

---

# RoFt-Mol: Benchmarking Robust Fine-Tuning with Molecular Graph Foundation Models

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

In the era of foundation models, fine-tuning pre-trained models for specific downstream tasks has become crucial. This drives the need for robust fine-tuning methods to address challenges such as model overfitting and sparse labeling. Molecular graph foundation models (MGFMs) face unique difficulties that complicate fine-tuning. These models are limited by smaller pre-training datasets and more severe data scarcity for downstream tasks, both of which require enhanced model generalization. Moreover, MGFMs must accommodate diverse objectives, including both regression and classification tasks. To better understand and improve fine-tuning techniques under these conditions, we classify eight fine-tuning methods into three mechanisms: weight-based, representation-based, and partial fine-tuning. We benchmark these methods on downstream regression and classification tasks across supervised and self-supervised pre-trained models in diverse labeling settings. This extensive evaluation provides valuable insights and informs the design of a refined robust fine-tuning method, ROFT-MOL. This approach combines the strengths of simple post-hoc weight interpolation with more complex weight ensemble fine-tuning methods, delivering improved performance across both task types while maintaining the ease of use inherent in post-hoc weight interpolation.

## 1 Introduction

In recent years, foundation models [1, 2] have achieved success in learning high-quality, general-purpose representations of images and text through pre-training on diverse datasets [3, 4, 5, 6, 7, 8]. To adapt these pre-trained models for downstream applications, additional training on task-specific data, known as fine-tuning, is often required. However, vanilla fine-tuning frequently encounters challenges, including model overfitting [9, 10, 11], catastrophic forgetting of pre-trained knowledge [12, 13, 14, 15], and distribution shifts between fine-tuned and test samples, which can lead to negative transfer [16, 17]. These challenges highlight the need for robust fine-tuning strategies [18, 19, 20, 21, 22, 23].

Recently, the advantages of foundation models have been extended to various scientific applications [24, 25, 26]. Among these, molecular graph foundation models (MGFMs) have gained significant attention for their promising potential in biochemistry [27, 28, 29, 30, 31, 32, 33, 34, 35, 36]. While MGFMs exhibit scaling behaviors similar to foundation models in other domains [37], they face unique challenges related to data and tasks.

A primary challenge stems from the significantly smaller pre-training datasets in this domain, typically consisting of at most  $O(100M)$  molecular samples, compared to the billions of samples used in other domains [38]. This limitation restricts the parameter scale of MGFMs ( $O(100M)$  parameters) and their generalization capacity [39, 40]. Furthermore, downstream tasks in this domain often involve limited data for fine-tuning, with datasets containing only tens or a few hundred labeled samples [41], exacerbating the difficulty of achieving robust model generalization. In addition to data constraints, many downstream tasks, such as molecular property prediction, are regression-based [42, 43]. These

tasks require models to capture fine-grained numerical patterns, which presents a distinct requirement compared to the coarse-grained feature reliance typical in classification tasks in CV and NLP. These factors collectively highlight the need for a careful examination of fine-tuning strategies for MGFMs and their appropriate improvement.

To answer this question, we introduce ROFT-MOL, a benchmark that evaluates existing fine-tuning methods across diverse molecular property prediction tasks. To explore factors influencing the fine-tuning (FT) performance of MGFMs, we categorize 8 FT methods into 3 distinct mechanisms: 1) *weight-based FT*, which ensembles the weights from both pre-trained and fine-tuned models, 2) *representation-based FT*, which regularizes the proximity between pre-trained and fine-tuned latent data representations, and 3) *partial FT*, which optimizes only a subset of the pre-trained model’s weights while keeping the rest frozen. To derive generalizable insights into how different fine-tuning mechanisms interact with pre-training strategies and downstream task types, we evaluate six diverse pre-trained models, spanning self-supervised and supervised learning, with pure graph-based, graph transformer based and multi-modal models in varying scales, then evaluate on a broad set of molecular property prediction tasks, including 8 classification and 4 regression tasks. To simulate the challenges encountered during the fine-tuning stages of MGFMs, we further consider the few-shot and out-of-distribution settings. Drawing from the broad range of pre-trained models and downstream tasks, we indeed find that the choice of best fine-tuning mechanism is highly determined by the *pre-training objective* and the *downstream task type*. We summarize high-level insights as follows, with further detailed results presented in Sec. 4. The bold text within brackets indicates the corresponding support in the experiment sections for clear cross-referencing:

- **Impact from Supervised vs. Self-supervised pre-trained models:** Supervised pre-training learns domain-specific information with task supervision, while self-supervised pre-training captures general-purpose knowledge through training on generic synthetic tasks. We observe that, in few shot fine-tuning, supervised pre-training generally yields better fine-tuning performance than self-supervised pre-training even when the pre-training tasks do not align well with the fine-tuning tasks. In contrast, for non-few-shot settings, supervised pre-training performs better only when the supervised pre-training tasks closely align with the downstream tasks [Q2].
- **Impact from Classification vs. Regression tasks:** Regression tasks need more precise numerical labels and finer molecule modeling. Therefore, MGFMs face less risk of overfitting in regression tasks compared to classification tasks, particularly in the few-shot setting [Q1].
- **Correspondence with different fine-tuning methods:** For self-supervised pre-trained models, *weight-based fine-tuning* often results in better performance by effectively integrating general knowledge from pre-training with task-specific knowledge from fine-tuning [Finding 1]. On the other hand, *partial fine-tuning* typically leads to underfitted molecular representations in few-shot fine-tuning, particularly for regression tasks [Finding 2]. For supervised pre-trained models, *representation-based fine-tuning* performs well due to the preservation of domain-relevant pre-trained representations [Finding 3].

Based on the findings, we argue that the first step in selecting or designing an effective fine-tuning strategy is to consider the pre-training strategies. Then after finding the suitable fine-tuning mechanisms, we need to take the type of downstream tasks into account. For instance, weight-based fine-tuning methods generally work the best under self-supervised pre-trained model, while simple post-hoc weight interpolation between pre-trained and fine-tuned model weights (WiSE-FT) performs well for classification tasks but struggles with regression tasks. In contrast, a more complex weight ensemble approach ( $L^2$ -SP) achieves better performance in regression tasks, though it comes with the cost of increased tuning complexity. Therefore, inspired by the rule, we propose a **new method**, **DWiSE-FT** that achieves strong performance for both regression and classification tasks as a weight-based solution for self-supervised pre-trained model. DWiSE-FT combines the strengths of WiSE-FT and  $L^2$ -SP, providing strong performance for both task types while maintaining the plug-and-play ease of post-hoc interpolation. The success of DWiSE-FT showcases that our benchmark identifies valuable insights in improving fine-tuning strategies given distinct MGFMs.

## 2 Finetuning Methods for Evaluation

In this section, we briefly introduce representative methodologies used in pre-training and fine-tuning for MGFMs.

**Self-supervised Pre-training** strategies have been proven to be effective in generating transferable molecular representations for downstream tasks [44]. In a high level, they can be divided into

reconstruction methods and *contrastive* methods. The generative-based strategies adopt mask-based graph reconstruction by utilizing graph autoencoders [28, 45, 46, 47], context predictions [27, 35] and generative language model pre-training [48, 49]. On the other hand, contrastive-based methods aim for maximizing the similarity between perturbed instance pairs [50, 30, 51, 52, 53, 54, 55, 56, 57, 58]. Moreover, the advancement of language models has prompted numerous studies to employ multi-modal frameworks. These approaches harness language models to enhance molecular understanding through techniques such as cross-modal contrastive learning and alignment [59, 60, 61, 62].

In this work, we select *GraphMAE* [28] as the representative of the reconstruction-based pre-trained model, which focuses on masked feature reconstruction with scaled cosine error that enabled robust training. Regarding the contrastive pre-trained model, we choose *Mole-BERT* [52] that combines the node-level masked atom modeling to predict the masked atom tokens and the graph-level contrastive learning through triplet loss and contrastive loss. Lastly, we choose *MoleculeSTM* [60] as the representative of multi-modal molecule structure-text model that jointly learning molecules’ chemical structures and textual descriptions via a contrastive learning strategy.

**Supervised Pre-training.** Recently, to leverage more diverse datasets and tasks, researchers started exploring the ability of supervised pre-training with multi-task learning for molecular representations [63, 31, 32]. We adopt pre-trained models trained on multi-task labeled samples in a supervised manner from the *Graphium* library [32]. In addition to the GNN-based backbone, more expressive architectures like Graph Transformer [64, 65, 66] have been proposed and can be used as the pre-trained backbone with supervised labels, which we adopt *GraphGPS* [65] as a representative.

**Fine-tuning’s** overall goal is to adapt the pre-trained model to downstream applications. Specifically, given a pre-trained GNN encoder  $f_{\theta}$  with parameters  $\theta$  initialized from the pretrained parameters  $\theta_{\text{pre}}$ , fine-tuning optimizes the encoder  $f_{\theta}$  and an additional prediction head  $g_{\phi}$  with parameters  $\phi$  over downstream molecules  $\{(\mathcal{G}_i, y_i)\}_{i=1}^N$ . The vanilla version, **full-FT**, optimizes the entire model weights following:

$$\min_{\{\theta, \phi\}} \sum_{i=1}^N \mathcal{L}(g_{\phi} \circ f_{\theta}(\mathcal{G}_i), y_i), \quad (1)$$

where  $\theta$  is initialized as  $\theta_{\text{pre}}$  and  $\mathcal{L}$  denotes the loss function for prediction tasks. As discussed, there are advanced fine-tuning strategies proposed on top of the full-FT framework. As shown in Fig. 1, we group them into three categories based on their mechanisms and benchmark representative methods for each category. More FT methods that fall into each category or others will be discussed in Appendix C.

• **Partial FT** strategies only optimizes partial weights of the pre-trained model, *i.e.*, a subset of weights within  $\{\theta, \phi\}$  will be updated following the same objective as Eq. 1. *Linear Probing (LP)* only trains the additional prediction head  $g$  during the FT. *Surgical FT* [12] updates only partial layers within the encoder. For instance, we can update the weights for  $k$ -th layer of the GNN encoder as  $\min_{\{\theta\}_k, \phi\}} \sum_{i=1}^N \mathcal{L}(g_{\phi} \circ f_{\theta}(\mathcal{G}_i), y_i)$ , where  $k$  is the hyperparameter that can be tuned. *LP-FT* [20] aims to address the issue of pre-trained feature distortion during the full-FT process. It first performs the LP step to the prediction head  $g_{\phi}$  while keeping the encoder  $f_{\theta}$  with fixed pre-trained parameters  $\theta_{\text{pre}}$ , followed by applying full-FT with the tuned prediction head.

• **Weight-based FT** strategies mainly update the entire model weights through combining pre-trained model weights and fine-tuned model weights. *WiSE-FT* [19] linearly interpolates between pre-training parameters  $\theta_{\text{pre}}$  and fine-tuning parameters  $\theta_{\text{ft}}$  using a mixing coefficient  $\alpha$ , to get the interpolated GNN  $f_{\theta_{\text{int}}}$  with weights  $\theta_{\text{int}} = (1 - \alpha) \cdot \theta_{\text{pre}} + \alpha \cdot \theta_{\text{ft}}$ . We first perform full-FT to obtain the adapted encoder  $f_{\theta_{\text{ft}}}$  and classifier  $g_{\phi}$ , then apply post-hoc weight ensembling to get  $f_{\theta_{\text{int}}}$ , with final predictions given by  $g_{\phi} \circ f_{\theta_{\text{int}}}(\mathcal{G}_i)$ .  $\alpha$ , as a hyperparameter, controls the weight ensemble. *L<sup>2</sup>-SP* [14] regularizes the fine-tuning model weights  $\theta$  closer to the pre-trained weights  $\theta_{\text{pre}}$  by  $\Omega(\theta, \phi) = \frac{\delta}{2} \|\theta - \theta_{\text{pre}}\|_2^2$ . We optimize for  $\theta$  and  $\phi$  by combining the prediction loss from Eq. 1 and  $\Omega(\theta, \phi)$  with tunable trade-off coefficient  $\delta$ .

• **Representation-based FT** methods mainly regulate the latent representation space during FT. *Feature-map* [13] adds distance regularization between the latent representations of pre-trained and fine-tuned models to the Full-FT loss. The regularization is defined as  $\Omega(\theta) = \delta \sum_{i=1}^N \frac{1}{2} \|f_{\theta}(\mathcal{G}_i) - f_{\theta_{\text{pre}}}(\mathcal{G}_i)\|_2^2$ , where  $\delta$  controls the regularization strength. *BSS* [17] aims at resolving the negative transfer issue through eliminating the spectral components corresponding to small singular values that are less transferable. The regularization is done as  $\Omega(\mathbf{F}) = \delta \sum_{i=1}^k \sigma_{-i}^2$ ,

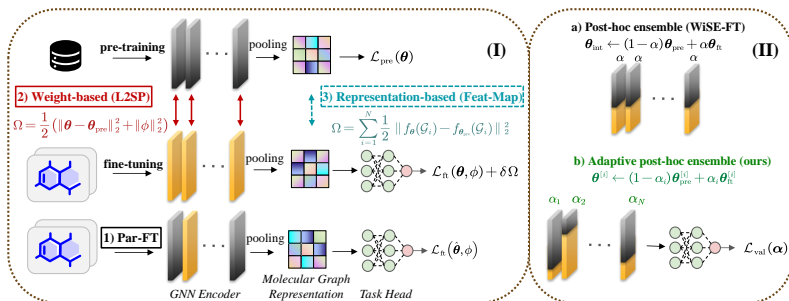


Figure 1: The overall framework of fine-tuning strategies evaluated in our benchmark, ROFT-MOL, and the proposed novel method, DWiSE-FT. **(I)** The GNN encoder is pre-trained on a large database by the objective  $\mathcal{L}_{pre}$ , and fine-tuned on the downstream dataset by  $\mathcal{L}_{ft}$  (c.f., Eq. 1). 1) Partial-FT, 2) Weight-based FT, and 3) Representation-based FT achieve robust fine-tuning by freezing partial pre-trained model weights, regularizing model weights and latent representations, respectively. **(II)** DWiSE-FT combines the strength of simple post-hoc weight interpolation with more elaborate weight ensemble, showing the improved performance while maintaining easy usage.

where  $\mathbf{F} = [f_{\theta}(\mathcal{G}_0), \dots, f_{\theta}(\mathcal{G}_b)]$  is the feature matrix of a batch of graphs and  $\sigma_{-i}$  are the  $i$ -th smallest singular values obtained from the SVD of  $\mathbf{F}$ . We can tune  $k$  and  $\delta$  to determine the number of singular values to penalize and the degree of penalty.

