

---

# Meta-Learning GNN Initializations for Low-Resource Molecular Property Prediction

---

Cuong Q. Nguyen<sup>1</sup> Constantine Kretsoulas<sup>2</sup> Kim M. Branson<sup>1</sup>

## Abstract

Building *in silico* models to predict chemical properties and activities is a crucial step in drug discovery. However, limited labeled data often hinders the application of deep learning in this setting. Meanwhile advances in meta-learning have enabled state-of-the-art performances in few-shot learning benchmarks, naturally prompting the question: *Can meta-learning improve deep learning performance in low-resource drug discovery projects?* In this work, we assess the transferability of graph neural networks initializations learned by the Model-Agnostic Meta-Learning (MAML) algorithm – and its variants FO-MAML and ANIL – for chemical properties and activities tasks. Using the ChEMBL20 dataset to emulate low-resource settings, our benchmark shows that meta-initializations perform comparably to or outperform multi-task pre-training baselines on 16 out of 20 in-distribution tasks and on all out-of-distribution tasks, providing an average improvement in AUPRC of 11.2% and 26.9% respectively. Finally, we observe that meta-initializations consistently result in the best performing models across fine-tuning sets with  $k \in \{16, 32, 64, 128, 256\}$  instances.

## 1. Introduction

Drug discovery is a multi-parameter optimization process requiring efficient exploration of chemical space for compounds with desired properties. In a typical project, medicinal chemists propose structural changes to compounds in an effort to improve their therapeutic effects without compromising other properties. Validating these changes are costly – e.g. compounds need to be purchased or synthesized, as-

says need to be developed and validated – and thus *in silico* models are often used to prioritize experiments. Following the Merck Molecular Activity Challenge, there has been significant interest in applying deep learning to property prediction. More recently, by directly learning molecular features from chemical graphs, novel architectures in the graph neural networks family have demonstrated improved predictions in quantum chemistry and various property prediction benchmarks (Lusci et al., 2013; Duvenaud et al., 2015; Kearnes et al., 2016; Gilmer et al., 2017; Feinberg et al., 2018; Yang et al., 2019).

The successes of deep learning, however, hinge on an abundance of data: For instance, ImageNet (Deng et al., 2009) contains over 14M images and the English Wikipedia database commonly used to pre-train language models has over 2,500M words. On the contrary, labeled scientific data in drug discovery projects often consists of many small, sparse, and heavily biased datasets, consequently limiting the applications of deep learning in this setting. Recent works approach this problem by using pre-training and multitask learning to leverage data from multiple sources (Ram-sundar et al., 2015; Wenzel et al., 2019; Hu et al., 2019).

In parallel, the problem of learning in low-data domain has been tackled vehemently by the few-shot learning community. A prominent solution is the meta-learning paradigm, which aims to learn a learner that is efficient at adapting to new task (Thrun & Pratt, 1998; Vilalta & Drissi, 2002; Vanschoren, 2018). Matching Networks (Vinyals et al., 2016), a member of this family, have been previously applied to property prediction in one-shot learning settings by Altae-Tran et al. (2017). A related approach is the Model-Agnostic Meta-Learning (MAML) algorithm (Finn et al., 2017), which has been particularly successful at producing state-of-the-arts results on few-shots classification, regression, and reinforcement learning benchmarks, resulting in numerous follow-ups that expand on this elegant framework.

In this work, we evaluate gated graph neural networks initializations learned by MAML, its first-order approximation (FO-MAML), and the Almost-No-Inner-Loop (ANIL) variant (Raghu et al., 2020) for transfer learning to low-resource molecular properties and activities tasks. Specifically we aim to answer the following questions:

<sup>1</sup>GlaxoSmithKline, Artificial Intelligence & Machine Learning

<sup>2</sup>GlaxoSmithKline, Data & Computational Sciences. Correspondence to: Cuong Q. Nguyen <cuong.q.nguyen@gsk.com>.

1. Does meta-initializations offer improvements over multitask pre-training in this setting?
2. How little data can meta-initializations learn efficiently from?

Using ChEMBL20 (Bento et al., 2014), performances of meta-initializations on in- and out-of-distribution tasks are benchmarked with multitask pre-training baselines, showing favorable performances across fine-tuning set sizes of  $k \in \{16, 32, 64, 128, 256\}$  instances.

