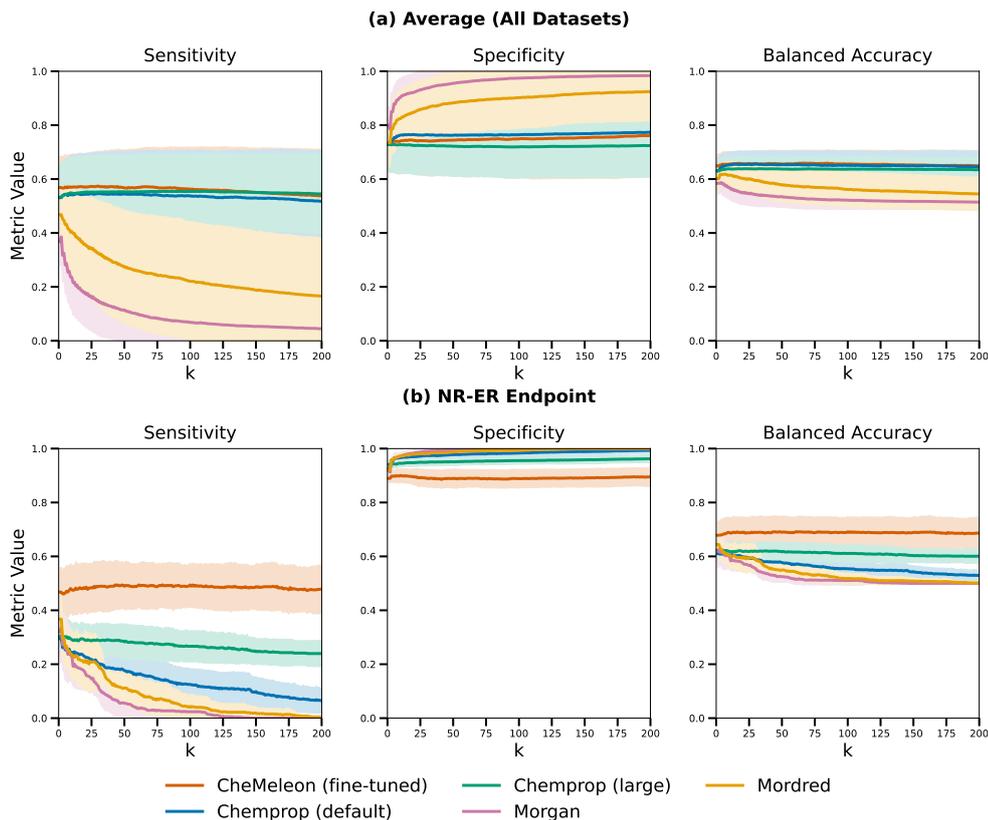
Model	Target	CheMeleon		Chemprop	
		MAE	RMSE	MAE	RMSE
multitask	mu	0.347	0.607	0.408	0.655
	alpha	0.294	0.824	0.371	0.680
	HOMO	0.00249	0.00426	0.00322	0.00499
	LUMO	0.00247	0.00417	0.00350	0.00519
	gap	0.00338	0.00600	0.00454	0.00709
	r2	18.1	36.0	24.6	40.3
	ZPVE	0.000463	0.000153	0.000416	0.000551
	Cv	0.134	0.331	0.181	0.288
	U0	3.73	15.64	3.55	5.37
	U298	3.75	15.75	3.54	5.38
	H298	3.77	15.84	3.54	5.38
	G298	3.53	14.51	3.50	5.36
individual	gap	0.00322	0.00588	0.00390	0.00653
	U0	2.17	14.13	1.45	2.90

**Table S3:** 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	0.7990	0.3080	0.3057
PCBA	Random	0.9098	0.2104	0.2151	0.9045	0.1971	0.2023
PCBA (None)	Random	0.9055	0.3693	0.3755	0.9004	0.3614	0.3565
PCBA	Scaffold	0.8895	0.2875	0.2913	0.8812	0.265	0.2705

## 26 S2 k-Nearest Neighbors Representation Probing

27 Following the conventions of main text Section 1.4, Figure S4 presents results for the  
 28 kNN probing study of the **CheMeleon** learned representation using an agglomerative  
 29 clustering-based split.