### 3 Experimental Settings in the Benchmark

In this section, we briefly introduce the experimental settings in this work. More detailed experimental settings can be found in Appendix F.

**Foundation Models.** For self-supervised pre-training, we adopt three open-source pre-trained checkpoints: *Mole-BERT*, *GraphMAE*, and *MoleculeSTM*. For supervised pre-training, we use models from the *Graphium* [32] library, which get pre-trained on the Toymix and Largemix datasets provided in this library. To differentiate between them, we refer to these models as *Graphium-Toy* and *Graphium-Large*. For larger graph transformer based model, we adopt the pre-trained checkpoint of *GraphGPS* [65] pre-trained on the PCQM4MV2 [67]. For details of datasets used in pre-training are in Appendix D. Furthermore, we include the traditional baseline XGBoost [68] for Fewshot scenarios to better compare with the foundation model in Appendix G.2.

**Downstream Datasets.** We use 8 classification and 4 regression datasets for downstream task evaluation. Detailed statistics and references for these tasks are in Appendix E.

**† Classification.** The BBBP dataset measures if a molecule will penetrate blood-brain barrier. The Tox21, ToxCast, and ClinTox datasets are related to toxicity qualitative measurements. The Sider dataset stores qualitative results of different types of adverse drug reactions. The MUV dataset is specifically designed for validation of virtual screening techniques. The HIV dataset provides qualitative activity results of the molecular ability to inhibit HIV replication. The BACE dataset contains qualitative binding results for a set of inhibitors of human  $\beta$ -secretase 1 (BACE-1).

**† Regression.** Esol is a dataset which measures aqueous solubility of molecules. The Lipo dataset measures the octanol-water partition coefficient. Cep is a subset of the Havard Clean Energy Project (CEP), which estimates the organic photovoltaic efficiency. Malaria measures the drug efficacy against the parasite that causes malaria.

**Dataset Splits.** For each downstream dataset, we experiment with *random*, *scaffold*, and *size* splits to create the Train/Val/Test subsets. Specifically, the random splitting shuffles the data, maintaining the Train/Val/Test sets as in-distribution (ID). The other two splitting methods simulate out-of-distribution (OOD) challenges in real-world applications. For scaffold splitting, we follow prior works [69], ensuring structural differences in molecular scaffolds across splits. Size splitting, following Zou et al. [70], arranges molecules in ascending order by size, evaluating model generalization across different molecule sizes.

**Size of fine-tuning samples.** In practice, molecular property prediction tasks can have very limited experimentally-validated data, e.g., with less than 100 samples [41]. Thus, we consider both *Non-Fewshot* and *Fewshot* settings to better simulate the label scarcity issue. In the Non-Fewshot setting,

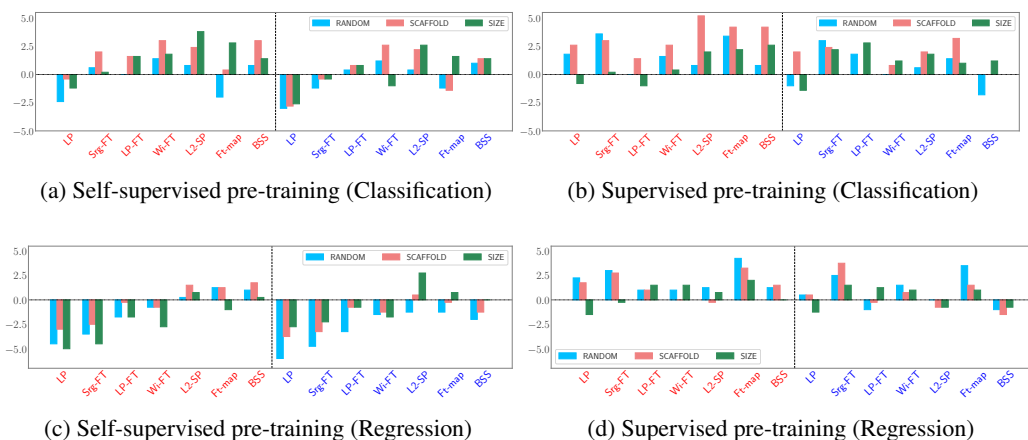


Figure 2: Average Rank improvements over Full-fine-tuning for 7 robust fine-tuning methods in self-supervised and supervised pre-training across 8 *classification* (a, b) datasets and across 4 *regression* (c, d) datasets. Each subfigure presents **few-shot-50** (left of the dashed line) and **few-shot-100** (right of the dashed line) settings, with **random**, **scaffold**, and **size** splits.

we use all available samples from the splitted train set. In the Fewshot settings, we sample subsets of 50, 100, and 500 molecules from the Train set for fine-tuning, while keeping the Val/Test sets unchanged to ensure a fair comparison. Note that we exclude MUV, Tox21, and ToxCast datasets for the Fewshot settings, as we cannot *randomly* select training samples while ensuring that all tasks have a specified number of labels simultaneously, due to the severe label scarcity issues in these datasets.

**Evaluation Metrics.** We use AUC to evaluate the performance for classification datasets and RMSE for regression datasets. We report the model performance over 5 random seeds and the test performance are reported based on the best validation performance. The AVG, AVG-F, AVG-R denote the average metrics, average metrics without max and min values, and average rank over all the datasets for each evaluated method, respectively.

Table 1: A summary of evaluated pre-trained models and their corresponding result tables for reference. “CLF” and “RGS” represent classification and regression tasks, respectively, while “NON” and “FEW” denote Non-Fewshot and Fewshot settings.

Objectives	Models	Reference Tables of Experimental Results			
		CLF-NON	CLF-FEW	RGS-NON	RGS-FEW
Self-Supervised	Mole-BERT	2	6	3	7
	GraphMAE	12	14	13	15
	MoleculeSTM	8	10	9	11
Supervised	Graphium-Toy	2	6	3	7
	Graphium-Large	8	10	9	11
	GraphGPS	12	14	13	15

## 4 Results and Analysis

This section mainly analyzes the experimental results from Mole-BERT and Graphium-Toy models as representatives of self-supervised and supervised pre-training. Table 1 is a summary of all pre-trained models we test on and their corresponding result tables for reference. Since we observe similar trends from pre-trained models of the same category, we will refer to them in our result analysis and compare over different pre-trained models in Sec. 4.3. Due to limited space, more findings with different fine-tuning methods and pre-trained models comparison can be found in Appendix G.

### 4.1 Self-supervised Pre-trained Models

**Q1: How does self-supervised pre-training influence downstream prediction tasks?**

**(1a) Regression tasks require more task-specific knowledge from downstream fine-tuning compared to classification tasks.**

When checking the few-shot results in Fig. 2a and 2c, full fine-tuning ranks the highest for regression tasks but only achieves mid-tier performance for classification tasks. This disparity likely arises from

Table 2: Robust fine-tuning performance on 8 **Classification** datasets (AUC metrics) in the **Non-Fewshot** setting, evaluated across 3 dataset splits (RANDOM, SCAFFOLD, SIZE), over **MOLE-BERT** and **GRAPHIUM-TOY** models. AVG, AVG-F, AVG-R denote the average AUC, average AUC without max and min values, and average rank over all the datasets for each method, respectively. Standard deviations across five replicates are shown. We **bold** and underline the best and second-best performances in each scenario.

SPLIT	METHODS	CLINTOX	BBBP	BACE	HIV	MUV	SIDER	TOX21	TOXCAST	AVG	AVG-F	AVG-R
SELF-SUPERVISED PRE-TRAINING (MOLE-BERT)												
SCAFFOLD	FULL-FT	77.70 ± 1.50	67.93 ± 3.85	80.12 ± 1.07	77.00 ± 0.80	80.50 ± 0.81	63.47 ± 0.77	78.31 ± 0.28	65.18 ± 0.35	73.78	74.37	3.75
	LP	66.49 ± 0.46	65.42 ± 0.26	78.70 ± 0.27	77.15 ± 0.12	79.27 ± 0.48	62.01 ± 0.60	78.12 ± 0.15	64.75 ± 0.17	71.49	71.77	6.12
	SURGICAL-FT	68.19 ± 1.58	67.70 ± 0.54	84.24 ± 0.37	76.65 ± 0.46	81.60 ± 1.02	64.61 ± 0.31	78.34 ± 0.10	65.21 ± 0.28	73.32	72.95	3.62
	LP-FT	70.35 ± 0.99	68.30 ± 0.65	81.90 ± 0.70	76.69 ± 0.40	77.65 ± 1.15	63.38 ± 0.67	77.60 ± 0.19	65.32 ± 0.24	72.65	72.65	4.88
	WISE-FT	73.59 ± 3.74	66.52 ± 3.29	82.73 ± 0.87	77.21 ± 0.69	81.92 ± 0.94	63.62 ± 0.62	78.05 ± 0.28	65.41 ± 0.25	73.63	73.78	3.38
	L <sup>2</sup> -SP	73.95 ± 1.86	67.86 ± 1.68	81.47 ± 0.80	76.63 ± 0.56	77.21 ± 0.72	65.27 ± 0.45	78.66 ± 0.17	63.55 ± 0.16	73.07	73.26	4.50
SIZE	FEATURE-MAP	70.65 ± 0.76	65.41 ± 2.37	73.44 ± 0.23	76.71 ± 0.26	80.03 ± 0.47	64.35 ± 0.17	76.61 ± 0.39	65.77 ± 0.15	71.62	71.43	5.25
	BSS	76.07 ± 3.23	67.47 ± 3.80	80.98 ± 1.27	77.12 ± 0.86	77.35 ± 1.76	63.88 ± 0.80	78.19 ± 0.40	65.00 ± 0.27	73.26	73.53	4.50
	FULL-FT	72.78 ± 1.74	87.37 ± 0.82	66.00 ± 1.99	79.85 ± 0.64	77.02 ± 2.15	52.46 ± 0.29	75.74 ± 0.48	63.13 ± 0.32	71.79	72.42	4.88
	LP	76.07 ± 0.32	82.73 ± 0.76	47.18 ± 0.45	78.16 ± 0.24	78.52 ± 1.60	51.25 ± 0.22	74.92 ± 0.22	63.33 ± 0.20	69.02	70.37	6.00
	SURGICAL-FT	73.55 ± 0.81	88.82 ± 0.53	66.43 ± 0.88	79.30 ± 0.87	80.52 ± 1.47	51.87 ± 0.23	76.32 ± 0.16	64.51 ± 0.20	72.66	73.44	3.50
	LP-FT	73.32 ± 0.93	83.42 ± 1.67	64.84 ± 1.38	79.10 ± 1.14	79.38 ± 1.86	52.82 ± 0.32	76.40 ± 0.28	63.37 ± 0.29	71.83	73.07	3.88
RANDOM	WISE-FT	73.45 ± 1.08	87.79 ± 1.53	66.58 ± 1.11	79.89 ± 1.75	78.41 ± 1.88	52.46 ± 0.49	76.46 ± 0.46	63.53 ± 0.65	72.32	73.05	3.00
	L <sup>2</sup> -SP	73.97 ± 0.88	87.15 ± 0.68	64.58 ± 1.93	80.05 ± 0.53	74.83 ± 1.06	52.37 ± 0.22	75.84 ± 0.28	60.63 ± 0.36	71.18	71.65	5.12
	FEATURE-MAP	74.61 ± 0.53	85.42 ± 0.31	51.23 ± 0.46	76.39 ± 0.91	75.20 ± 2.27	51.96 ± 0.26	76.81 ± 0.25	63.42 ± 0.76	69.38	69.73	5.00
	BSS	73.99 ± 0.77	86.84 ± 1.00	66.97 ± 1.58	79.64 ± 1.44	73.42 ± 2.60	53.50 ± 0.66	75.69 ± 0.26	62.41 ± 0.69	71.56	72.02	4.62
SUPERVISED PRE-TRAINING (GRAPHIUM-TOY)												
SCAFFOLD	FULL-FT	81.27 ± 3.88	69.17 ± 1.32	79.75 ± 1.07	76.42 ± 0.72	76.84 ± 1.80	63.63 ± 0.06	78.12 ± 0.46	66.37 ± 0.26	73.95	74.45	3.75
	LP	80.48 ± 0.00	66.90 ± 0.00	80.44 ± 0.00	75.83 ± 0.00	73.35 ± 0.00	62.03 ± 0.00	79.02 ± 0.00	66.09 ± 0.00	73.02	73.61	5.12
	SURGICAL-FT	86.17 ± 0.00	73.71 ± 0.00	84.16 ± 0.00	77.47 ± 0.00	78.87 ± 0.00	64.02 ± 0.00	78.23 ± 0.00	67.34 ± 0.00	76.25	76.63	1.38
	LP-FT	83.67 ± 3.53	69.98 ± 0.83	79.28 ± 0.32	76.17 ± 2.01	77.82 ± 1.15	61.20 ± 0.00	76.94 ± 0.00	66.28 ± 0.00	73.92	74.41	4.62
	WISE-FT	85.40 ± 1.61	71.89 ± 1.79	78.13 ± 0.92	76.69 ± 1.76	74.37 ± 1.79	63.58 ± 0.00	77.98 ± 0.33	66.48 ± 0.43	74.31	74.26	3.62
	L <sup>2</sup> -SP	76.83 ± 8.87	67.35 ± 0.82	78.17 ± 0.02	73.69 ± 0.03	62.35 ± 0.15	62.21 ± 0.45	76.27 ± 0.32	62.75 ± 0.88	69.95	69.87	6.62
SIZE	FEATURE-MAP	90.13 ± 2.12	70.99 ± 0.27	83.17 ± 0.49	73.61 ± 0.03	78.74 ± 0.76	62.12 ± 0.02	79.99 ± 0.12	65.03 ± 0.08	75.47	75.25	3.50
	BSS	79.99 ± 5.89	67.10 ± 0.93	78.12 ± 2.32	72.50 ± 0.51	61.20 ± 0.08	61.13 ± 0.95	76.69 ± 0.64	65.45 ± 0.89	70.27	70.18	7.38
	FULL-FT	85.96 ± 4.28	87.62 ± 0.90	67.41 ± 2.44	81.47 ± 1.94	72.03 ± 2.55	54.72 ± 0.01	69.71 ± 0.37	61.31 ± 0.37	72.53	72.98	3.88
	LP	81.84 ± 0.02	78.09 ± 0.00	58.08 ± 0.01	77.48 ± 0.00	69.46 ± 0.00	53.59 ± 0.00	73.65 ± 0.00	61.25 ± 0.00	69.18	69.67	5.38
	SURGICAL-FT	86.59 ± 0.01	89.07 ± 0.00	70.94 ± 0.01	82.50 ± 0.00	74.47 ± 0.00	56.24 ± 0.00	72.30 ± 0.00	62.74 ± 0.00	74.36	74.92	1.62
	LP-FT	86.78 ± 2.69	88.02 ± 1.50	63.72 ± 1.85	82.57 ± 0.46	73.51 ± 1.77	52.40 ± 0.00	68.23 ± 0.87	60.85 ± 0.00	72.01	72.61	4.00
RANDOM	WISE-FT	82.44 ± 3.02	87.76 ± 0.5	72.89 ± 0.66	81.37 ± 1.07	73.67 ± 3.44	55.87 ± 0.01	68.85 ± 0.84	60.61 ± 0.53	72.93	73.31	3.62
	L <sup>2</sup> -SP	71.03 ± 3.67	81.32 ± 1.51	68.82 ± 0.06	70.66 ± 0.00	64.69 ± 0.32	52.08 ± 0.84	70.91 ± 0.34	56.50 ± 0.01	67.00	67.10	6.88
	FEATURE-MAP	82.48 ± 3.25	87.70 ± 0.64	69.56 ± 0.20	67.23 ± 1.93	71.49 ± 0.13	54.43 ± 0.03	74.12 ± 0.09	58.73 ± 0.04	70.72	70.60	4.38
	BSS	72.42 ± 0.03	82.92 ± 1.60	62.76 ± 4.23	72.81 ± 0.66	65.79 ± 5.31	52.89 ± 1.12	71.91 ± 0.44	57.79 ± 1.80	67.41	67.25	6.25