## 2. Background

**MAML, FO-MAML, and ANIL** MAML’s approach to few-shot learning is to directly optimize for a set of initial parameters that is efficient at learning from new data. The algorithm consists of an outer loop that learns an initialization  $\theta_0$ , and an inner loop that adapts  $\theta_0$  to new tasks. In this setting, a set of tasks  $\{T_1, T_2, \dots, T_K\}$  – denoted as  $\mathcal{T}^{tr}$  – is available to obtain  $\theta_0$ , from which we would like to learn the set of tasks  $\mathcal{T}^{test}$ . Following the nomenclature in Finn et al. (2017), we call the process of obtaining  $\theta_0$  *meta-training*, and the process of adapting to  $\mathcal{T}^{test}$  *meta-testing*. More formally, we define a task  $T_j$  with  $K$  instances as  $T = \{(x_i, y_i) \mid i \in \{1, \dots, K\}\}$ , which is divided into a training set  $D_{T_j}^{tr}$  and a test set  $D_{T_j}^{test}$ , also referred to in the literature as the support and query set, respectively. The inner loop adaptation to  $T_j$  for a neural network  $f$  parameterized by  $\theta$  using gradient descent is expressed as

$$\theta_N^j = \theta_{N-1}^j - \alpha \nabla_{\theta} \mathcal{L}_{D_{T_j}^{tr}}(f_{\theta_{N-1}^j})$$

where  $\theta_N^j$  denotes the parameters of  $f$  after  $N$  steps toward task  $T_j$ ,  $\alpha$  is the inner loop learning rate, and  $\mathcal{L}_{D_{T_j}^{tr}}$  is the loss on the training set of task  $T_j$ . The loss is calculated using  $f$  after  $N - 1$  updates. The inner loop is repeated for a batch of  $B$  tasks sampled from  $\mathcal{T}^{tr}$ .

For the outer loop, the meta-loss is defined as the sum of task-specific losses after inner loop updates:

$$\mathcal{L}_{meta}(\theta_0) = \sum_{j=1}^B \mathcal{L}_{D_{T_j}^{test}}(f_{\theta_N^j})$$

The task-specific loss  $\mathcal{L}_{D_{T_j}^{test}}$  is calculated on the test set of task  $T_j$ . We then minimize the meta-loss using stochastic gradient descent to optimize the initialization  $\theta_0$ , with updates expressed by

$$\theta_0 \leftarrow \theta_0 - \eta \nabla_{\theta} \mathcal{L}_{meta}(\theta_0)$$

where  $\eta$  is the outer loop learning rate. Intuitively, the meta-loss  $\mathcal{L}_{meta}(\theta_0)$  measures how well  $\theta_0$  adapts to new tasks,

Table 1. Distribution of task types in each split.  $A$ ,  $T$ ,  $P$ ,  $B$ , and  $F$  denote ADME, Toxicity, Physicochemical, Binding, and Functional as found in ChEMBL20.  $\mathcal{T}^{tr}$  and  $\mathcal{T}^{val}$  only contain  $B$  and  $F$  tasks, while  $\mathcal{T}^{test}$  contains all 5 task types.

	$A$	$T$	$P$	$B$	$F$
$\mathcal{T}^{tr}$	0	0	0	126	737
$\mathcal{T}^{val}$	0	0	0	10	10
$\mathcal{T}^{test}$	1	1	1	10	10

and minimizing this loss enables the algorithm to learn good initial parameters.

Updating  $\theta_0$  is computationally expensive since it requires the use of second-order derivatives to compute  $\nabla_{\theta} \mathcal{L}_{meta}(\theta_0)$ . FO-MAML sidesteps this problem by omitting the second-order terms, effectively ignoring the inner loop gradients. On the other hand, Raghu et al. (2020) proposes the ANIL algorithm, which reduces the number of second-order gradients required by limiting inner loop adaptation to only the penultimate layer of the network. ANIL and FO-MAML have both demonstrated significant speedup over MAML.

**Graph Neural Networks** The graph neural networks framework enables representation learning on graph structured data by learning node-level representations which are aggregated to form graph-level representations. Throughout our experiments, we use a variant of the Gated Graph Neural Network (GGNN) architecture (Li et al., 2017), a member of the message passing neural network (MPNN) family (Gilmer et al., 2017). Similar to other MPNNs, the GGNN architecture operates in two phases: a message passing phase and a readout phase. For an undirected graph  $\mathcal{G}$  with  $V$  nodes where each node has  $F$  features, the message passing phase updates the hidden representation of node  $v$  at layer  $t$  according to

$$m_v^{t+1} = A_{e_{vw}} h_v^t + \sum_{w \in N(v)} A_{e_{vw}} h_w^t$$

$$h_v^{t+1} = \text{GRU}(h_v^t, m_v^{t+1})$$

where  $A_{e_{vw}} \in \mathbb{R}^{F \times F}$  is an edge-specific learnable weight matrix,  $N(v)$  denotes neighbors of  $v$ , GRU is the Gated Recurrent Unit (Cho et al., 2014), and  $m_v \in \mathbb{R}^F$  is a message used to update the hidden representation of node  $v$  denoted by  $h_v \in \mathbb{R}^F$ . Computing the message  $m_v$  is often interpreted as aggregating information across central and neighboring nodes. A deviation from Li et al. (2017) comes in our choice to remove weight sharing between GRUs in different layers. Following  $T$  updates, the readout phase pools node representations according to

$$\hat{y} = \text{MLP}\left(\sum_{v \in G} h_v^T\right)$$

Table 2. Performance on in-distribution tasks measured in AUPRC. The top and bottom halves of the table are tasks with type  $B$  and  $F$ , respectively. Mean and standard deviation are obtained from 25 repeats (see Evaluation in Section 3 for details). The best and second best values are in bold and regular text, respectively. Statistically significant difference from the next best is denoted by (\*).