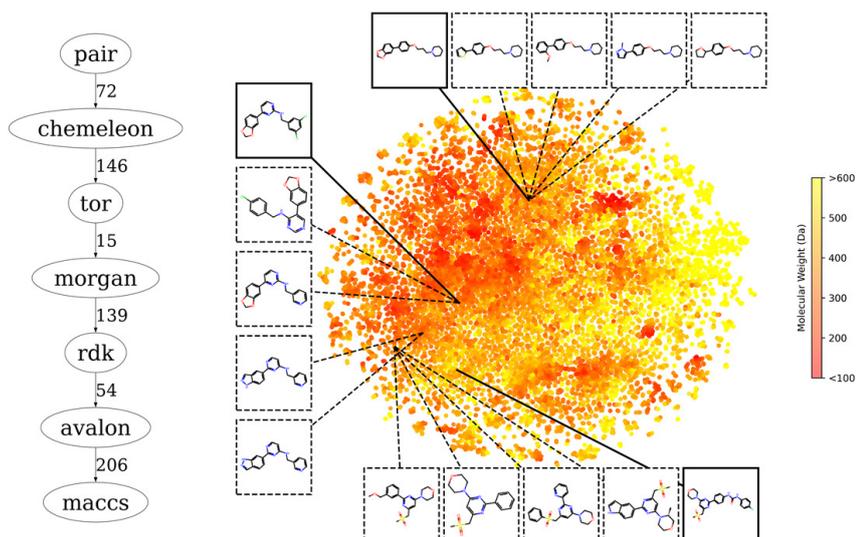
**Supplementary Figure S4:** kNN probing of fixed and learned molecule representations for the agglomerative clustering split. Upper row shows average results over 20 ToxCast endpoints. Lower row shows the results on the NR-ER endpoint, which was also used in the comparison by Ball et al. [6]. The results were obtained with a 5-fold cross validation using an agglomerative clustering split.

30 These results reinforce the conclusions of the main text, again showing that  
 31 **CheMeleon** offers favorable performance improvements and trade-offs relative to  
 32 standard models. The ablation model further reiterates that performance improve-  
 33 ments are attributable to descriptor-based pre-training rather than an architectural  
 34 bias.

### 35 S3 Foundation Fingerprint

36 Adapting the provided source code from the O’Boyle and Sayle [7] study, we gener-  
 37 ated the learned representation (LR) for each of the molecules provided in the Single  
 38 Assay benchmark. We calculated the cosine distance between the embeddings within  
 39 each series to arrive at a **CheMeleon**-based sorting order. This was compared with

40 the sort order derived from Atom Pair fingerprints [8], Topological-torsion fingerprints [9], Morgan fingerprints [10], RDKit fingerprints [11], Avalon fingerprints [12], and  
 41 MACCS keys [13] as they collectively constitute the standard set of molecular fingerprints used for molecular representation, such as for feeding the Random Forest  
 42 model used here. This provides insights into how the model has learned to separate common chemical moieties. The same **CheMeleon** fingerprints are then subjected to  
 43 the methods of Orlov et al. [14] to generate a two-dimensional projection, and a few  
 44 of the series present in the benchmark are highlighted. An overview of the complete  
 45 fingerprint analysis is presented in Figure S5.



**Supplementary Figure S5:** (left) Hasse diagram showing the difference in performance on the O’Boyle and Sayle [7] Single Assay benchmark when using **CheMeleon** as a fingerprint in comparison to the traditional Atom Pair fingerprints [8] (“pair”), Topological-torsion fingerprints [9] (“tor”), Morgan fingerprints [10] (“morgan”), RDKit fingerprints [11] (“rdk”), Avalon fingerprints [12] (“avalon”), and MACCS keys [13] (“maccs”). A 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 [7]. (right) t-SNE projection computed using `scikit-learn` [15] with perplexity hyperparameter tuned 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. [14]. 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.

49 The t-SNE projection shows that **CheMeleon** can separate large and small  
50 molecules in its feature space and that the three highlighted assays are far apart from  
51 each other. Each assay also provides further chemical insight. At the top, all five  
52 species are small modifications of the original structure at the same site and are thus  
53 all projected into similar locations in chemical space. The assay on the left shows that  
54 modifications removing the methylenedioxy group cause two subsets of the series to be  
55 projected into *different* regions of chemical space. The final assay (shown at the bot-  
56 tom) follows a similar trend. In this case, however, the lead compound is not projected  
57 into the same chemical space as the other four compounds in the series, reflecting the  
58 extent of the initial structural modification.