Table 3: Robust fine-tuning performance on 4 **Regression** datasets (RMSE metrics) in the **Non-Fewshot** setting, evaluated across 3 dataset splits (RANDOM, SCAFFOLD, SIZE), over **MOLE-BERT** and **GRAPHIUM-TOY** models. AVG-R, AVG-R\* denote the average rank and the rank based on the average normalized performance over all the datasets for each method, respectively. Standard deviations across five replicates are shown. We **bold** and underline the best and second-best performances in each scenario.

SPLIT	METHODS	SELF-SUPERVISED PRE-TRAINING (MOLE-BERT)					SUPERVISED PRE-TRAINING (GRAPHIUM-TOY)				
		ESOL	LIPO	MALARIA	CYP	AVG-R	ESOL	LIPO	MALARIA	CYP	AVG-R*
SCAFFOLD	FULL-FT	1.126 ± 0.014	<u>0.728 ± 0.011</u>	1.152 ± 0.015	<b>1.377 ± 0.015</b>	3.75	3	0.911 ± 0.041	0.709 ± 0.009	1.110 ± 0.009	1.419 ± 0.014
	LP	1.614 ± 0.010	0.870 ± 0.003	1.110 ± 0.002	2.006 ± 0.002	7.00	8	0.973 ± 0.000	0.881 ± 0.000	1.105 ± 0.000	1.826 ± 0.000
	SURGICAL-FT	1.166 ± 0.017	0.783 ± 0.003	1.120 ± 0.014	1.601 ± 0.006	5.25	6	<u>0.892 ± 0.000</u>	0.709 ± 0.000	1.105 ± 0.000	1.419 ± 0.000
	LP-FT	<b>1.070 ± 0.021</b>	0.730 ± 0.002	1.144 ± 0.022	1.397 ± 0.013	3.50	4	0.922 ± 0.004	0.735 ± 0.019	<u>1.080 ± 0.005</u>	<b>1.368 ± 0.037</b>
	WISE-FT	1.264 ± 0.055	0.768 ± 0.010	<b>1.072 ± 0.001</b>	1.470 ± 0.029	4.00	2	<b>0.888 ± 0.014</b>	0.708 ± 0.008	1.128 ± 0.021	1.490 ± 0.024
	L <sup>2</sup> -SP	1.099 ± 0.030	0.742 ± 0.008	1.101 ± 0.001	1.631 ± 0.006	3.75	3	0.948 ± 0.022	0.729 ± 0.015	1.141 ± 0.015	1.606 ± 0.013
SIZE	FEATURE-MAP	1.403 ± 0.012	0.842 ± 0.004	1.083 ± 0.002	1.787 ± 0.003	5.75	7	0.895 ± 0.016	<b>0.688 ± 0.018</b>	<b>1.074 ± 0.000</b>	1.472 ± 0.010
	BSS	1.110 ± 0.022	<b>0.726 ± 0.004</b>	1.125 ± 0.018	1.385 ± 0.018	3.00	1	0.896 ± 0.018	0.718 ± 0.018	1.130 ± 0.005	1.408 ± 0.039
	FULL-FT	1.419 ± 0.044	0.745 ± 0.008	0.896 ± 0.007	<b>1.893 ± 0.035</b>	3.25	3	1.070 ± 0.082	0.719 ± 0.010	0.886 ± 0.007	1.906 ± 0.006
	LP	2.073 ± 0.012	0.912 ± 0.004	0.921 ± 0.008	2.381 ± 0.006	8.00	8	1.115 ± 0.000	0.829 ± 0.000	0.907 ± 0.000	2.246 ± 0.000
	SURGICAL-FT	1.685 ± 0.060	0.775 ± 0.007	0.890 ± 0.005	2.145 ± 0.022	5.00	6	<b>0.993 ± 0.000</b>	0.719 ± 0.000	<b>0.860 ± 0.000</b>	1.906 ± 0.000
	LP-FT	1.440 ± 0.081	0.735 ± 0.013	0.893 ± 0.007	1.905 ± 0.016	3.50	2	1.038 ± 0.038	0.694 ± 0.012	0.883 ± 0.005	1.913 ± 0.031
RANDOM	WISE-FT	1.814 ± 0.092	0.831 ± 0.007	<b>0.873 ± 0.005</b>	1.951 ± 0.024	4.50	5	1.100 ± 0.005	<b>0.691 ± 0.015</b>	0.894 ± 0.007	1.943 ± 0.039
	L <sup>2</sup> -SP	1.438 ± 0.046	0.799 ± 0.002	0.888 ± 0.005	2.101 ± 0.016	4.00	4	1.053 ± 0.026	0.720 ± 0.015	0.904 ± 0.002	2.122 ± 0.018
	FEATURE-MAP	1.656 ± 0.025	0.880 ± 0.011	0.893 ± 0.002	2.252 ± 0.008	6.25	7	<b>0.993 ± 0.034</b>	0.724 ± 0.009	0.884 ± 0.001	1.970 ± 0.013
	BSS	<b>1.375 ± 0.019</b>	<b>0.731 ± 0.007</b>	0.887 ± 0.010	1.900 ± 0.016	1.50	1	1.043 ± 0.022	0.703 ± 0.016	0.905 ± 0.005	<b>1.890 ± 0.071</b>

the distinct nature of these tasks. Classification tasks typically require coarser-grained features, as exemplified by the Tox21 dataset. In this case, determining toxicity may largely rely on recognizing certain functional groups, such as toxicophores or structural alerts [71]. In contrast, regression tasks demand finer-grained features. For example, predicting precise solubility involves factors such as partial charge distribution, conformational flexibility, and hydrogen bond patterns, among others [72]. Consequently, models fine-tuned for regression tasks must acquire more downstream knowledge during the fine-tuning process and are generally less prone to overfitting compared to those used for classification tasks.

### (1b) Molecular representations learned from self-supervised pre-training are not informative enough for downstream tasks.

As shown in Tables 2 and 3, LP is consistently the worst performing method for self-supervised pre-trained models across all data splits, even under the few-shot fine-tuning in Fig. 2a and 2c. Furthermore, this behavior is widely observed across all tested self-supervised models as GraphMAE and MoleculeSTM, which contrasts the observations in CV where LP demonstrates robust OOD performance by preserving high quality and generalizable features from pre-trained embeddings [19, 20]. We attribute this to the misalignment between general-purpose representations produced by self-supervised pre-training and the features required by the specific molecular tasks. Consequently,

relying solely on tuning the classifier  $g_\phi$  is insufficient to extract meaningful predictions from these non-informative representations.

Below, we summarize insightful findings from the performance of different fine-tuning strategies.

• **Finding 1. Under few-shot fine-tuning, weight-based fine-tuning strategies stand out with WiSE-FT for classification tasks and  $L^2$ -SP for regression tasks.**

Among various fine-tuning methods, weight-based approaches consistently outperform others across a wide range of experiments, regardless of the few-shot sample sizes (Fig. 2a and 2c). Self-supervised models are known to capture general-purpose knowledge for substructure discovery[39]. During fine-tuning, combining pre-trained and fine-tuned weights proves effective in extracting molecular patterns relevant to downstream tasks. Notably, WiSE-FT demonstrates superior performance on classification datasets, whereas  $L^2$ -SP excels in regression tasks. This finding is also supported by MoleculeSTM in table 11 where  $L^2$ -SP remains as top method under all few-shot regression tasks and WiSE-FT excels under Fewshot-50 classification. Essentially, WiSE-FT applies a straightforward post-hoc linear interpolation between pre-trained and fine-tuned models, governed by a single coefficient. In contrast,  $L^2$ -SP implicitly determines the weight combination through the training loss [15, 14], aligning with statement (1a) that regression tasks typically demand more nuanced modeling.

• **Finding 2. Partial FT results in underfitted molecular representations under Fewshot settings, which is more severe for regression tasks compared to classification.**

For the non-few-shot fine-tuning (Tables 2 and 3), surgical FT and LP-FT improve over full FT in both classification and regression tasks. However, in few-shot fine-tuning, both methods rank as the worst methods. This is likely because partial fine-tuning underfits and bias towards the limited samples. This issue is more pronounced in regression tasks.

## 4.2 Supervised Pre-trained Models

*Q2: How does supervised pre-training influence downstream tasks?*

We first discuss the **task similarity** between the datasets used in the pre-training and downstream fine-tuning process. As introduced in Appendix. D, the ToyMix dataset used for supervised pre-training contains QM9, Tox21 and Zinc12K. The predictions from QM9 are not directly related to our downstream tasks, but may involve indirect correlations, as the quantum chemical properties provided by QM9 are highly valuable for characterizing molecular features. **Tox21** is an overlapping dataset that also exists as one of the downstream datasets. Its tasks in predicting qualitative toxicity measurements are *highly related* to the downstream **ClinTox** and **ToxCast** datasets, and also *correlate* to the **Sider** dataset which contains evaluation in drug side effects. Lastly, Zinc12K, which is to predict the constrained solubility, is relevant to the **Esol** and **Lipo** datasets that involve solubility predictions. Other downstream tasks *do not share* the same tasks with pre-training *directly*. Then we observe the following rules.

**(2a) Under few-shot fine-tuning, supervised pre-training models generally yield higher fine-tuning performance compared to self-supervised pre-training, regardless of the pre-training and fine-tuning task correlations.**

Supervised pre-training brings more benefits to downstream tasks than self-supervised pre-training in few-shot situations when checking Tables 6 and 7. Besides, the benefits are less relevant to the task similarity in contrast to the non-few-shot cases. For example, the improvements are also observed in HIV and Cep datasets even their tasks do not share with pre-training tasks directly. This implies that learned domain-specific knowledge still offer better insights than generic knowledge when fine-tuning guidance is minimal.

**(2b) Under non-few-shot fine-tuning, fine-tuning performance given supervised pre-training outperforms self-supervised pre-training when its objectives closely align with downstream tasks, while task misalignment may harm performance.**

From Tables 2 and 3, we observe consistent fine-tuning performance improvements over self-supervised pre-training on highly task-correlated downstream datasets including ClinTox, Esol, Lipo and Tox21. Even when pre-training involves regression tasks and downstream tasks are classification, performance gains occur if the physical meanings align. For datasets that do not directly share tasks with pre-training, we observe mixed performance on Sider, Malaria, and Cep datasets, and even worse performance on HIV and MUV datasets. This observation contrasts to few-shot cases in (2a),



Table 4: DWiSE-FT performance on 4 **Regression** datasets (RMSE metrics) in the **Fewshot** setting with 50, 100 samples, evaluated across 3 dataset splits (RANDOM, SCAFFOLD, SIZE) given **MOLE-BERT** model. AVG-R denote the average rank. Standard deviations across five replicates are shown. We **bold** and underline the best and second-best performances in each scenario.