CHEMBL ID	k-NN	FINETUNE-ALL	FINETUNE-TOP	FO-MAML	ANIL	MAML
2363236	0.316 ± 0.007	0.328 ± 0.028	0.329 ± 0.023	<b>0.337 ± 0.019</b>	0.325 ± 0.008	0.332 ± 0.013
1614469	0.438 ± 0.023	0.470 ± 0.034	0.490 ± 0.033	0.489 ± 0.019	0.446 ± 0.044	<b>0.507 ± 0.030</b>
2363146	0.559 ± 0.026	<b>0.626 ± 0.037</b>	<b>0.653 ± 0.029</b>	0.555 ± 0.017	0.506 ± 0.034	0.595 ± 0.051
2363366	0.511 ± 0.050	0.567 ± 0.039	0.551 ± 0.048	0.546 ± 0.037	<b>0.570 ± 0.031</b>	<b>0.598 ± 0.041</b>
2363553	<b>0.739 ± 0.007</b>	0.724 ± 0.015	<b>0.737 ± 0.023</b>	0.694 ± 0.011	0.686 ± 0.020	0.691 ± 0.013
1963818	0.607 ± 0.041	<b>0.708 ± 0.036</b>	0.595 ± 0.142	0.677 ± 0.026	0.692 ± 0.081	<b>0.745 ± 0.048</b>
1963945	0.805 ± 0.031	<b>0.848 ± 0.034</b>	0.835 ± 0.036	0.779 ± 0.039	0.753 ± 0.033	0.836 ± 0.023
1614423	0.503 ± 0.044	0.628 ± 0.058	0.642 ± 0.063	<b>0.760 ± 0.024</b>	0.730 ± 0.077	<b>0.837 ± 0.036*</b>
2114825	0.679 ± 0.027	0.739 ± 0.050	0.732 ± 0.051	<b>0.837 ± 0.042</b>	0.759 ± 0.078	<b>0.885 ± 0.014*</b>
1964116	0.709 ± 0.042	0.758 ± 0.044	0.769 ± 0.048	0.895 ± 0.023	<b>0.903 ± 0.016</b>	<b>0.912 ± 0.013</b>
2155446	0.471 ± 0.008	0.473 ± 0.017	0.476 ± 0.013	<b>0.497 ± 0.024</b>	0.478 ± 0.020	<b>0.500 ± 0.017</b>
1909204	0.538 ± 0.023	0.589 ± 0.031	0.577 ± 0.039	<b>0.592 ± 0.043</b>	0.547 ± 0.029	<b>0.601 ± 0.027</b>
1909213	0.694 ± 0.009	<b>0.742 ± 0.015</b>	<b>0.759 ± 0.012</b>	0.698 ± 0.024	0.694 ± 0.025	0.729 ± 0.013
3111197	0.617 ± 0.028	0.663 ± 0.066	0.673 ± 0.071	0.636 ± 0.036	<b>0.737 ± 0.035</b>	<b>0.746 ± 0.045</b>
3215171	0.480 ± 0.042	0.552 ± 0.043	0.551 ± 0.045	<b>0.729 ± 0.031</b>	0.700 ± 0.050	<b>0.764 ± 0.019</b>
3215034	0.474 ± 0.072	0.540 ± 0.156	0.455 ± 0.189	<b>0.819 ± 0.048</b>	0.681 ± 0.042	0.805 ± 0.046
1909103	0.881 ± 0.026	<b>0.936 ± 0.013</b>	<b>0.921 ± 0.020</b>	0.877 ± 0.046	0.730 ± 0.055	0.900 ± 0.032
3215092	0.696 ± 0.038	0.777 ± 0.039	0.791 ± 0.042	<b>0.877 ± 0.028</b>	0.834 ± 0.026	<b>0.907 ± 0.017</b>
1738253	0.710 ± 0.048	0.860 ± 0.029	0.861 ± 0.025	<b>0.885 ± 0.033</b>	0.758 ± 0.111	<b>0.908 ± 0.011</b>
1614549	0.710 ± 0.035	0.850 ± 0.041	0.860 ± 0.051	<b>0.930 ± 0.022</b>	0.860 ± 0.034	<b>0.947 ± 0.014</b>
AVG. RANK	5.4	3.5	3.5	3.1	4.0	<b>1.7</b>

Table 3. Performance on out-of-distribution tasks measured in AUPRC. Mean and standard deviations are obtained from 25 repeats (see Evaluation in 3 for details). Notations are the same as Table 2.