## 59 S4 Training Details

60 All models were trained using PyTorch Lightning [16] and associated open-source  
61 machine learning and cheminformatics packages. We used eight Nvidia 2080 Ti GPUs  
62 for foundation model training and an Nvidia Quadro RTX 4000 laptop GPU for fine-  
63 tuning. The architecture for **CheMeleon** was as described in the main text with  
64 hyperparameters chosen on the basis of heuristics.

65 To pre-train **CheMeleon**, a set of 1 million molecules was randomly selected from  
66 the PubChem database. They were further pre-processed to remove SMILES strings  
67 that were longer than 150 characters, had multiple fragments (such as salts), or could  
68 not be converted into a valid molecular graph using the RDKit [11]. From there,  
69 descriptors were calculated and saved to disk, as described in the main text. Note that  
70 we used [mordred-community](#), a community-maintained fork of the original Mordred  
71 descriptor calculator, which is no longer maintained.

72 Unfortunately, many of the reportedly best models from the literature, such as  
73 those referenced in the main text, are not available for us to run, preventing direct,  
74 fair, and statistically rigorous comparisons between models. **GROVER** [17], one of  
75 the first foundation models in this space, no longer provides model weights and has not  
76 been maintained since publication. **MolE** [18] provides weights for a toy model used  
77 to demonstrate reproducibility, but not the full model. **MolGPS** [19] and Beaini et al.  
78 [20] pre-trained foundation models following the concatenated datasets approach dis-  
79 cussed previously, and released only the training datasets. **GraphQPT** [21] leveraged  
80 QM descriptors as targets during pre-training but has not yet made model checkpoints  
81 available.

82 Among the models which were available, we selected a representative set of the  
83 best performers. Every reference model benchmarked as part of this study used the  
84 following architecture and training approach, which were decided in deference to the  
85 suggested configuration from each of the original studies where such an indication was  
86 given, and held consistent otherwise:

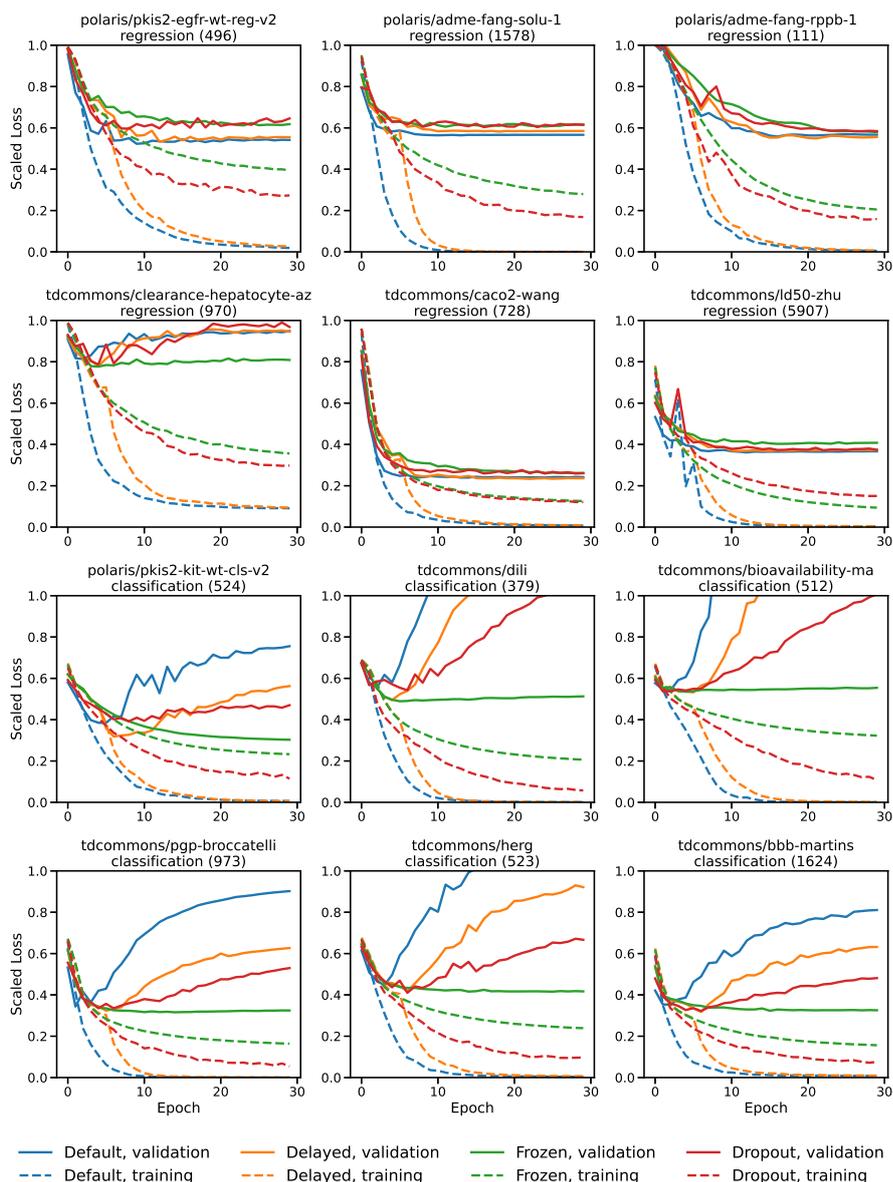
- 87 • **Random Forest**: increased number of estimators to 500 as suggested in Svetnik  
88 et al. [22] and set the maximum number of features to  $\log_2$  of their total number
- 89 • **MolFormer**: tuned the entire network with a single additional readout layer
- 90 • **fastprop**: default settings as set by Burns and Green [23]
- 91 • **Chemprop**: default settings as set by Heid et al. [3]