SPLIT	METHODS	FEWSHOT 50					FEWSHOT 100				
		ESOL	LIPO	MALARIA	CEP	AVG	ESOL	LIPO	MALARIA	CEP	AVG
RANDOM	WiSE-FT	1.384 ± 0.047	1.212 ± 0.020	1.276 ± 0.007	2.410 ± 0.051	3.75	1.189 ± 0.030	1.142 ± 0.025	<b>1.256 ± 0.006</b>	2.211 ± 0.028	3.00
	$L^2$ -SP	1.272 ± 0.029	1.196 ± 0.019	1.277 ± 0.006	2.280 ± 0.031	3.00	1.161 ± 0.016	1.149 ± 0.007	1.260 ± 0.004	<u>2.131 ± 0.014</u>	3.25
	Top	<b>1.329 ± 0.021</b>	<b>1.164 ± 0.010</b>	<b>1.271 ± 0.007</b>	<b>2.275 ± 0.022</b>	1.25	<b>1.120 ± 0.038</b>	<b>1.139 ± 0.017</b>	<b>1.256 ± 0.006</b>	<b>2.131 ± 0.014</b>	1.50
	DWiSE-FT	1.378 ± 0.055	1.189 ± 0.020	1.273 ± 0.009	<b>2.222 ± 0.059</b>	2.00	1.132 ± 0.025	<b>1.138 ± 0.028</b>	<b>1.256 ± 0.004</b>	<b>2.129 ± 0.020</b>	1.25
SCAFFOLD	WiSE-FT	1.842 ± 0.056	1.177 ± 0.009	<u>1.162 ± 0.004</u>	2.454 ± 0.043	3.50	1.544 ± 0.063	1.041 ± 0.017	<u>1.151 ± 0.007</u>	2.301 ± 0.042	3.50
	$L^2$ -SP	1.699 ± 0.049	1.086 ± 0.009	<u>1.162 ± 0.002</u>	2.331 ± 0.024	2.50	1.473 ± 0.009	0.961 ± 0.003	1.153 ± 0.002	2.201 ± 0.038	2.50
	Top	<u>1.680 ± 0.042</u>	<b>1.036 ± 0.007</b>	<b>1.159 ± 0.000</b>	<b>2.292 ± 0.026</b>	1.25	<b>1.436 ± 0.054</b>	<b>0.937 ± 0.008</b>	<b>1.149 ± 0.003</b>	<b>2.187 ± 0.034</b>	1.25
	DWiSE-FT	<b>1.616 ± 0.047</b>	1.110 ± 0.013	1.173 ± 0.005	2.306 ± 0.030	2.50	1.485 ± 0.041	0.979 ± 0.014	1.158 ± 0.009	<b>2.149 ± 0.040</b>	2.75
SIZE	WiSE-FT	2.615 ± 0.072	1.391 ± 0.042	0.929 ± 0.004	2.762 ± 0.053	4.00	2.216 ± 0.056	1.124 ± 0.031	0.917 ± 0.004	2.543 ± 0.027	3.75
	$L^2$ -SP	2.393 ± 0.068	1.306 ± 0.037	0.915 ± 0.002	<b>2.497 ± 0.019</b>	2.50	<u>1.731 ± 0.071</u>	<b>1.025 ± 0.028</b>	<u>0.905 ± 0.002</u>	<u>2.424 ± 0.024</u>	1.75
	Top	<u>2.369 ± 0.075</u>	<u>1.297 ± 0.040</u>	<b>0.911 ± 0.002</b>	<b>2.497 ± 0.019</b>	1.50	<u>1.731 ± 0.071</u>	<b>1.025 ± 0.028</b>	<b>0.898 ± 0.003</b>	<u>2.424 ± 0.024</u>	1.50
	DWiSE-FT	<b>1.488 ± 0.101</b>	<b>1.113 ± 0.021</b>	<u>0.913 ± 0.007</u>	2.539 ± 0.023	1.75	<b>1.469 ± 0.052</b>	1.031 ± 0.022	0.920 ± 0.006	<b>2.390 ± 0.025</b>	2.25

which entails that downstream task specific knowledge can be learned given sufficient guidance on top of generic knowledge from self-supervised pre-training.

Below are some detailed findings with different fine-tuning methods given supervised pre-training.

• **Finding 3. Fine-tuning strategies that regularizes towards pre-trained molecular representations rank top, while weight-based methods are suboptimal.**

From non-few-shot (Tables 2 and 3) and few-shot fine-tuning (Figs. 2b and 2d) in both supervised models with ToyMix and LargeMix, surgical FT and Feature-map tend to be the top-ranking methods. However, best performing weight-based methods for self-supervised pre-training, only show mediocre performance here. This can also be observed in the larger-scale GraphGPS model as discussed in Appendix G.1. In addition, the other representation-based method BSS shows limited performance compared to Feature-map, which directly regularizes the distance to pre-trained representations. These observations suggest that given the task alignment between supervised pre-training and downstream fine-tuning, pre-trained representations tend to contain transferable features for downstream tasks. Consequently, controlling the degree to preserve pre-trained representations is the key to downstream fine-tuning performance.

### 4.3 Discussions over Pre-trained Models

Our extensive evaluation shows that the ranking of fine-tuning techniques remains consistent across pre-trained models within the same category, either supervised or self-supervised, regardless of model architecture, scale, or pre-training dataset. This suggests that our guidance for selecting fine-tuning methods based on the pre-training paradigm is broadly applicable and generalizable across diverse model designs. For instance, self-supervised models such as Mole-BERT and MoleculeSTM tend to benefit more from weight-based fine-tuning, while supervised models like Graphium and GraphGPS perform better with feature-based approaches.

## 5 Methodology Exploration

Based on findings in Sec. 4, we observe that weight-based fine-tuning generally performs well under self-supervised pre-training. However, the top strategy varies: WiSE-FT excels in classification tasks, while  $L^2$ -SP is more effective for regression tasks. This motivates us to further explore the connections and trade-offs between these methods to identify potential improvements. In this section, we introduce DWiSE-FT, an extension of the weight ensemble method unifying the strengths from WiSE-FT and  $L^2$ -SP. DWiSE-FT demonstrates top-ranking results through efficient post-processing that better suits the practical fine-tuning needs.

### 5.1 Motivation

As introduced in Sec. 2, WiSE-FT adopts the post-hoc linear interpolation between the pre-trained and fine-tuned model weights as  $(1 - \alpha) \cdot \theta_{\text{pre}} + \alpha \cdot \theta_{\text{ft}}$ . Although  $L^2$ -SP does not explicitly have weight interpolation in the form, the optimal weight  $\tilde{\theta}_{\text{ft}}$  from the weight-regularized loss  $\tilde{\mathcal{L}}(\theta)$  is indeed the linear interpolation of the optimal model from full FT  $\theta_{\text{ft}}^*$  and the pre-trained model  $\theta_{\text{pre}}$ .

**Proposition 1.** Given  $\tilde{\mathcal{L}}(\theta) = \mathcal{L}(\theta) + \frac{\delta}{2} \|\theta - \theta_{\text{pre}}\|_2^2$ , we define the optimal weights as  $\tilde{\theta}_{\text{ft}} = \text{argmin}_{\theta} \tilde{\mathcal{L}}(\theta)$  and  $\theta_{\text{ft}}^* = \text{argmin}_{\theta} \mathcal{L}(\theta)$ .

$$\mathbf{Q}^T \tilde{\theta}_{\text{ft}} = (\mathbf{I} + \delta \mathbf{H})^{-1} \mathbf{I} \mathbf{Q}^T \theta_{\text{ft}}^* + \delta (\mathbf{I} + \delta \mathbf{H})^{-1} \mathbf{Q}^T \theta_{\text{pre}} \quad (2)$$

where  $\mathbf{H}$  is the hessian matrix of  $\mathcal{L}$  evaluated at  $\theta_{\text{ft}}^*$  and  $\mathbf{H} = \mathbf{Q} \mathbf{A} \mathbf{Q}^T$ .



316 Namely,  $L^2$ -SP can be seen as a more tailored weight ensemble method, employing variable mixing  
 317 coefficients for different weights. This approach balances the influence of the prediction loss and the  
 318 degree of weight regularization, unlike the fixed interpolation controlled by  $\alpha$  across all weights in  
 319 WiSE-FT. By accounting for subtle differences in loss values,  $L^2$ -SP is better suited for regression  
 320 tasks, which are more sensitive to numerical variations.

321 While  $L^2$ -SP excels on regression datasets, its regularization coefficient is less interpretable and  
 322 necessitates retraining when experimenting with different values. In contrast, WiSE-FT offers a  
 323 simpler and more flexible approach, performing post-hoc interpolation without additional training  
 324 once the model is fine-tuned once. Furthermore, the mixing coefficient  $\alpha$  is both easy to adjust and  
 325 straightforward to interpret. Therefore, our goal is to find a method that benefits from both WiSE-FT  
 326 and  $L^2$ -SP to accommodate regression and classification tasks at the same time.

## 327 5.2 Algorithm

328 We propose DWiSE-FT that shares the framework of using the  $\alpha$  to control the weight ensemble  
 329 between the pre-trained model and fine-tuned model. The key idea, inspired by Eq. 4 is to enable  
 330 different  $\alpha$  values when ensembling the weights for different encoder layers as shown in Fig. 1. Given  
 331 the pre-trained model with parameters  $\theta_{\text{pre}}$  and model after full fine-tuning with parameters  $\theta_{\text{ft}}$ , The  
 332 interpolated model has weights  $\theta^{[i]}$  with mixing coefficient  $\alpha_i$  for the  $i$ -th layer as:

$$\theta^{[i]} = (1 - \alpha_i) \cdot \theta_{\text{pre}}^{[i]} + \alpha_i \cdot \theta_{\text{ft}}^{[i]} \quad (3)$$

333 This approach naturally incorporates the characteristics of  $L^2$ -SP and even surgical FT: The weight  
 334 ensemble in DWiSE-FT offers the flexibility through varying mixing layer-wise coefficients between  
 335 the pre-trained and fine-tuned models, addressing the limitations of WiSE-FT. Additionally, we enable  
 336 the selection of  $\alpha$  through optimization via validation loss gradient inspired by the Gradient-based  
 337 Neural Architecture Search (NAS) [73].

## 339 5.3 Experiment results

340 Regarding the classification datasets, DWiSE-FT should have the performance at least as good as  
 341 WiSE-FT since WiSE-FT is the special case of DWiSE-FT with one fixed mixing coefficient. We  
 342 evaluate DWiSE-FT to see how it improves upon WiSE-FT and matches the superior performance  
 343 of  $L^2$ -SP for regression tasks under few-shot fine-tuning. Please note that, due to space constraints,  
 344 we only present the experiments for few-shot fine-tuning with 50 and 100 samples in the main text.  
 345 The complete table is available in Appendix E, Table 17. In Table 4, we compare DWiSE-FT’s  
 346 performance against WiSE-FT,  $L^2$ -SP, and the best-performing method in each setting. Specifically,  
 347 we find that DWiSE-FT consistently outperforms WiSE-FT. Furthermore, DWiSE-FT often surpasses  
 348  $L^2$ -SP or at least maintains comparable results in most scenarios. Additionally, in some cases,  
 349 DWiSE-FT even exceeds the performance of the best-performing methods. Therefore, DWiSE-FT  
 350 can be a great candidate for fine-tuning on regression datasets in practice since it guarantees top  
 351 performance with easier usage.

## 352 6 Conclusion

353 This work benchmarks totally 8 fine-tuning methods, categorizing them into three groups, and evaluate  
 354 them across 12 downstream datasets under 36 different experimental settings covering 3 dataset  
 355 splits, 4 training sample sizes, and 6 molecular pre-trained models. The design of these settings  
 356 reflects practical demands of molecular representation fine-tuning under 1) diversified foundation  
 357 model with both supervised and self-supervised pre-training, 2) wide range of downstream tasks  
 358 in both classification and regression that has not been widely studied by previous literature and  
 359 3) scarcely labeled molecules for fine-tuning. The study analyzes what is needed when facing  
 360 classification vs. regression tasks and when given supervised vs. self-supervised pre-training. Then,  
 361 we provide insights in best performing fine-tuning methods accordingly under aforementioned  
 362 scenarios. Additionally, we propose an extended fine-tuning method DWiSE-FT, driven by our  
 363 observations, that maintains top-ranking results through a more efficient and automated design for  
 364 certain fine-tuning scenarios. This highlights the value of our benchmark in offering valuable insights  
 365 for both fine-tuning methodology design and practical guidance in molecular representation learning.

## References

- [1] Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, et al. On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*, 2021.
- [2] Ce Zhou, Qian Li, Chen Li, Jun Yu, Yixin Liu, Guangjing Wang, Kai Zhang, Cheng Ji, Qiben Yan, Lifang He, et al. A comprehensive survey on pretrained foundation models: A history from bert to chatgpt. *arXiv preprint arXiv:2302.09419*, 2023.
- [3] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021.
- [4] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4015–4026, 2023.
- [5] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 1(2):3, 2022.
- [6] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- [7] Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, et al. Sparks of artificial general intelligence: Early experiments with gpt-4. *arXiv preprint arXiv:2303.12712*, 2023.
- [8] Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, et al. A survey of large language models. *arXiv preprint arXiv:2303.18223*, 2023.
- [9] Jeremy Howard and Sebastian Ruder. Universal language model fine-tuning for text classification. *arXiv preprint arXiv:1801.06146*, 2018.
- [10] Hao Li, Pratik Chaudhari, Hao Yang, Michael Lam, Avinash Ravichandran, Rahul Bhotika, and Stefano Soatto. Rethinking the hyperparameters for fine-tuning. In *International Conference on Learning Representations*, 2020.
- [11] Simon Kornblith, Jonathon Shlens, and Quoc V Le. Do better imagenet models transfer better? In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2661–2671, 2019.
- [12] Yoonho Lee, Annie S Chen, Fahim Tajwar, Ananya Kumar, Huaxiu Yao, Percy Liang, and Chelsea Finn. Surgical fine-tuning improves adaptation to distribution shifts. In *The Eleventh International Conference on Learning Representations*, 2022.
- [13] Xingjian Li, Haoyi Xiong, Hanchao Wang, Yuxuan Rao, Liping Liu, and Jun Huan. Delta: Deep learning transfer using feature map with attention for convolutional networks. In *International Conference on Learning Representations*, 2019.
- [14] LI Xuhong, Yves Grandvalet, and Franck Davoine. Explicit inductive bias for transfer learning with convolutional networks. In *International Conference on Machine Learning*, pages 2825–2834. PMLR, 2018.
- [15] Ekdeep Singh Lubana, Puja Trivedi, Danai Koutra, and Robert Dick. How do quadratic regularizers prevent catastrophic forgetting: The role of interpolation. In *Conference on Lifelong Learning Agents*, pages 819–837. PMLR, 2022.