CHEMBL ID	k-NN	FINETUNE-ALL	FINETUNE-TOP	FO-MAML	ANIL	MAML
1804798	0.338 ± 0.020	0.351 ± 0.026	0.357 ± 0.031	0.360 ± 0.017	0.361 ± 0.029	<b>0.367 ± 0.024</b>
2095143	0.256 ± 0.054	0.147 ± 0.046	0.281 ± 0.082	<b>0.562 ± 0.034</b>	<b>0.564 ± 0.037</b>	0.522 ± 0.054
918058	0.407 ± 0.138	0.559 ± 0.098	0.609 ± 0.076	0.506 ± 0.096	0.415 ± 0.163	<b>0.694 ± 0.082</b>
AVG. RANK	5.7	4.7	3.3	3.0	2.7	1.7

to calculate the neural network output  $\hat{y}$ . Using sum as the readout operation is the second deviation from Li et al. (2017), and has been shown to have maximal expressive power over mean and max aggregators (Xu et al., 2018).

### 3. Experimental Settings

**ChEMBL20 Dataset** We evaluate the effectiveness of meta-initializations for low-resource tasks using a subset of ChEMBL20. More specifically, the dataset processed by Mayr et al. (2018) is filtered for tasks with at least 128 instances. The resulting dataset contains 902 binary classification tasks from 5 distinct task types: ADME ( $A$ ), Toxicity ( $T$ ), Physicochemical ( $P$ ), Binding ( $B$ ), and Functional ( $F$ ).

The tasks are further divided into  $\mathcal{T}^{tr}$ ,  $\mathcal{T}^{val}$ , and  $\mathcal{T}^{test}$ .  $\mathcal{T}^{val}$  consists of 10 randomly selected  $B$  and  $F$  tasks.  $\mathcal{T}^{test}$  consists of all  $A$ ,  $T$ , and  $P$  tasks in addition to 10 random  $B$  and  $F$  tasks. The rest of  $B$  and  $F$  tasks are included in  $\mathcal{T}^{tr}$

for meta-training. A summary of task type distribution is shown in Table 1. For baselines,  $\mathcal{T}^{tr}$  and  $\mathcal{T}^{val}$  are combined and split into  $D_{baseline}^{tr}$  for training and  $D_{baseline}^{val}$  for early stopping. This setup gives the baselines access to more tasks than MAML, FO-MAML, and ANIL during training.

Each molecule is represented as an undirected graph where nodes and edges are atoms and bonds. We use the OpenEye Toolkit to generate 75 atomic features for each node, similar to those provided by DeepChem (Ramsundar et al., 2019).

**Baselines** We include the Finetune-All, Finetune-Top, and k-NN baselines as proposed by Triantafillou et al. (2019). All baselines start with training a multi-task GGNN on  $D_{baseline}^{tr}$ . The k-NN baseline uses the activations from the penultimate layer of pre-trained model to perform classification from 3 nearest neighbors. Finetune-Top reinitializes and trains the penultimate layer while Finetune-All updates all parameters in the model.

To ensure the baselines are competitive, we perform hyperparameter tuning using the Tree-of-Parzen Estimator implementation of Hyperopt (Bergstra et al., 2015) to optimize performance on  $D_{baseline}^{val}$ . Appendix 1 provides details of the process and the resulting hyperparameters.

**Meta-Learning** The same GGNN architecture as the baselines is used for all three meta-learning algorithms. Training hyperparameters are hand-tuned for performance on  $\mathcal{T}^{val}$  (see Appendix 2 for details). We use the Learn2Learn (Arnold et al., 2019) and PyTorch (Paszke et al., 2019) libraries for our implementation.

**Evaluation** For each  $T_j$  in  $\mathcal{T}^{test}$ , we fine-tune initializations on  $k$  randomly selected instances from  $D_{T_j}^{tr}$  using the Adam optimizer with learning rate of  $10^{-4}$  and batch size of  $b = \min(64, k)$ . We use  $D_{T_j}^{val}$  for early stopping with patience of 10 epochs and collect performances on  $D_{T_j}^{test}$ . We use  $B$  and  $T$  tasks to assess *in-distribution* performance, and  $A$ ,  $T$ , and  $P$  tasks for *out-of-distribution* performance. For each method, the procedure is repeated 25 times with 5 different sets of  $k$  instances and 5 random seeds.

## 4. Results & Discussions

**Performances on  $\mathcal{T}^{test}$**  The performance of each method on 23 test tasks when  $k = 128$  instances are used for fine-tuning is reported in Table 2 and 3. Since random splits have been shown to be overly optimistic in scientific applications (Kearnes et al., 2017; Wu et al., 2018), we emphasize relative ranking over absolute performance throughout our benchmark. We observe that meta-initializations generally exhibit similar or better performances over baselines despite having been trained on fewer tasks. For in-distribution tasks, MAML performs comparably to or outperforms other methods on 16 out of 20 tasks, 2 of which shows significant improvement over the next best method, making it the top performer with an average rank of 1.7. ANIL and FO-MAML, while benefitting from a shorter training time (Appendix 3), rank 3.1 and 4.0 on average, respectively. Similar to observations by Triantafillou et al. (2019), Finetune-All and Finetune-Top baselines prove to be strong competitors, both ranking above ANIL in our benchmark. Given their significantly shorter training time, we suspect both baselines to remain crucial in compute-limited settings. In out-of-distribution settings, meta-initializations outperform baselines on all 3 tasks. Again, MAML is ranked as the best method, followed by ANIL and FO-MAML. Overall, compared to the best baselines, meta-initializations learned by MAML provide an average increase in AUPRC of 11.2% for in-distribution tasks and 26.9% out-of-distributions tasks.