92 • **minimol**: learned representation held constant, tuned a multi-layer skip connection  
93 MLP with regularization and hidden size of 512 as described by Kläser et al. [24]

94 When using **CheMeleon** for a benchmark, an equivalent directed message-passing  
95 neural network (D-MPNN) is initialized with the pre-trained weights and biases ( $\approx 8.7$   
96 million parameters) and followed by a feedforward neural network (FNN) initialized  
97 from scratch (single hidden layer,  $\approx 525$  thousand parameters). The entire model is  
98 then fine-tuned end-to-end using stochastic gradient descent. Despite the distinct ini-  
99 tialization and disparate parameter counts, treating the two submodels equally during  
100 fine-tuning is typically the most effective approach, as we demonstrate in Section S5.

## 101 S5 Analysis of Fine-Tuning Strategies

102 The choice of fine-tuning strategy can significantly impact the performance of a pre-  
103 trained model on downstream tasks [25, 26]. We evaluated several fine-tuning strategies  
104 across the 28 Polaris benchmarks. These include the default end-to-end fine-tuning,  
105 delayed unfreezing of the D-MPNN, permanently freezing the D-MPNN, and applying  
106 dropout regularization. Averaged learning curves obtained by applying these strate-  
107 gies to 12 representative benchmarks from the Polaris benchmark set are shown in  
108 Figure S6.



**Supplementary Figure S6:** Learning curves for the fine-tuning of **CheMeleon** on selected Polaris benchmarks. Four different strategies were tested: 1. *Default*: fine-tuning the entire model using default **Chemprop** settings, 2. *Delayed*: freezing the D-MPNN weights for the first 5 epochs, 3. *Frozen*: freezing the D-MPNN weights for the entire process, and 4. *Dropout*: applying a 20% dropout rate to both the D-MPNN and FNN layers. The curves are averaged over 5 runs with different random splits (train:validation = 4:1). The subplot titles indicate the benchmark name, the task type, and the training set size.

109 The first two rows of Figure S6 correspond to regression tasks, while the third and  
110 fourth rows correspond to classification tasks. Under the default fine-tuning strategy,  
111 the validation loss decreases rapidly in all cases, typically reaching a minimum within  
112 fewer than 10 epochs. Overall, **ChemEleon**'s pretrained representations appear par-  
113 ticularly well-suited for regression tasks, with the default strategy yielding the best or  
114 near-best performance across all regression benchmarks.