- 413 [16] Zirui Wang, Zihang Dai, Barnabás Póczos, and Jaime Carbonell. Characterizing and avoiding  
414 negative transfer. In *Proceedings of the IEEE/CVF conference on computer vision and pattern*  
415 *recognition*, pages 11293–11302, 2019.
- 416 [17] Xinyang Chen, Sinan Wang, Bo Fu, Mingsheng Long, and Jianmin Wang. Catastrophic  
417 forgetting meets negative transfer: Batch spectral shrinkage for safe transfer learning. *Advances*  
418 *in Neural Information Processing Systems*, 32, 2019.
- 419 [18] Zhiqiang Shen, Zechun Liu, Jie Qin, Marios Savvides, and Kwang-Ting Cheng. Partial is  
420 better than all: Revisiting fine-tuning strategy for few-shot learning. In *Proceedings of the*  
421 *AAAI conference on artificial intelligence*, volume 35, pages 9594–9602, 2021.
- 422 [19] Mitchell Wortsman, Gabriel Ilharco, Jong Wook Kim, Mike Li, Simon Kornblith, Rebecca  
423 Roelofs, Raphael Gontijo Lopes, Hannaneh Hajishirzi, Ali Farhadi, Hongseok Namkoong,  
424 et al. Robust fine-tuning of zero-shot models. In *Proceedings of the IEEE/CVF conference on*  
425 *computer vision and pattern recognition*, pages 7959–7971, 2022.
- 426 [20] Ananya Kumar, Aditi Raghunathan, Robbie Jones, Tengyu Ma, and Percy Liang. Fine-  
427 tuning can distort pretrained features and underperform out-of-distribution. In *International*  
428 *Conference on Learning Representations*, 2022.
- 429 [21] Yang Shu, Xingzhuo Guo, Jialong Wu, Ximei Wang, Jianmin Wang, and Mingsheng Long.  
430 Clipood: Generalizing clip to out-of-distributions. In *International Conference on Machine*  
431 *Learning*, pages 31716–31731. PMLR, 2023.
- 432 [22] Anders Johan Andreassen, Yasaman Bahri, Behnam Neyshabur, and Rebecca Roelofs. The  
433 evolution of out-of-distribution robustness throughout fine-tuning. *Transactions on Machine*  
434 *Learning Research*, 2021.
- 435 [23] Polina Kirichenko, Pavel Izmailov, and Andrew Gordon Wilson. Last layer re-training is  
436 sufficient for robustness to spurious correlations. In *The Eleventh International Conference on*  
437 *Learning Representations*, 2022.
- 438 [24] Tobias Golling, Lukas Heinrich, Michael Kagan, Samuel Klein, Matthew Leigh, Margarita  
439 Osadchy, and John Andrew Raine. Masked particle modeling on sets: towards self-supervised  
440 high energy physics foundation models. *Machine Learning: Science and Technology*, 5(3):  
441 035074, 2024.
- 442 [25] Henry W Leung and Jo Bovy. Towards an astronomical foundation model for stars with  
443 a transformer-based model. *Monthly Notices of the Royal Astronomical Society*, 527(1):  
444 1494–1520, 2024.
- 445 [26] Tung Nguyen, Johannes Brandstetter, Ashish Kapoor, Jayesh K Gupta, and Aditya Grover.  
446 Climax: A foundation model for weather and climate. *arXiv preprint arXiv:2301.10343*, 2023.
- 447 [27] W Hu, B Liu, J Gomes, M Zitnik, P Liang, V Pande, and J Leskovec. Strategies for pre-training  
448 graph neural networks. In *International Conference on Learning Representations (ICLR)*,  
449 2020.
- 450 [28] Zhenyu Hou, Xiao Liu, Yukuo Cen, Yuxiao Dong, Hongxia Yang, Chunjie Wang, and Jie Tang.  
451 Graphmae: Self-supervised masked graph autoencoders. In *Proceedings of the 28th ACM*  
452 *SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 594–604, 2022.
- 453 [29] Jun Xia, Chengshuai Zhao, Bozhen Hu, Zhangyang Gao, Cheng Tan, Yue Liu, Siyuan Li, and  
454 Stan Z Li. Mole-bert: Rethinking pre-training graph neural networks for molecules. In *The*  
455 *Eleventh International Conference on Learning Representations*, 2023.
- 456 [30] Susheel Suresh, Pan Li, Cong Hao, and Jennifer Neville. Adversarial graph augmentation to  
457 improve graph contrastive learning. *Advances in Neural Information Processing Systems*, 34:  
458 15920–15933, 2021.
- 459 [31] Nima Shoghi, Adeesh Kolluru, John R Kitchin, Zachary Ward Ulissi, C Lawrence Zitnick,  
460 and Brandon M Wood. From molecules to materials: Pre-training large generalizable mod-  
461 els for atomic property prediction. In *The Twelfth International Conference on Learning*  
462 *Representations*, 2023.

- [32] Dominique Beaini, Shenyang Huang, Joao Alex Cunha, Zhiyi Li, Gabriela Moisesescu-Pareja, Oleksandr Dymov, Samuel Maddrell-Mander, Callum McLean, Frederik Wenkel, Luis Müller, et al. Towards foundational models for molecular learning on large-scale multi-task datasets. In *The Twelfth International Conference on Learning Representations*, 2023.
- [33] Shuxin Zheng, Jiyan He, Chang Liu, Yu Shi, Ziheng Lu, Weitao Feng, Fusong Ju, Jiaxi Wang, Jianwei Zhu, Yaosen Min, et al. Towards predicting equilibrium distributions for molecular systems with deep learning. *arXiv preprint arXiv:2306.05445*, 2023.
- [34] Jerret Ross, Brian Belgodere, Vijil Chenthamarakshan, Inkit Padhi, Youssef Mroueh, and Payel Das. Large-scale chemical language representations capture molecular structure and properties. *Nature Machine Intelligence*, 4(12):1256–1264, 2022.
- [35] Yu Rong, Yatao Bian, Tingyang Xu, Weiyang Xie, Ying Wei, Wenbing Huang, and Junzhou Huang. Self-supervised graph transformer on large-scale molecular data. *Advances in neural information processing systems*, 33:12559–12571, 2020.
- [36] Haitao Mao, Zhikai Chen, Wenzhuo Tang, Jianan Zhao, Yao Ma, Tong Zhao, Neil Shah, Mikhail Galkin, and Jiliang Tang. Position: Graph foundation models are already here. In *Forty-first International Conference on Machine Learning*, 2024.
- [37] Dingshuo Chen, Yanqiao Zhu, Jieyu Zhang, Yuanqi Du, Zhixun Li, Qiang Liu, Shu Wu, and Liang Wang. Uncovering neural scaling laws in molecular representation learning. *Advances in Neural Information Processing Systems*, 36, 2024.
- [38] Ruoxi Sun, Hanjun Dai, and Adams Wei Yu. Does gnn pretraining help molecular representation? *Advances in Neural Information Processing Systems*, 35:12096–12109, 2022.
- [39] Hanchen Wang, Jean Kaddour, Shengchao Liu, Jian Tang, Joan Lasenby, and Qi Liu. Evaluating self-supervised learning for molecular graph embeddings. *Advances in Neural Information Processing Systems*, 36, 2024.
- [40] Mohammad Sadegh Akhondzadeh, Vijay Lingam, and Aleksandar Bojchevski. Probing graph representations. In *International Conference on Artificial Intelligence and Statistics*, pages 11630–11649. PMLR, 2023.
- [41] Kevin Tirta Wijaya, Minghao Guo, Michael Sun, Hans-Peter Seidel, Wojciech Matusik, and Vahid Babaei. Two-stage pretraining for molecular property prediction in the wild. *arXiv preprint arXiv:2411.03537*, 2024.
- [42] Zhenqin Wu, Bharath Ramsundar, Evan N Feinberg, Joseph Gomes, Caleb Geniesse, Aneesh S Pappu, Karl Leswing, and Vijay Pande. Moleculenet: a benchmark for molecular machine learning. *Chemical science*, 9(2):513–530, 2018.
- [43] Yuanyuan Hou, Shiyu Wang, Bing Bai, HC Stephen Chan, and Shuguang Yuan. Accurate physical property predictions via deep learning. *Molecules*, 27(5):1668, 2022.
- [44] Ziwen Zhao, Yuhua Li, Yixiong Zou, Ruixuan Li, and Rui Zhang. A survey on self-supervised pre-training of graph foundation models: A knowledge-based perspective. *arXiv preprint arXiv:2403.16137*, 2024.
- [45] Qiaoyu Tan, Ninghao Liu, Xiao Huang, Soo-Hyun Choi, Li Li, Rui Chen, and Xia Hu. S2gae: Self-supervised graph autoencoders are generalizable learners with graph masking. In *Proceedings of the sixteenth ACM international conference on web search and data mining*, pages 787–795, 2023.
- [46] Chun Wang, Shirui Pan, Guodong Long, Xingquan Zhu, and Jing Jiang. Mgae: Marginalized graph autoencoder for graph clustering. In *Proceedings of the 2017 ACM on Conference on Information and Knowledge Management*, pages 889–898, 2017.
- [47] Shirui Pan, Ruiqi Hu, Guodong Long, Jing Jiang, Lina Yao, and Chengqi Zhang. Adversarially regularized graph autoencoder for graph embedding. In *International Joint Conference on Artificial Intelligence 2018*, pages 2609–2615. Association for the Advancement of Artificial Intelligence (AAAI), 2018.

- [48] Ziniu Hu, Yuxiao Dong, Kuansan Wang, Kai-Wei Chang, and Yizhou Sun. Gpt-gnn: Generative pre-training of graph neural networks. In *Proceedings of the 26th ACM SIGKDD international conference on knowledge discovery & data mining*, pages 1857–1867, 2020.
- [49] Zaixi Zhang, Qi Liu, Hao Wang, Chengqiang Lu, and Chee-Kong Lee. Motif-based graph self-supervised learning for molecular property prediction. *Advances in Neural Information Processing Systems*, 34:15870–15882, 2021.
- [50] Petar Veličković, William Fedus, William L Hamilton, Pietro Liò, Yoshua Bengio, and R Devon Hjelm. Deep graph infomax. In *the International Conference on Learning Representations (ICLR)*, 2018.
- [51] Yuning You, Tianlong Chen, Yongduo Sui, Ting Chen, Zhangyang Wang, and Yang Shen. Graph contrastive learning with augmentations. *Advances in neural information processing systems*, 33:5812–5823, 2020.
- [52] Jun Xia, Chengshuai Zhao, Bozhen Hu, Zhangyang Gao, Cheng Tan, Yue Liu, Siyuan Li, and Stan Z Li. Mole-bert: Rethinking pre-training graph neural networks for molecules. 2023.
- [53] Yuyang Wang, Jianren Wang, Zhonglin Cao, and Amir Barati Farimani. Molecular contrastive learning of representations via graph neural networks. *Nature Machine Intelligence*, 4(3): 279–287, 2022.
- [54] Jinhua Zhu, Yingce Xia, Lijun Wu, Shufang Xie, Tao Qin, Wengang Zhou, Houqiang Li, and Tie-Yan Liu. Unified 2d and 3d pre-training of molecular representations. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 2626–2636, 2022.
- [55] Yuning You, Tianlong Chen, Yang Shen, and Zhangyang Wang. Graph contrastive learning automated. In *the International Conference on Machine Learning (ICML)*, pages 12121–12132. PMLR, 2021.
- [56] Jiezhong Qiu, Qibin Chen, Yuxiao Dong, Jing Zhang, Hongxia Yang, Ming Ding, Kuansan Wang, and Jie Tang. Gcc: Graph contrastive coding for graph neural network pre-training. In *Proceedings of the 26th ACM SIGKDD international conference on knowledge discovery & data mining*, pages 1150–1160, 2020.
- [57] Sihang Li, Xiang Wang, An Zhang, Yingxin Wu, Xiangnan He, and Tat-Seng Chua. Let invariant rationale discovery inspire graph contrastive learning. In *the International Conference on Machine Learning (ICML)*, pages 13052–13065. PMLR, 2022.
- [58] Minghao Xu, Hang Wang, Bingbing Ni, Hongyu Guo, and Jian Tang. Self-supervised graph-level representation learning with local and global structure. In *the International Conference on Machine Learning (ICML)*, pages 11548–11558. PMLR, 2021.
- [59] Bing Su, Dazhao Du, Zhao Yang, Yujie Zhou, Jiangmeng Li, Anyi Rao, Hao Sun, Zhiwu Lu, and Ji-Rong Wen. A molecular multimodal foundation model associating molecule graphs with natural language. *arXiv preprint arXiv:2209.05481*, 2022.
- [60] Shengchao Liu, Weili Nie, Chengpeng Wang, Jiarui Lu, Zhuoran Qiao, Ling Liu, Jian Tang, Chaowei Xiao, and Animashree Anandkumar. Multi-modal molecule structure-text model for text-based retrieval and editing. *Nature Machine Intelligence*, 5(12):1447–1457, 2023.
- [61] Philipp Seidl, Andreu Vall, Sepp Hochreiter, and Günter Klambauer. Enhancing activity prediction models in drug discovery with the ability to understand human language. In *International Conference on Machine Learning*, pages 30458–30490. PMLR, 2023.
- [62] Zhiyuan Liu, Sihang Li, Yanchen Luo, Hao Fei, Yixin Cao, Kenji Kawaguchi, Xiang Wang, and Tat-Seng Chua. Molca: Molecular graph-language modeling with cross-modal projector and uni-modal adapter. In *The Conference on Empirical Methods in Natural Language Processing*, 2023.