**Effect of Fine-tuning Set Size** From  $\mathcal{T}^{test}$ , we select all tasks with at least 256 instances in  $D_{T_j}^{tr}$ , resulting in 18 tasks available for evaluation (as opposed to 9 when a threshold of

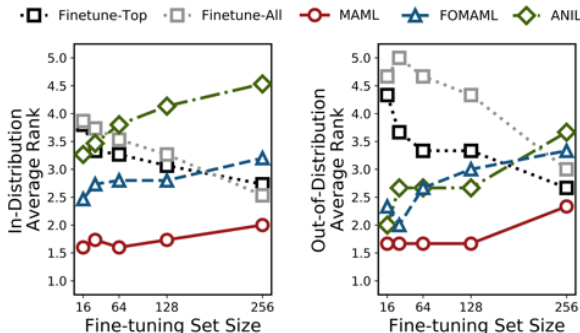


Figure 1. Average ranks of each method performance after fine-tuning with  $k \in \{16, 32, 64, 128, 256\}$  instances for in- (left) and out-of-distribution tasks (right). Rankings are based on mean AUPRC measured from five random seeds. MAML is consistently ranked as the best method across all  $k$  for in-distribution and out-of-distribution tasks, respectively.

512 instances is used). The average ranking of each method after fine-tuning on  $k \in \{16, 32, 64, 128, 256\}$  instances is reported in Figure 1 (see Appendix 4 for performance on each task). As the best performing baselines from the previous experiment, Finetune-Top and Finetune-All are selected for comparison. We observe that the baselines benefit greatly from having more data, with Finetune-All rising from fifth to second in in-distribution tasks and Finetune-Top rising from fourth to second in out-of-distribution tasks. Nonetheless, MAML remains the best method, consistently ranked first across fine-tuning set sizes for both sets of tasks.

## 5. Conclusion & Future Directions

In this work, we explore meta-learning as a tool for learning to predict chemical properties and activities in low-resource settings. Emulating this setting using the ChEMBL20 dataset, we demonstrate that GGNN’s initializations learned by MAML perform comparably to or outperform multitask pre-training baselines on 16 out of 20 in-distribution tasks and on all 3 out-of-distribution tasks. Improved performances of meta-initializations are further shown to remain consistent across fine-tuning sets of size  $k \in \{16, 32, 64, 128, 256\}$ .

While the ChEMBL20 dataset enables differentiating between in- and out-of-distribution tasks, we recognize that its chemical space is biased towards compounds which have been reviewed and selected for publications. Moreover, our benchmark does not include initializations obtained using self- and un-supervised approaches such as those described in Veličković et al. (2018), Hu et al. (2019), and Sun et al. (2020). We leave experiments with additional datasets and methods to future work. Overall, we believe our contributions open opportunities in applying deep learning to ongoing drug discovery projects where limited data is available.

## Acknowledgment

We would like to thank Jijie Zhang, Stephen Young, Robert Woodruff, and Darren Green for helpful discussions during the preparation of this manuscript.