115 Notably, for regression tasks we observe little to no degradation in validation loss  
116 after the minimum is reached, even as the training loss continues to decrease, often  
117 approaching zero. Although this behavior may appear counterintuitive, it has been  
118 reported previously in large, overparameterized neural networks trained with stochas-  
119 tic gradient descent, particularly when initialized from pretrained representations and  
120 applied to small datasets [27, 28]. One possible explanation is an effective separation  
121 between network components that capture underlying structure and those that pri-  
122 marily memorize label noise, such that memorization does not substantially impair  
123 generalization to unseen inputs.

124 In contrast, classification tasks exhibit more pronounced overfitting, making early  
125 stopping necessary to prevent degradation in validation performance. In these cases,  
126 freezing the D-MPNN for the first few epochs sometimes lowers the minimum valida-  
127 tion loss, although the effect is inconsistent across benchmarks. The transition from  
128 pre-training to fine-tuning likely induces a more severe change in the effective loss  
129 landscape for classification tasks than for regression tasks, which may reduce transfer  
130 efficiency. However, further investigation is required to test this hypothesis.

131 Permanently freezing the backbone consistently underperforms in regression tasks.  
132 In classification tasks, it typically reaches a minimum validation loss comparable to  
133 the other strategies, but requires substantially longer training time to do so. Finally,  
134 introducing a 20% dropout rate has little effect on performance for either regression  
135 or classification tasks.

136 Overall, although delayed unfreezing can be beneficial in specific classification  
137 scenarios, the default strategy of treating the pretrained D-MPNN and the newly ini-  
138 tialized FNN uniformly during fine-tuning emerges as a robust and effective approach  
139 for **ChemEleon**. This approach is used throughout the present study and suggested  
140 for all downstream user applications of **ChemEleon**.

141 **References**

- 142 [1] Gaulton, A., Bellis, L.J., Bento, A.P., Chambers, J., Davies, M., Hersey, A., Light,  
143 Y., McGlinchey, S., Michalovich, D., Al-Lazikani, B., Overington, J.P.: ChEMBL:  
144 a large-scale bioactivity database for drug discovery. *Nucleic Acids Research*  
145 **40**(D1), 1100–1107 (2011) <https://doi.org/10.1093/nar/gkr777>
- 146 [2] Tilborg, D., Alenicheva, A., Grisoni, F.: Exposing the limitations of molecu-  
147 lar machine learning with activity cliffs. *Journal of Chemical Information and*  
148 *Modeling* **62**(23), 5938–5951 (2022) <https://doi.org/10.1021/acs.jcim.2c01073>
- 149 [3] Heid, E., Greenman, K.P., Chung, Y., Li, S.-C., Graff, D.E., Vermeire, F.H.,  
150 Wu, H., Green, W.H., McGill, C.J.: Chemprop: A machine learning package for  
151 chemical property prediction. *Journal of Chemical Information and Modeling*  
152 **64**(1), 9–17 (2023) <https://doi.org/10.1021/acs.jcim.3c01250>
- 153 [4] Wu, Z., Ramsundar, B., Feinberg, E.N., Gomes, J., Geniesse, C., Pappu, A.S.,  
154 Leswing, K., Pande, V.: Moleculenet: a benchmark for molecular machine learn-  
155 ing. *Chemical Science* **9**(2), 513–530 (2018) <https://doi.org/10.1039/c7sc02664a>
- 156 [5] Ash, J.R., Wognum, C., Rodríguez-Pérez, R., Aldeghi, M., Cheng, A.C., Clevert,  
157 D.-A., Engkvist, O., Fang, C., Price, D.J., Hughes-Oliver, J.M., Walters, W.P.:  
158 Practically significant method comparison protocols for machine learning in small  
159 molecule drug discovery. *Journal of Chemical Information and Modeling* **65**(18),  
160 9398–9411 (2025) <https://doi.org/10.1021/acs.jcim.5c01609>
- 161 [6] Ball, N., Madden, J., Pains, A., Mathea, M., Palmer, A.D., Sperber, S., Har-  
162 tung, T., Ravenzwaay, B.: Key read across framework components and biology  
163 based improvements. *Mutation Research/Genetic Toxicology and Environmen-  
164 tal Mutagenesis* **853**, 503172 (2020) [https://doi.org/10.1016/j.mrgentox.2020.](https://doi.org/10.1016/j.mrgentox.2020.503172)  
165 [503172](https://doi.org/10.1016/j.mrgentox.2020.503172)
- 166 [7] O’Boyle, N.M., Sayle, R.A.: Comparing structural fingerprints using a literature-  
167 based similarity benchmark. *Journal of Cheminformatics* **8**(1) (2016) [https://doi.](https://doi.org/10.1186/s13321-016-0148-0)  
168 [org/10.1186/s13321-016-0148-0](https://doi.org/10.1186/s13321-016-0148-0)
- 169 [8] Carhart, R.E., Smith, D.H., Venkataraghavan, R.: Atom pairs as molecular  
170 features in structure-activity studies: definition and applications. *Journal of*  
171 *Chemical Information and Computer Sciences* **25**(2), 64–73 (1985) [https://doi.](https://doi.org/10.1021/ci00046a002)  
172 [org/10.1021/ci00046a002](https://doi.org/10.1021/ci00046a002)
- 173 [9] Nilakantan, R., Bauman, N., Dixon, J.S., Venkataraghavan, R.: Topological tor-  
174 sion: a new molecular descriptor for sar applications. comparison with other  
175 descriptors. *Journal of Chemical Information and Computer Sciences* **27**(2), 82–85  
176 (1987) <https://doi.org/10.1021/ci00054a008>
- 177 [10] Morgan, H.L.: The generation of a unique machine description for chemical