- [63] Johannes Gasteiger, Muhammed Shuaibi, Anuroop Sriram, Stephan Günnemann, Zachary Ward Ulissi, C Lawrence Zitnick, and Abhishek Das. Gemnet-oc: Developing graph neural networks for large and diverse molecular simulation datasets. *Transactions on Machine Learning Research*, 2022.
- [64] Chengxuan Ying, Tianle Cai, Shengjie Luo, Shuxin Zheng, Guolin Ke, Di He, Yanming Shen, and Tie-Yan Liu. Do transformers really perform badly for graph representation? *Advances in neural information processing systems*, 34:28877–28888, 2021.
- [65] Ladislav Rampásek, Michael Galkin, Vijay Prakash Dwivedi, Anh Tuan Luu, Guy Wolf, and Dominique Beaini. Recipe for a general, powerful, scalable graph transformer. *Advances in Neural Information Processing Systems*, 35:14501–14515, 2022.
- [66] Grégoire Mialon, Dexiong Chen, Margot Selosse, and Julien Mairal. Graphit: Encoding graph structure in transformers. *arXiv preprint arXiv:2106.05667*, 2021.
- [67] Weihua Hu, Matthias Fey, Hongyu Ren, Maho Nakata, Yuxiao Dong, and Jure Leskovec. Ogb-lsc: A large-scale challenge for machine learning on graphs. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2)*.
- [68] Tianqi Chen and Carlos Guestrin. Xgboost: A scalable tree boosting system. In *Proceedings of the 22nd acm sigkdd international conference on knowledge discovery and data mining*, pages 785–794, 2016.
- [69] Bharath Ramsundar, Peter Eastman, Pat Walters, and Vijay Pande. *Deep learning for the life sciences: applying deep learning to genomics, microscopy, drug discovery, and more*. "O'Reilly Media, Inc.", 2019.
- [70] Deyu Zou, Shikun Liu, Siqi Miao, Victor Fung, Shiyu Chang, and Pan Li. Gess: Benchmarking geometric deep learning under scientific applications with distribution shifts. In *The Thirty-eight Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2024.
- [71] Pankaj Kumar Singh, Arvind Negi, Pawan Kumar Gupta, Monika Chauhan, and Raj Kumar. Toxicophore exploration as a screening technology for drug design and discovery: techniques, scope and limitations. *Archives of toxicology*, 90:1785–1802, 2016.
- [72] Bernard Faller and Peter Ertl. Computational approaches to determine drug solubility. *Advanced drug delivery reviews*, 59(7):533–545, 2007.
- [73] Xuanyi Dong and Yi Yang. Searching for a robust neural architecture in four gpu hours. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 1761–1770, 2019.
- [74] James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al. Overcoming catastrophic forgetting in neural networks. *Proceedings of the national academy of sciences*, 114(13):3521–3526, 2017.
- [75] Jaejun Lee, Raphael Tang, and Jimmy Lin. What would elsa do? freezing layers during transformer fine-tuning. *arXiv preprint arXiv:1911.03090*, 2019.
- [76] Vinay V Ramasesh, Ethan Dyer, and Maithra Raghu. Anatomy of catastrophic forgetting: Hidden representations and task semantics. *arXiv preprint arXiv:2007.07400*, 2020.
- [77] Cian Eastwood, Ian Mason, Christopher KI Williams, and Bernhard Schölkopf. Source-free adaptation to measurement shift via bottom-up feature restoration. *arXiv preprint arXiv:2107.05446*, 2021.
- [78] Utku Evci, Vincent Dumoulin, Hugo Larochelle, and Michael C Mozer. Head2toe: Utilizing intermediate representations for better transfer learning. In *International Conference on Machine Learning*, pages 6009–6033. PMLR, 2022.

- [79] Niv Cohen, Rinon Gal, Eli A Meirom, Gal Chechik, and Yuval Atzmon. “this is my unicorn, fluffy”: Personalizing frozen vision-language representations. In *European conference on computer vision*, pages 558–577. Springer, 2022.
- [80] Jonathan Frankle, Gintare Karolina Dziugaite, Daniel Roy, and Michael Carbin. Linear mode connectivity and the lottery ticket hypothesis. In *International Conference on Machine Learning*, pages 3259–3269. PMLR, 2020.
- [81] Dongyue Li and Hongyang Zhang. Improved regularization and robustness for fine-tuning in neural networks. *Advances in Neural Information Processing Systems*, 34:27249–27262, 2021.
- [82] Henry Gouk, Timothy M Hospedales, and Massimiliano Pontil. Distance-based regularisation of deep networks for fine-tuning. *arXiv preprint arXiv:2002.08253*, 2020.
- [83] Junjiao Tian, Zecheng He, Xiaoliang Dai, Chih-Yao Ma, Yen-Cheng Liu, and Zsolt Kira. Trainable projected gradient method for robust fine-tuning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7836–7845, 2023.
- [84] Xingjian Li, Haoyi Xiong, Hanchao Wang, Yuxuan Rao, Liping Liu, Zeyu Chen, and Jun Huan. Delta: Deep learning transfer using feature map with attention for convolutional networks. *arXiv preprint arXiv:1901.09229*, 2019.
- [85] Xuhong Li, Yves Grandvalet, Rémi Flamary, Nicolas Courty, and Dejing Dou. Representation transfer by optimal transport. *arXiv preprint arXiv:2007.06737*, 2020.
- [86] Jiying Zhang, Xi Xiao, Long-Kai Huang, Yu Rong, and Yatao Bian. Fine-tuning graph neural networks via graph topology induced optimal transport. *arXiv preprint arXiv:2203.10453*, 2022.
- [87] Zhi Kou, Kaichao You, Mingsheng Long, and Jianmin Wang. Stochastic normalization. *Advances in Neural Information Processing Systems*, 33:16304–16314, 2020.
- [88] Beliz Gunel, Jingfei Du, Alexis Conneau, and Ves Stoyanov. Supervised contrastive learning for pre-trained language model fine-tuning. *arXiv preprint arXiv:2011.01403*, 2020.
- [89] Jincheng Zhong, Ximei Wang, Zhi Kou, Jianmin Wang, and Mingsheng Long. Bi-tuning of pre-trained representations. *arXiv preprint arXiv:2011.06182*, 2020.
- [90] Yifan Zhang, Bryan Hooi, Dapeng Hu, Jian Liang, and Jiashi Feng. Unleashing the power of contrastive self-supervised visual models via contrast-regularized fine-tuning. *Advances in Neural Information Processing Systems*, 34:29848–29860, 2021.
- [91] Haolin Pan, Yong Guo, Qinyi Deng, Haomin Yang, Jian Chen, and Yiqun Chen. Improving fine-tuning of self-supervised models with contrastive initialization. *Neural Networks*, 159: 198–207, 2023.
- [92] Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola, Aaron Maschinot, Ce Liu, and Dilip Krishnan. Supervised contrastive learning. *Advances in neural information processing systems*, 33:18661–18673, 2020.
- [93] Sachin Goyal, Ananya Kumar, Sankalp Garg, Zico Kolter, and Aditi Raghunathan. Finetune like you pretrain: Improved finetuning of zero-shot vision models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 19338–19347, 2023.
- [94] Changdae Oh, Junhyuk So, Hoyoon Byun, YongTaek Lim, Minchul Shin, Jong-June Jeon, and Kyungwoo Song. Geodesic multi-modal mixup for robust fine-tuning. *Advances in Neural Information Processing Systems*, 36, 2024.
- [95] Renhong Huang, Jiarong Xu, Xin Jiang, Chenglu Pan, Zhiming Yang, Chunping Wang, and Yang Yang. Measuring task similarity and its implication in fine-tuning graph neural networks. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 12617–12625, 2024.



- [96] Yifei Sun, Qi Zhu, Yang Yang, Chunping Wang, Tianyu Fan, Jiajun Zhu, and Lei Chen. Fine-tuning graph neural networks by preserving graph generative patterns. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 9053–9061, 2024.
- [97] Shengrui Li, Xueting Han, and Jing Bai. Adapterggnn: Parameter-efficient fine-tuning improves generalization in gnns. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 13600–13608, 2024.
- [98] Teague Sterling and John J Irwin. Zinc 15–ligand discovery for everyone. *Journal of chemical information and modeling*, 55(11):2324–2337, 2015.
- [99] Weihua Hu, Bowen Liu, Joseph Gomes, Marinka Zitnik, Percy Liang, Vijay Pande, and Jure Leskovec. Strategies for pre-training graph neural networks. *arXiv preprint arXiv:1905.12265*, 2019.
- [100] Sunghwan Kim, Jie Chen, Tiejun Cheng, Asta Gindulyte, Jia He, Siqian He, Qingliang Li, Benjamin A Shoemaker, Paul A Thiessen, Bo Yu, et al. Pubchem in 2021: new data content and improved web interfaces. *Nucleic acids research*, 49(D1):D1388–D1395, 2021.
- [101] Raghunathan Ramakrishnan, Pavlo O Dral, Matthias Rupp, and O Anatole Von Lilienfeld. Quantum chemistry structures and properties of 134 kilo molecules. *Scientific data*, 1(1):1–7, 2014.
- [102] Vijay Prakash Dwivedi, Chaitanya K Joshi, Anh Tuan Luu, Thomas Laurent, Yoshua Bengio, and Xavier Bresson. Benchmarking graph neural networks. *Journal of Machine Learning Research*, 24(43):1–48, 2023.
- [103] Aravind Subramanian, Rajiv Narayan, Steven M Corsello, David D Peck, Ted E Natoli, Xiaodong Lu, Joshua Gould, John F Davis, Andrew A Tubelli, Jacob K Asiedu, et al. A next generation connectivity map: L1000 platform and the first 1,000,000 profiles. *Cell*, 171(6):1437–1452, 2017.
- [104] Maho Nakata and Tomomi Shimazaki. Pubchemqc project: a large-scale first-principles electronic structure database for data-driven chemistry. *Journal of chemical information and modeling*, 57(6):1300–1308, 2017.
- [105] Yanli Wang, Tiejun Cheng, and Stephen H Bryant. Pubchem bioassay: a decade’s development toward open high-throughput screening data sharing. *SLAS DISCOVERY: Advancing Life Sciences R&D*, 22(6):655–666, 2017.
- [106] Ines Filipa Martins, Ana L Teixeira, Luis Pinheiro, and Andre O Falcao. A bayesian approach to in silico blood-brain barrier penetration modeling. *Journal of chemical information and modeling*, 52(6):1686–1697, 2012.
- [107] Ann M Richard, Richard S Judson, Keith A Houck, Christopher M Grulke, Patra Volarath, Inthirany Thillainadarajah, Chihae Yang, James Rathman, Matthew T Martin, John F Wambaugh, et al. Toxcast chemical landscape: paving the road to 21st century toxicology. *Chemical research in toxicology*, 29(8):1225–1251, 2016.
- [108] Michael Kuhn, Ivica Letunic, Lars Juhl Jensen, and Peer Bork. The sider database of drugs and side effects. *Nucleic acids research*, 44(D1):D1075–D1079, 2016.
- [109] Kaitlyn M Gayvert, Neel S Madhukar, and Olivier Elemento. A data-driven approach to predicting successes and failures of clinical trials. *Cell chemical biology*, 23(10):1294–1301, 2016.
- [110] Sebastian G Rohrer and Knut Baumann. Maximum unbiased validation (muv) data sets for virtual screening based on pubchem bioactivity data. *Journal of chemical information and modeling*, 49(2):169–184, 2009.
- [111] Daniel Zaharevitz. Aids antiviral screen data, 2015.
- [112] Govindan Subramanian, Bharath Ramsundar, Vijay Pande, and Rajiah Aldrin Denny. Computational modeling of  $\beta$ -secretase 1 (bace-1) inhibitors using ligand based approaches. *Journal of chemical information and modeling*, 56(10):1936–1949, 2016.

- 703 [113] John S Delaney. Esol: estimating aqueous solubility directly from molecular structure. *Journal*  
704 *of chemical information and computer sciences*, 44(3):1000–1005, 2004.
- 705 [114] Anna Gaulton, Louisa J Bellis, A Patricia Bento, Jon Chambers, Mark Davies, Anne Hersey,  
706 Yvonne Light, Shaun McGlinchey, David Michalovich, Bissan Al-Lazikani, et al. ChEMBL: a  
707 large-scale bioactivity database for drug discovery. *Nucleic acids research*, 40(D1):D1100–  
708 D1107, 2012.
- 709 [115] Francisco-Javier Gamo, Laura M Sanz, Jaume Vidal, Cristina De Cozar, Emilio Alvarez,  
710 Jose-Luis Lavandera, Dana E Vanderwall, Darren VS Green, Vinod Kumar, Samiul Hasan,  
711 et al. Thousands of chemical starting points for antimalarial lead identification. *Nature*, 465  
712 (7296):305–310, 2010.
- 713 [116] Johannes Hachmann, Roberto Olivares-Amaya, Sule Atahan-Evrenk, Carlos Amador-Bedolla,  
714 Roel S Sánchez-Carrera, Aryeh Gold-Parker, Leslie Vogt, Anna M Brockway, and Alán  
715 Aspuru-Guzik. The harvard clean energy project: large-scale computational screening and  
716 design of organic photovoltaics on the world community grid. *The Journal of Physical*  
717 *Chemistry Letters*, 2(17):2241–2251, 2011.

## NeurIPS Paper Checklist

### 1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [\[Yes\]](#)

Justification: The main claims are justified by the experimental results and discussion in sec 4. Also, we refer the claims to the later findings in introduction.

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

### 2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [\[Yes\]](#)

Justification: We include a section in appendix discussing our limitations and future works.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

### 3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [\[Yes\]](#)

Justification: For the proposition 1 included in the paper, we have the complete proof in the appendix.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

#### 4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [\[Yes\]](#)

Justification: We detail the experimental settings in the sec 3 and more hyperparameter tuning and dataset details in Appendix.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
  - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
  - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
  - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
  - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

#### 5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: The access to datasets and codes are provided and we include the detailed settings in the paper main text and appendix.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

## 6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: We clarify the dataset splits, hyperparameters and evaluation in the main text and appendix.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

## 7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: We include the standard deviation for all the experiments.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).

- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

## 8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: We include the computing resources in the appendix.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

## 9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

Answer: [Yes]

Justification: We follow the code of ethics.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

## 10. Broader impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]

Justification: We include the broader impact discussion in appendix.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.

- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

## 11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: Our paper poses no such risks.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

## 12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: We cite the sources of the datasets that are used.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets) has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.



- If this information is not available online, the authors are encouraged to reach out to the asset’s creators.

### 13. New assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [Yes]

Justification: The new assets introduced in the paper are well documented.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

### 14. Crowdsourcing and research with human subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: The paper does not involve crowdsourcing nor research with human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

### 15. Institutional review board (IRB) approvals or equivalent for research with human subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: The paper does not involve crowdsourcing nor research with human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

### 16. Declaration of LLM usage

1030 Question: Does the paper describe the usage of LLMs if it is an important, original, or  
1031 non-standard component of the core methods in this research? Note that if the LLM is used  
1032 only for writing, editing, or formatting purposes and does not impact the core methodology,  
1033 scientific rigorousness, or originality of the research, declaration is not required.

1034 Answer: [NA]

1035 Justification: The core method development in this research does not involve LLMs as any  
1036 important, original, or non-standard components.

1037 Guidelines:

- 1038 • The answer NA means that the core method development in this research does not  
1039 involve LLMs as any important, original, or non-standard components.
- 1040 • Please refer to our LLM policy (<https://neurips.cc/Conferences/2025/LLM>)  
1041 for what should or should not be described.

## 1042 A Proof of proposition 1

1043 **Proposition 2.** Given  $\tilde{\mathcal{L}}(\theta) = \mathcal{L}(\theta) + \frac{\delta}{2}\|\theta - \theta_{pre}\|_2^2$ , we define the optimal weights as  $\tilde{\theta}_{ft} =$   
 1044  $\operatorname{argmin}_{\theta} \tilde{\mathcal{L}}(\theta)$  and  $\theta_{ft}^* = \operatorname{argmin}_{\theta} \mathcal{L}(\theta)$ .

$$\mathbf{Q}^T \tilde{\theta}_{ft} = (\mathbf{\Lambda} + \delta \mathbf{I})^{-1} \mathbf{\Lambda} \mathbf{Q}^T \theta_{ft}^* + \delta (\mathbf{\Lambda} + \delta \mathbf{I})^{-1} \mathbf{Q}^T \theta_{pre} . \quad (4)$$

1045 where  $\mathbf{H}$  is the hessian matrix of  $\mathcal{L}$  evaluated at  $\theta_{ft}^*$  and  $\mathbf{H} = \mathbf{Q} \mathbf{\Lambda} \mathbf{Q}^T$ .

1046 *Proof.* Based on the quadratic approximation, we can approximate  $\mathcal{L}(\theta)$  as follows:

$$\begin{aligned} \mathcal{L}(\theta) &= \mathcal{L}(\theta_{ft}^*) + \mathcal{L}'(\theta_{ft}^*)(\theta - \theta_{ft}^*) + \frac{1}{2}(\theta - \theta_{ft}^*)^T \mathbf{H}(\theta - \theta_{ft}^*) \\ &= \mathcal{L}(\theta_{ft}^*) + \frac{1}{2}(\theta - \theta_{ft}^*)^T \mathbf{H}(\theta - \theta_{ft}^*) \end{aligned}$$

1047 since  $\mathcal{L}'(\theta_{ft}^*) = 0$  as  $\theta_{ft}^*$  is the minimum. Then, we add the weight regularization term, such that

$$\tilde{\mathcal{L}}(\theta) = \mathcal{L}(\theta_{ft}^*) + \frac{1}{2}(\theta - \theta_{ft}^*)^T \mathbf{H}(\theta - \theta_{ft}^*) + \frac{\delta}{2}\|\theta - \theta_{pre}\|_2^2$$

1048 Then, we solve for  $\tilde{\theta}_{ft}$  by setting  $\nabla \tilde{\mathcal{L}}(\theta) = 0$

$$\begin{aligned} \mathbf{H}(\tilde{\theta}_{ft} - \theta_{ft}^*) + \delta(\tilde{\theta}_{ft} - \theta_{pre}) &= 0 \\ (\mathbf{H} + \delta \mathbf{I})\tilde{\theta}_{ft} &= \mathbf{H}\theta_{ft}^* + \delta\theta_{pre} \\ \tilde{\theta}_{ft} &= (\mathbf{H} + \delta \mathbf{I})^{-1}(\mathbf{H}\theta_{ft}^* + \delta\theta_{pre}) \\ \tilde{\theta}_{ft} &= (\mathbf{Q} \mathbf{\Lambda} \mathbf{Q}^T + \delta \mathbf{I})^{-1}(\mathbf{Q} \mathbf{\Lambda} \mathbf{Q}^T \theta_{ft}^* + \delta\theta_{pre}) \\ \tilde{\theta}_{ft} &= (\mathbf{Q}(\mathbf{\Lambda} + \delta \mathbf{I})\mathbf{Q}^T)^{-1}(\mathbf{Q} \mathbf{\Lambda} \mathbf{Q}^T \theta_{ft}^* + \delta\theta_{pre}) \\ \mathbf{Q}^T \tilde{\theta}_{ft} &= (\mathbf{\Lambda} + \delta \mathbf{I})^{-1} \mathbf{\Lambda} \mathbf{Q}^T \theta_{ft}^* + \delta(\mathbf{\Lambda} + \delta \mathbf{I})^{-1} \mathbf{Q}^T \theta_{pre} \end{aligned}$$

1049

□

## 1050 B Limitations and Future Works

1051 We acknowledge certain limitations in this current work and highlight potential improvements for  
 1052 future research. Firstly, this study primarily focuses on the *property prediction tasks* of *small*  
 1053 *molecules* using *2D-graph* based foundation models. Exploring a broader array of foundation models  
 1054 across a wider range of applications—such as covering more areas like DNA, proteins, and materials,  
 1055 addressing various scientific tasks like linker design and chemical reactions, and incorporating  
 1056 diverse data formats like 3D geometric data—is highly worthwhile. Secondly, although we attempt  
 1057 to include many representative fine-tuning methods from various categories in this study, additional  
 1058 fine-tuning methods from different categories, as discussed in Appendix C, deserve investigation. For  
 1059 instance, future research could explore whether graph-specific fine-tuning methods offer additional  
 1060 benefits over non-graph fine-tuning approaches across various settings we design. Thirdly, the method  
 1061 DWiSE-FT introduced here is an extension and combination of existing methods directly motivated  
 1062 by our benchmark findings for specific fine-tuning scenarios. Future work may involve more thorough  
 1063 exploration into fine-tuning methodology design inspired by our current findings, and aiming to  
 1064 develop approaches effective across a broader range of fine-tuning scenarios.

1065 Regarding the broader impact, we recognize our work can be beneficial to the drug discovery and  
 1066 material science, but people should be aware of the misuse of molecular property prediction tasks to  
 1067 harmful chemical production.

## 1068 C Additional Discussions of Related Works

1069 In this section, we additionally discuss more related works about fine-tuning (FT) techniques. De-  
 1070 signing advanced fine-tuning strategies first gained attention in the computer vision (CV) and natural

1071 language processing (NLP) domains, leading to the development of various research directions. We  
1072 categorize the mainstream approaches into the following groups.

1073 **Partial model FT.** Numerous studies demonstrate that freezing certain parameters while fine-tuning  
1074 only specific components of the pre-trained model can help mitigate overfitting during the fine-tuning  
1075 process [74, 75, 76, 77, 78, 79]. Specifically, Linear Probing (LP) only trains the additional prediction  
1076 head during FT. Surgical FT [12] selectively fine-tunes a subset of layers based on the specific  
1077 mechanism of distribution shifts.

1078 **Weight-based FT** strategies mainly control the model weights during the FT. Specifically, WiSE-  
1079 FT [19], grounded on the linear mode connectivity [80], linearly interpolates between pre-training  
1080 parameters and fine-tuning parameters by a mixing coefficient.  $L^2$ -SP [14] regularizes the fine-tuning  
1081 model weights using  $L^2$  distance to constrain the parameters around pre-trained ones. REGSL [81]  
1082 further introduces a layer-wise parameter regularization, where the constraint strength gradually  
1083 reduces from the top to bottom layers. MARS-SP [82] adopts the projected gradient method (PGM)  
1084 to constrain the fine-tuning model weights within a small sphere centered on the pre-trained ones.  
1085 More recently, TPGM [83] further incorporates trainable weight projection radii constraint for each  
1086 layer, inspired by MARS-SP, to support layer-wise regularization optimization.

1087 **Representation-based FT** methods mainly regulate the latent representation space during FT. Feature-  
1088 map [13] adds distance regularization between the latent representations of pre-trained and fine-  
1089 tuned models to the Full-FT loss. DELTA [84] specifically constrains feature maps with the pre-  
1090 trained activations selected by channel-wise attention. BSS [17] penalizes the spectral components  
1091 corresponding to small singular values that are less transferable to prevent negative transfer. Li et al.  
1092 [85] proposes to transfer representations by encouraging small deviations from the reference one  
1093 through an regularizer based on optimal transport. Inspired by this, GTOT-Tuning [86] presents  
1094 optimal transport-based fine-tuning framework. LP-FT [20] first performs LP to prediction head  
1095 while keeping the pre-trained encoder fixed, followed by applying full-FT with the tuned prediction  
1096 head.

1097 **Architecture Refinement.** Besides the weight and representation based FT, StochNorm [87] refactors  
1098 the widely used Batch Normalization (BN) module and proposes Stochastic Normalization, to transfer  
1099 more pre-trained knowledge during the fine-tuning process and mitigate over-fitting.

1100 **Contrastive-based FT.** As discussed in Sec. 2, contrastive-based strategies have been widely demon-  
1101 strated to be effective in the pre-training stage. There are other works which explore its effectiveness  
1102 in the fine-tuning process. Gunel et al. [88], Bi-tuning [89], Core-tuning [90] and COIN [91] intro-  
1103 duce supervised contrastive learning [92] to better leverage the label information in the target datasets  
1104 with more discriminative representations as a result. More recently, FLYP [93] shows that simply  
1105 finetuning a classifier via the same contrastive loss as pre-training leads to superior performance in  
1106 finetuning image-text models. Oh et al. [94] fine-tunes the model with contrastive loss on additional  
1107 hard negative samples, which are generated by geodesic multi-modal Mixup, for robust fine-tuning in  
1108 multi-modal models.

1109 **Graph-specific fine-tuning techniques.** Apart from the CV and NLP domains, several fine-tuning  
1110 techniques specifically designed for the Graph-ML domain have recently been proposed. GTOT-  
1111 Tuning [86] achieves efficient knowledge transfer from the pre-trained models by an optimal transport-  
1112 based FT framework. Bridge-Tune [95] introduces an intermediate step that bridges pre-training  
1113 and downstream tasks by considering the task similarity between them. G-tuning [96] tunes the  
1114 pre-trained GNN so that it can reconstruct the generative patterns (graphons) of the downstream  
1115 graphs. Li et al. [97] leverages expressive adapters for GNNs, to boost adaptation to the downstream  
1116 tasks.

## 1117 D Pre-training Datasets Detail

1118 For self-supervised pre-training, *Mole-BERT* and *GraphMAE* are pre-trained over 2M molecules sam-  
1119 pled from the ZINC15 database [98], following previous works [99]. *MoleculeSTM* is initially trained  
1120 on PubChemSTM, a large multimodal dataset comprising over 280,000 chemical structure-text pairs  
1121 constructed from the PubChem database [100].

1122 For supervised pre-training, we use the models from the *Graphium* [32] library, which get pre-trained  
1123 on the Toymix and Largemix datasets provided in this library. The ToyMix dataset [32], totally 2.61M

graph-level data points, contains QM9 [101], Tox21 [42] and Zinc12K [102]. Specifically, QM9 consists of 19 graph-level quantum properties associated to an energy-minimized 3D conformation of the molecules. Zinc12K is to predict the constrained solubility which is the term  $\log P - SA - \text{cycle}$  (octanol-water partition coefficients,  $\log P$ , penalized by the synthetic accessibility score,  $SA$ , and number of long cycles,  $\text{cycle}$ ). The Largemix dataset, totally 343.4M graph-level data points and 197.7M node-level data points, contains four different datasets with tasks taken from quantum chemistry (PCQM4M\_G25\_N4), bio-assays (PCBA1328) and transcriptomics (L1000 VCAP and MCF7). Specifically, L1000 VCAP and MCF7 are from the LINCS L1000 database [103], which is generated using high-throughput transcriptomics. VCAP and MCF7 are, respectively, prostate cancer and human breast cancer cell lines. The PCQM4M\_G25\_N4 dataset is sourced from the PubChemQC project [104] that computed DFT properties on the energy-minimized conformation of 3.8M small molecules from PubChem. The PCBA1328 dataset, originally sourced from Wang et al. [105], comprises 1,328 assays and 1.56M molecules and contains information about a molecule’s biological activity across various assay settings. The pretraining dataset for GraphGPS is PCQM4Mv2, which is a large-scale molecular dataset containing 3.75M graphs curated from PubChemQC. The task is to regress the HOMO-LUMO gap, a quantum physical property originally calculated using Density Functional Theory.

## E Dataset Statistics

The statistics and references of the downstream datasets included in this work are shown in Table 5.

Table 5: Summary for the molecular datasets used for downstream FT, where “# TASKS” and “# MOLECULES” denote the number of tasks and molecules of each dataset, respectively.