## References

- Altae-Tran, H., Ramsundar, B., Pappu, A. S., and Pande, V. Low Data Drug Discovery with One-Shot Learning. *ACS Central Science*, 3(4):283–293, April 2017. ISSN 2374-7943. doi: 10.1021/acscentsci.6b00367.
- Arnold, S. M. R., Mahajan, P., Datta, D., and Ian, B. learn2learn, September 2019. URL <https://github.com/learnables/learn2learn>.
- Bento, A. P., Gaulton, A., Hersey, A., Bellis, L. J., Chambers, J., Davies, M., Krüger, F. A., Light, Y., Mak, L., McGlinchey, S., Nowotka, M., Papadatos, G., Santos, R., and Overington, J. P. The ChEMBL bioactivity database: an update. *Nucleic Acids Research*, 42(D1): D1083–D1090, January 2014. ISSN 0305-1048. doi: 10.1093/nar/gkt1031.
- Bergstra, J., Komer, B., Eliasmith, C., Yamins, D., and Cox, D. D. Hyperopt: a Python library for model selection and hyperparameter optimization. *Computational Science & Discovery*, 8(1):014008, July 2015. ISSN 1749-4699. doi: 10.1088/1749-4699/8/1/014008.
- Cho, K., Merriënboer, B. v., Gulcehre, C., Bougares, F., Schwenk, H., and Bengio, Y. Learning phrase representations using RNN encoder-decoder for statistical machine translation. *Conference on Empirical Methods in Natural Language Processing (EMNLP 2014)*, 2014.
- Deng, J., Dong, W., Socher, R., Li, L.-J., Li, K., and Fei-Fei, L. ImageNet: A Large-Scale Hierarchical Image Database. In *CVPR09*, 2009.
- Duvenaud, D., Maclaurin, D., Aguilera-Iparraguirre, J., Gómez-Bombarelli, R., Hirzel, T., Aspuru-Guzik, A., and Adams, R. P. Convolutional Networks on Graphs for Learning Molecular Fingerprints. *arXiv:1509.09292 [cs, stat]*, November 2015. arXiv: 1509.09292.
- Feinberg, E. N., Sur, D., Wu, Z., Husic, B. E., Mai, H., Li, Y., Sun, S., Yang, J., Ramsundar, B., and Pande, V. S. PotentialNet for Molecular Property Prediction. *ACS Central Science*, 4(11):1520–1530, November 2018. ISSN 2374-7943. doi: 10.1021/acscentsci.8b00507.
- Finn, C., Abbeel, P., and Levine, S. Model-Agnostic Meta-Learning for Fast Adaptation of Deep Networks. In Precup, D. and Teh, Y. W. (eds.), *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pp. 1126–1135, International Convention Centre, Sydney, Australia, August 2017. PMLR.
- Gilmer, J., Schoenholz, S. S., Riley, P. F., Vinyals, O., and Dahl, G. E. Neural message passing for Quantum chemistry. In *Proceedings of the 34th International Conference on Machine Learning - Volume 70, ICML’17*, pp. 1263–1272, Sydney, NSW, Australia, August 2017. JMLR.org.
- Hu, W., Liu, B., Gomes, J., Zitnik, M., Liang, P., Pande, V., and Leskovec, J. Pre-training Graph Neural Networks. *arXiv:1905.12265 [cs, stat]*, May 2019. arXiv: 1905.12265.
- Kearnes, S., McCloskey, K., Berndl, M., Pande, V., and Riley, P. Molecular graph convolutions: moving beyond fingerprints. *Journal of Computer-Aided Molecular Design*, 30(8):595–608, August 2016. ISSN 1573-4951. doi: 10.1007/s10822-016-9938-8.
- Kearnes, S., Goldman, B., and Pande, V. Modeling Industrial ADMET Data with Multitask Networks. *arXiv:1606.08793 [stat]*, January 2017. arXiv: 1606.08793.
- Kingma, D. P. and Ba, J. Adam: A Method for Stochastic Optimization. *arXiv:1412.6980 [cs]*, January 2017. arXiv: 1412.6980.
- Li, Y., Tarlow, D., Brockschmidt, M., and Zemel, R. Gated Graph Sequence Neural Networks. *arXiv:1511.05493 [cs, stat]*, September 2017. arXiv: 1511.05493.
- Lusci, A., Pollastri, G., and Baldi, P. Deep Architectures and Deep Learning in Chemoinformatics: The Prediction of Aqueous Solubility for Drug-Like Molecules. *Journal of Chemical Information and Modeling*, 53(7):1563–1575, July 2013. ISSN 1549-9596. doi: 10.1021/ci400187y.
- Mayr, A., Klambauer, G., Unterthiner, T., Steijaert, M., Wegner, J. K., Ceulemans, H., Clevert, D.-A., and Hochreiter, S. Large-scale comparison of machine learning methods for drug target prediction on ChEMBL. *Chemical Science*, 9(24):5441–5451, June 2018. ISSN 2041-6539. doi: 10.1039/C8SC00148K.
- Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., Desmaison, A., Kopf, A., Yang, E., DeVito, Z., Raison, M., Tejani, A., Chilamkurthy, S., Steiner, B., Fang, L., Bai, J., and Chintala, S. PyTorch: An Imperative Style, High-Performance Deep Learning Library. In Wallach, H., Larochelle, H., Beygelzimer, A., Alché-Buc, F., Fox, E., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems 32*, pp. 8024–8035. Curran Associates, Inc., 2019.