- 178 structures-a technique developed at chemical abstracts service. *Journal of Chem-*  
179 *ical Documentation* **5**(2), 107–113 (1965) <https://doi.org/10.1021/c160017a018>
- 180 [11] Landrum, G.: rdkit. Zenodo (2025). <https://doi.org/10.5281/ZENODO.591637> .  
181 <https://zenodo.org/doi/10.5281/zenodo.591637>
- 182 [12] Gedeck, P., Rohde, B., Bartels, C.: Qsar - how good is it in practice? comparison  
183 of descriptor sets on an unbiased cross section of corporate data sets. *Journal of*  
184 *Chemical Information and Modeling* **46**(5), 1924–1936 (2006) [https://doi.org/10.](https://doi.org/10.1021/ci050413p)  
185 [1021/ci050413p](https://doi.org/10.1021/ci050413p)
- 186 [13] Durant, J.L., Leland, B.A., Henry, D.R., Nourse, J.G.: Reoptimization of mdl  
187 keys for use in drug discovery. *Journal of Chemical Information and Computer*  
188 *Sciences* **42**(6), 1273–1280 (2002) <https://doi.org/10.1021/ci010132r>
- 189 [14] Orlov, A.A., Akhmetshin, T.N., Horvath, D., Marcou, G., Varnek, A.: From high  
190 dimensions to human insight: Exploring dimensionality reduction for chemical  
191 space visualization. *Molecular Informatics* **44**(1) (2024) [https://doi.org/10.1002/](https://doi.org/10.1002/minf.202400265)  
192 [minf.202400265](https://doi.org/10.1002/minf.202400265)
- 193 [15] Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O.,  
194 Blondel, M., Müller, A., Nothman, J., Louppe, G., Prettenhofer, P., Weiss, R.,  
195 Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot,  
196 M., Duchesnay, E.: Scikit-learn: Machine learning in python (2012) [https://doi.](https://doi.org/10.48550/ARXIV.1201.0490)  
197 [org/10.48550/ARXIV.1201.0490](https://doi.org/10.48550/ARXIV.1201.0490)
- 198 [16] Falcon, W., The PyTorch Lightning team: PyTorch Lightning. [https://doi.org/](https://doi.org/10.5281/zenodo.3828935)  
199 [10.5281/zenodo.3828935](https://doi.org/10.5281/zenodo.3828935) . <https://github.com/Lightning-AI/lightning>
- 200 [17] Rong, Y., Bian, Y., Xu, T., Xie, W., Wei, Y., Huang, W., Huang, J.: Self-  
201 supervised graph transformer on large-scale molecular data (2020) [https://doi.](https://doi.org/10.48550/ARXIV.2007.02835)  
202 [org/10.48550/ARXIV.2007.02835](https://doi.org/10.48550/ARXIV.2007.02835)
- 203 [18] Méndez-Lucio, O., Nicolaou, C.A., Earnshaw, B.: MolE: a foundation model for  
204 molecular graphs using disentangled attention. *Nature Communications* **15**(1)  
205 (2024) <https://doi.org/10.1038/s41467-024-53751-y>
- 206 [19] Sypetkowski, M., Wenkel, F., Poursafaei, F., Dickson, N., Suri, K., Fradkin, P.,  
207 Beaini, D.: On the scalability of GNNs for molecular graphs. In: *The Thirty-eighth*  
208 *Annual Conference on Neural Information Processing Systems* (2024). [https://](https://openreview.net/forum?id=klqhrq7fvB)  
209 [openreview.net/forum?id=klqhrq7fvB](https://openreview.net/forum?id=klqhrq7fvB)
- 210 [20] Beaini, D., Huang, S., Cunha, J.A., Li, Z., Moisescu-Pareja, G., Dymov, O.,  
211 Maddrell-Mander, S., McLean, C., Wenkel, F., Müller, L., Mohamud, J.H., Parviz,  
212 A., Craig, M., Koziarski, M., Lu, J., Zhu, Z., Gabellini, C., Klaser, K., Dean, J.,  
213 Wognum, C., Sypetkowski, M., Rabusseau, G., Rabbany, R., Tang, J., Morris,  
214 C., Ravanelli, M., Wolf, G., Tossou, P., Mary, H., Bois, T., Fitzgibbon, A.W.,