- Raghu, A., Raghu, M., Bengio, S., and Vinyals, O. Rapid Learning or Feature Reuse? Towards Understanding the Effectiveness of MAML. In *International Conference on Learning Representations*, 2020.
- Ramsundar, B., Kearnes, S., Riley, P., Webster, D., Konerding, D., and Pande, V. Massively Multitask Networks for Drug Discovery. *arXiv:1502.02072 [cs, stat]*, February 2015. arXiv: 1502.02072.
- Ramsundar, B., Eastman, P., Walters, P., Pande, V., Leswing, K., and Wu, Z. *Deep Learning for the Life Sciences*. O'Reilly Media, 2019.
- Sun, F.-Y., Hoffmann, J., Verma, V., and Tang, J. InfoGraph: Unsupervised and Semi-supervised Graph-Level Representation Learning via Mutual Information Maximization. *arXiv:1908.01000 [cs, stat]*, January 2020. arXiv: 1908.01000.
- Thrun, S. and Pratt, L. (eds.). *Learning to Learn*. Springer US, 1998. ISBN 978-0-7923-8047-4. doi: 10.1007/978-1-4615-5529-2.
- Triantafillou, E., Zhu, T., Dumoulin, V., Lamblin, P., Evci, U., Xu, K., Goroshin, R., Gelada, C., Swersky, K., Manzagol, P.-A., and Larochelle, H. Meta-Dataset: A Dataset of Datasets for Learning to Learn from Few Examples. *arXiv:1903.03096 [cs, stat]*, October 2019. arXiv: 1903.03096.
- Vanschoren, J. Meta-Learning: A Survey. *arXiv:1810.03548 [cs, stat]*, October 2018. arXiv: 1810.03548.
- Veličković, P., Fedus, W., Hamilton, W. L., Liò, P., Bengio, Y., and Hjelm, R. D. Deep Graph Infomax. *arXiv:1809.10341 [cs, math, stat]*, December 2018. arXiv: 1809.10341.
- Vilalta, R. and Drissi, Y. A Perspective View and Survey of Meta-Learning. *Artificial Intelligence Review*, 18(2): 77–95, June 2002. ISSN 1573-7462. doi: 10.1023/A:1019956318069.
- Vinyals, O., Blundell, C., Lillicrap, T., kavukcuoglu, k., and Wierstra, D. Matching Networks for One Shot Learning. In Lee, D. D., Sugiyama, M., Luxburg, U. V., Guyon, I., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems 29*, pp. 3630–3638. Curran Associates, Inc., 2016.
- Wenzel, J., Matter, H., and Schmidt, F. Predictive Multitask Deep Neural Network Models for ADME-Tox Properties: Learning from Large Data Sets. *Journal of Chemical Information and Modeling*, 59(3):1253–1268, March 2019. ISSN 1549-9596. doi: 10.1021/acs.jcim.8b00785.
- Wu, Z., Ramsundar, B., Feinberg, E. N., Gomes, J., Geniesse, C., Pappu, A. S., Leswing, K., and Pande, V. MoleculeNet: a benchmark for molecular machine learning. *Chemical Science*, 9(2):513–530, January 2018. ISSN 2041-6539. doi: 10.1039/C7SC02664A.
- Xu, K., Hu, W., Leskovec, J., and Jegelka, S. How Powerful are Graph Neural Networks? *arXiv:1810.00826 [cs, stat]*, October 2018. arXiv: 1810.00826.
- 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., and Barzilay, R. Analyzing Learned Molecular Representations for Property Prediction. *Journal of Chemical Information and Modeling*, 59(8):3370–3388, August 2019. ISSN 1549-9596. doi: 10.1021/acs.jcim.9b00237.

---

# Appendix: Meta-Learning GNN Initializations for Low-Resource Molecular Property Prediction

---

Cuong Q. Nguyen<sup>1</sup> Constantine Kretsoulas<sup>2</sup> Kim M. Branson<sup>1</sup>

## 1. Baselines Hyperparameter Tuning

Using Hyperopt, we allow a maximum of 50 evaluations and provide the following search space:

- Number of GGNN layers: {3, 7, 9}
- Fully connected layer dimension: {1024, 2048}
- Batch size: {128, 256, 512}
- Learning rate:  $10^{\{-4.0, -3.75, -3.5, -3.25\}}$

The resulting architecture has 7 GGNN layers, 1 fully-connected layer with 1024 units, and Dropout applied with a probability of 0.2 at every layer except for the output layer. We use the Adam optimizer with a learning rate of  $10^{-3.75}$ , batch size of 512, and patience of 20 epochs for early stopping during pre-training.

## 2. Meta-Learning Hyperparameters

For MAML and ANIL we use an inner loop learning rate of 0.05, 2 inner gradient steps, and inner batch size of 32, while the outer loop has a learning rate of 0.003 and a batch size of 32. FO-MAML uses an outer loop learning rate of 0.0015.

## 3. Training Time

Training time was measured as the total time required to reach best performance on the  $D_{baseline}^{val}$  for baselines and  $\mathcal{T}^{val}$  for MAML, FO-MAML, and ANIL on 1 NVIDIA Tesla V100 GPU. We report the recorded times in Table ???. The mean and standard deviation are calculated by repeating the training process with five random seeds.

Table 1. Wall clock time to train each method

	TIME (HOURS)	SPEEDUP
MAML	$57.9 \pm 0.8$	1×
ANIL	$48.0 \pm 0.6$	1.2×
FO-MAML	$27.0 \pm 0.9$	2.1×
MULTI-TASK	$1.4 \pm 0.1$	41.4×

---

<sup>1</sup>GlaxoSmithKline, Artificial Intelligence & Machine Learning <sup>2</sup>GlaxoSmithKline, Data & Computational Sciences. Correspondence to: Cuong Q. Nguyen <cuong.q.nguyen@gsk.com>.

#### 4. Effect of Fine-tuning Set Size on Performances

We report the performances of each method on individual tasks below. Figure ?? show in-distribution tasks, while Figure ?? shows out-of-distribution tasks.