- 215 Banaszewski, B., Martin, C., Masters, D.: Towards foundational models for molec-  
216 ular learning on large-scale multi-task datasets. In: The Twelfth International  
217 Conference on Learning Representations (2024). [https://openreview.net/forum?](https://openreview.net/forum?id=Zc2aIcucwc)  
218 [id=Zc2aIcucwc](https://openreview.net/forum?id=Zc2aIcucwc)
- 219 [21] Fallani, A., Nugmanov, R., Arjona-Medina, J., Wegner, J.K., Tkatchenko, A.,  
220 Chernichenko, K.: Pretraining graph transformers with atom-in-a-molecule quan-  
221 tum properties for improved admet modeling. *Journal of Cheminformatics* **17**(1)  
222 (2025) <https://doi.org/10.1186/s13321-025-00970-0>
- 223 [22] Svetnik, V., Liaw, A., Tong, C., Culberson, J.C., Sheridan, R.P., Feuston, B.P.:  
224 Random Forest: A classification and regression tool for compound classification  
225 and qsar modeling. *Journal of Chemical Information and Computer Sciences*  
226 **43**(6), 1947–1958 (2003) <https://doi.org/10.1021/ci034160g>
- 227 [23] Burns, J.W., Green, W.H.: Generalizable, fast, and accurate DeepQSPR with  
228 fastprop. *Journal of Cheminformatics* **17**(1) (2025) <https://doi.org/10.1186/s13321-025-01013-4>
- 230 [24] Kläser, K., Banaszewski, B., Maddrell-Mander, S., McLean, C., Müller, L., Parviz,  
231 A., Huang, S., Fitzgibbon, A.: *MiniMol*: A parameter-efficient foundation model  
232 for molecular learning (2024) <https://doi.org/10.48550/ARXIV.2404.14986>
- 233 [25] Peters, M.E., Ruder, S., Smith, N.A.: To tune or not to tune? adapting pre-  
234 trained representations to diverse tasks. In: Proceedings of the 4th Workshop on  
235 Representation Learning for NLP (RepL4NLP@ACL), pp. 7–14. Association for  
236 Computational Linguistics, Florence, Italy (2019). <https://doi.org/10.18653/v1/W19-4302>
- 238 [26] Kumar, A., Raghunathan, A., Jones, R., Ma, T., Liang, P.: Fine-tuning can distort  
239 pretrained features and underperform out-of-distribution (2022) <https://doi.org/10.48550/ARXIV.2202.10054>
- 241 [27] Belkin, M., Hsu, D., Ma, S., Mandal, S.: Reconciling modern machine-learning  
242 practice and the classical bias–variance trade-off. *Proceedings of the National  
243 Academy of Sciences* **116**(32), 15849–15854 (2019) <https://doi.org/10.1073/pnas.1903070116>
- 245 [28] Zhang, C., Bengio, S., Hardt, M., Recht, B., Vinyals, O.: Understanding deep  
246 learning (still) requires rethinking generalization. *Communications of the ACM*  
247 **64**(3), 107–115 (2021) <https://doi.org/10.1145/3446776>