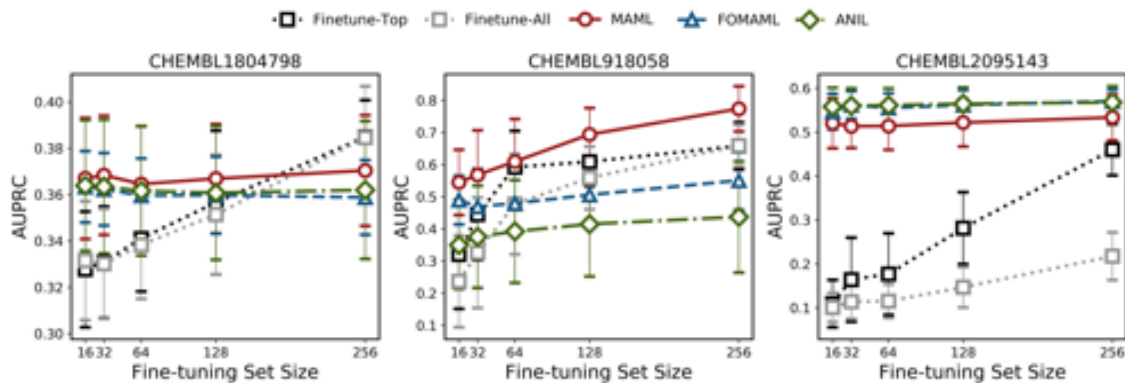


Figure 1. Performances on out-of-distribution tasks



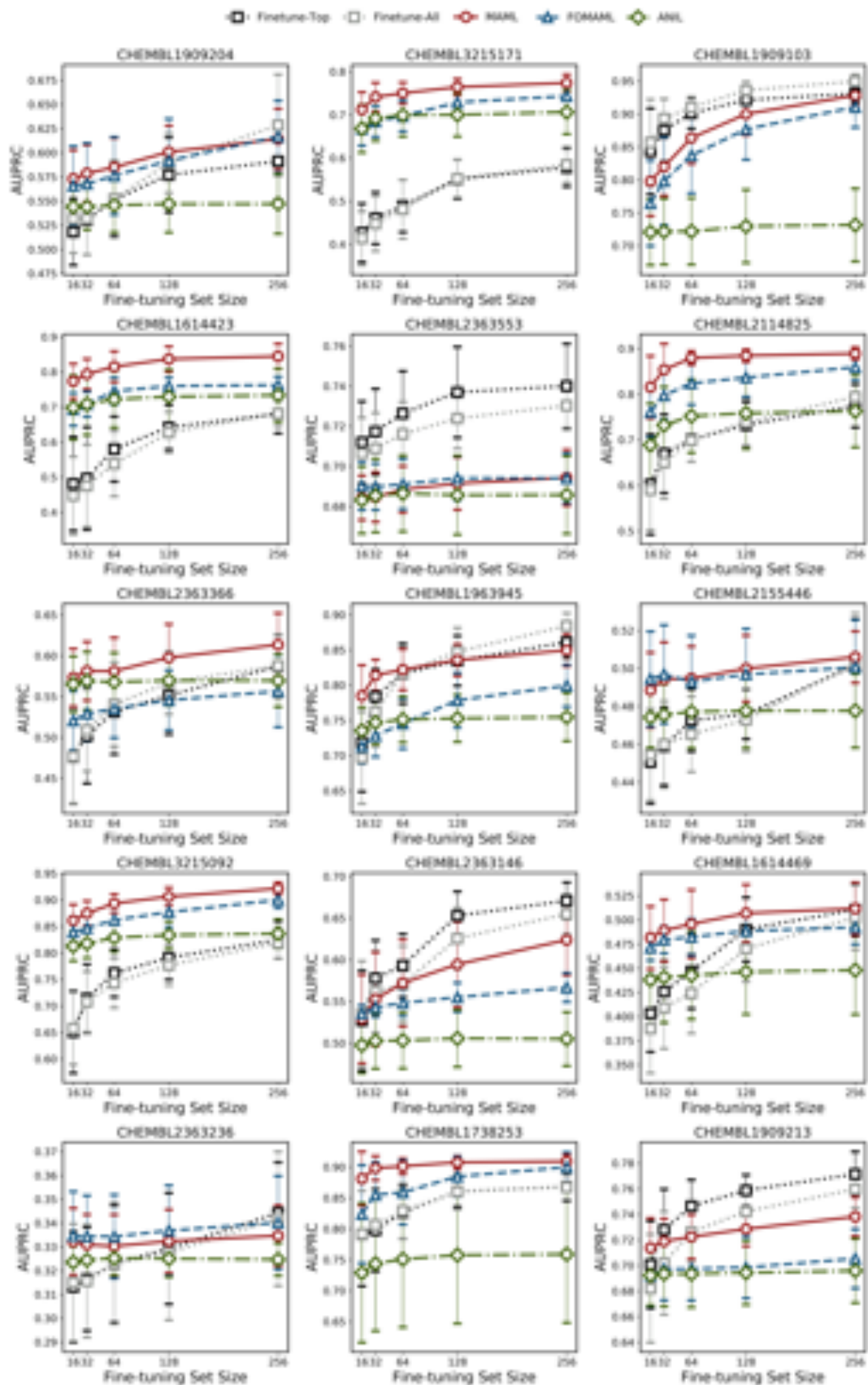


Figure 2. Performances on in-distribution tasks