DATASET	EVALUATION METRICS	TASK	# TASKS	# MOLECULES
BBBP [106]	AUC	CLASSIFICATION	1	2,039
Tox21	AUC	CLASSIFICATION	12	7,831
ToxCast [107]	AUC	CLASSIFICATION	617	8,576
SIDER [108]	AUC	CLASSIFICATION	27	1,427
CLINTOX [109]	AUC	CLASSIFICATION	2	1,478
MUV [110]	AUC	CLASSIFICATION	17	93,087
HIV ZAHAREVITZ [111]	AUC	CLASSIFICATION	1	41,127
BACE [112]	AUC	CLASSIFICATION	1	1,513
ESOL [113]	RMSE	REGRESSION	1	1,128
LIPO [114]	RMSE	REGRESSION	1	4,200
MALARIA [115]	RMSE	REGRESSION	1	9,999
CEP [116]	RMSE	REGRESSION	1	29,978

1142

## F Details of Experimental Implementation

**Pre-training Implementations.** For self-supervised pre-training, we use the open-source pre-trained checkpoints of Mole-BERT<sup>1</sup> and GraphMAE<sup>2</sup>. For supervised pre-training, we follow the same training pipeline as proposed in the Graphium<sup>3</sup>. We drop out the task head MLPs used for supervised pre-training during the downstream fine-tuning process, keeping only the graph encoder component. Note that we keep the architecture of the GNN encoder and the graph pooling strategy the same across the three pre-training models. Specifically, we use a 5-layer Graph Isomorphism Networks (GINs) with 300 hidden dimension and mean pooling as the readout function.

**Fine-tuning Implementations.** We keep the same training configurations across all the downstream datasets, pre-training models, and fine-tuning strategies, following Hu et al. [27]. Specifically, for each distinct setting, we fine-tune the pre-training models with 5 random seeds (0-4). We use a batch

<sup>1</sup><https://github.com/junxia97/Mole-BERT>

<sup>2</sup><https://github.com/THUDM/GraphMAE>

<sup>3</sup><https://github.com/datamol-io/graphium>

size of 32 and a dropout rate of 0.5. For each dataset, We train models for 100 epochs and report the test performance when the optimal validation performance is achieved.

**Hyperparameter Tuning.** We set learning rate to be 0.001 for all the methods and train for 100 epochs. Below is the detailed sets of hyperparameters we tuned for each fine-tuning strategy.

- *Surgical FT*: We tune  $k$  as which layer in GNN encoder to be updated from  $\{0, 1, 2, 3, 4\}$  since our backbone architecture is a 5-layer GIN.
- *WiSE-FT*: We tune the mixing coefficient  $\alpha$  from  $\{0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9\}$  to control the weight ensemble from pre-trained model and fine-tuned model. A larger  $\alpha$  indicates the weights are adopted more from the fine-tuned model.
- *$L^2$ -SP/ BSS/ Feature-map*: For these three methods that involve an additional regularization term in the loss, we tune the regularization coefficient  $\delta$  from  $\{1, 0.1, 0.01, 0.001, 0.0001\}$  to control the degree of regularization. For BSS, we follow the original paper and set  $k$  to be 1 meaning that we are regularizing the smallest singular value.
- *LP-FT*: We train the LP step before full fine-tuning for 100 epochs and then use the updated prediction head as initialization for the full-FT afterwards for 100 epochs. The training all use the default learning rate 0.001.
- *Full FT/ LP*: There is no additional hyperparameter tuning, where we use the default fine-tuning setting.
- *DWiSE-FT*: We tune the initialization of  $\alpha_i$  for each layer  $i$ , where we use the same value to initialize for all layers from  $\{0.9, 0.7, 0.5\}$  and the learning rate for validation loss descent from  $\{0.001, 0.005, 0.01\}$ . We tune  $\alpha$  over validation sets over 200 epochs.

Indeed, from the DWiSE-FT experiments with different starting points of mixing coefficients, the variance of final results is small since it will converge towards the optimal value of mixing coefficients regardless of the initial starting point given a reasonable training time.

**Computing Resources** The experiments are run on NVIDIA RTX A6000 with 48G memory.

## G Further Result Discussions

### G.1 Comparisons over pre-trained models

We mainly select the pretrained models based on their pre-training objective divided as supervised and self-supervised learning as discussed in 2. Then, among each category of pretrained models, we diversify with different architecture, model size and detailed training objective or pretraining dataset to discover the effect to the downstream finetuning method selection.

In the following, we will briefly discuss some more results that are not included in the main text with more pretrained models we tried. Detailed tables can be found in Appendix H

In general, we found the trend discussed in the main text about the difference of supervised pretrained model and self-supervised pretrained model hold in most cases. Especially, how they prefer over the representation based finetuning techniques or the weight based finetuning techniques remain consistent. However, some small variations may happen regarding the model size and architecture. For instance, for smaller model like 5 layer base GIN model, it is less likely to overfit on fewshot dataset compared to the larger scale graph transformer model. Also, the model expressiveness and capability will vary with different model scale. Therefore, we can compare the rank of different finetuning methods under pretrained models with the same scale, while it is not directly comparable if the model scale is significantly different.

For instance, both the Graphium model and the GraphGPS demonstrate superior performance from the representation based method like feature-map and BSS compared to other techniques. However, in contrast to the Graphium-Toy model results in the main text that feature-map perform better than BSS especially under the very few shot scenarios. In the GraphGPS results, we find that feature-map tend to be better with more finetuning samples and BSS tends to be better than feature-map in the fewshot cases. This might be due to the variation in the model size that leads to more overfitting, where BSS regularize over noisy feature space through penalizing smaller eigenvalues can be more crucial in reducing overfitting compared to feature-map. Also, we experience a change in pretrained

dataset compared to the ToyMix and LargeMix in the Graphium model, where the PCQM4Mv2 is less diversified. This might also cause the degraded performance of feature-map under GraphGPS with fewshot scenario since the learned representation from pretraining might not directly fit the downstream task. When there are more samples available, there might be a larger overlap with the learned representation space. Furthermore, we also observe the worse performance of LP and LP-FT under the larger model which coincides with findings in the main text from Graphium models.

Lastly, note that in the few-shot setting, GraphSTM underperforms other evaluated models in self-supervised pretraining. This is mainly because GraphSTM’s GNN encoder was specifically pretrained with graph-text alignment to enhance multi-modal tasks like structure-text retrieval. Therefore, the encoder would retain features optimized for cross-modal alignment rather than purely graph structural information. Since the downstream tasks in our benchmark do not involve text, the randomly initialized task head struggles to effectively utilize these features with limited data, whereas other models provide a more direct and task-relevant representation, leading to better performance in low-data scenarios.

## G.2 Comparisons over traditional method

To further understand the effect from foundation model pre-training and fine-tuning process, we include the XGBoost algorithm as a representative for the traditional method. Specifically, we tested the XGBoost algorithm under the Fewshot setting with 50, 100 and 500 samples to see whether it can surpass the performance of foundation model when the training data is scarce. The featurizer being used for the XGBoost model is the Extended Connectivity Circular Fingerprints adopted from the MoleculeNet paper. Then, we keep the exact same splits with the other experiments under random, scaffold and size split. From the result in table 16, we can conclude that foundation model result (e.g.) from Mole-BERT surpass the performance in XGBoost on almost all the settings. This indicates the benefit from the pretraining and finetuning framework and the value of our work in selecting the best finetuning technique given different pretraining situation.

## G.3 Additional findings

### • Finding 4. LP with pre-trained molecular representations from supervised pre-training surpasses full FT under few-shot fine-tuning, except for size splits.

For few-shot fine-tuning with 50 and 100 samples (*c.f.*, Fig. 2b and 2d), LP surpasses full FT in random and scaffold splits, differing from self-supervised pre-training discussed in (1a). This again supports the claim that directly adopting molecular representations from supervised pre-training retain useful knowledge for downstream tasks. But interestingly, this does not hold for size splits. We believe it is due to the susceptibility of graph level tasks under size shift, as noted in prior OOD studies [70]. Namely, the prediction head tends to overfit to the mapping from representations to output labels with molecules in a specific range of sizes, and thus cannot generalize to OOD molecules of different sizes.

### • Finding 5. Regulating feature representations brings significant benefits under few-shot fine-tuning but has only a marginal impact in non-few-shot fine-tuning.

Representation-based methods incorporates additional representation regularization in addition to full FT. BSS aims to eliminate noisy or non-transferable dimensions by regularizing small singular values of representations and Feature-map enforces a close distance of the fine-tuned representations to the pre-trained representations. Since the baseline full FT performs well under non-few-shot settings (*c.f.*, Tables 2 and 3), and pre-trained molecular representations are unsatisfying as discussed in Q1, having fine-tuned representations to unsatisfying pre-trained representations does not lead to any benefits. While under few-shot fine-tuning, representation regularization prevents overfitting with limited samples on top of full FT to some extend.

## H Additional Experimental Results

In this section, we present complementary baseline results that are not shown in the main text due to space limit. Table 1 is a summary of all pre-trained models we test on and their corresponding result tables for reference.



Table 6: Robust fine-tuning performance on 5 **Classification** datasets (AUC metrics) in the **Fewshot** setting (covering FEWSHOT-50, FEWSHOT-100, FEWSHOT-500), evaluated across 3 dataset splits (RANDOM, SCAFFOLD, SIZE), over **MOLE-BERT** and **GRAPHIUM-TOY** models. We **bold** and **underline** the best and second-best performances in each scenario.

SPLIT	METHODS	SELF-SUPERVISED PRE-TRAINING (MOLE-BERT)										SUPERVISED PRE-TRAINING (GRAPHIUM-TOY)									
		CLINTox	BBBP	BACE	HIV	SIDER	AVG	AVG-F	AVG-R			CLINTox	BBBP	BACE	HIV	SIDER	AVG	AVG-F	AVG-R		
RANDOM	FULL-FT	74.45 ± 2.10	88.56 ± 0.83	75.80 ± 0.43	57.41 ± 0.69	52.22 ± 0.48	69.69	69.22	4.40		70.14 ± 0.52	77.57 ± 0.01	80.45 ± 0.00	63.57 ± 0.00	55.57 ± 0.00	69.46	70.43	6.00			
	LP	73.50 ± 1.31	82.05 ± 0.37	75.04 ± 0.58	53.34 ± 2.39	51.40 ± 0.11	67.87	68.63	6.80		<u>84.09 ± 0.00</u>	81.04 ± 0.00	81.57 ± 0.00	49.05 ± 0.00	55.62 ± 0.00	70.27	72.74	4.20			
	SURGICAL-FT	<b>77.91 ± 1.25</b>	85.41 ± 0.66	<b>75.94 ± 0.40</b>	57.90 ± 0.40	51.90 ± 0.18	69.83	70.58	3.80		<u>77.51 ± 0.00</u>	84.90 ± 0.00	<b>81.93 ± 0.00</b>	64.72 ± 0.00	56.40 ± 0.00	73.14	74.76	2.40			
	LP-FT	72.66 ± 0.74	<b>88.99 ± 0.14</b>	75.18 ± 0.48	57.38 ± 0.37	51.68 ± 0.16	70.18	70.07	4.40		69.84 ± 0.00	80.15 ± 0.00	78.64 ± 0.00	<b>65.82 ± 0.00</b>	53.56 ± 0.00	69.60	71.43	6.00			
	WiSE-FT	76.12 ± 1.87	<u>88.72 ± 1.05</u>	75.59 ± 0.51	58.59 ± 0.77	<u>52.23 ± 0.50</u>	70.25	70.10	3.00		81.94 ± 0.03	83.74 ± 0.00	78.47 ± 0.00	63.17 ± 0.00	<u>56.44 ± 0.00</u>	72.75	74.53	4.40			
	L2-SP	79.27 ± 1.05	<u>88.50 ± 1.25</u>	75.17 ± 0.90	59.09 ± 1.33	<b>52.27 ± 0.52</b>	70.26	70.18	3.60		72.36 ± 1.46	81.07 ± 0.13	79.75 ± 0.50	63.68 ± 0.92	<u>55.48 ± 0.00</u>	70.45	71.90	5.20			
	FEATURE-MAP	74.43 ± 2.07	88.40 ± 0.84	73.84 ± 0.66	57.93 ± 1.13	51.82 ± 0.31	69.28	68.73	6.40		<b>84.80 ± 0.129</b>	<b>85.33 ± 0.021</b>	81.53 ± 0.194	60.64 ± 0.016	<b>56.40 ± 0.005</b>	73.76	75.66	2.60			
	BSS	75.31 ± 0.21	88.69 ± 0.54	75.50 ± 0.58	<b>59.19 ± 1.58</b>	52.13 ± 0.37	70.16	70.00	3.60		74.14 ± 2.15	77.94 ± 0.35	78.82 ± 1.14	64.45 ± 1.10	55.57 ± 0.05	70.18	72.18	5.20			
	FULL-FT	60.18 ± 1.70	59.58 ± 1.79	68.88 ± 2.31	55.47 ± 6.57	53.12 ± 0.45	59.47	58.44	6.00		61.94 ± 0.00	62.14 ± 0.00	76.51 ± 0.94	63.74 ± 0.00	54.02 ± 0.00	63.67	62.61	7.40			
	LP	60.36 ± 0.84	57.58 ± 0.82	70.25 ± 1.28	57.45 ± 5.76	51.76 ± 0.37	59.48	58.46	6.40		70.10 ± 0.00	57.74 ± 0.00	76.54 ± 0.00	65.43 ± 0.00	55.88 ± 0.00	66.94	66.57	4.80			
SCAFFOLD	SURGICAL-FT	60.80 ± 1.05	<b>60.86 ± 0.98</b>	71.16 ± 0.84	58.60 ± 6.33	52.24 ± 0.21	60.73	60.09	4.00		<u>71.30 ± 0.00</u>	63.24 ± 0.00	76.34 ± 0.00	66.81 ± 0.00	<u>56.56 ± 0.00</u>	66.85	67.12	4.40			
	LP-FT	59.59 ± 1.11	60.36 ± 1.20	71.57 ± 0.82	56.18 ± 0.07	<b>53.31 ± 0.29</b>	60.20	58.71	4.40		65.20 ± 0.00	63.16 ± 0.00	77.15 ± 0.00	66.60 ± 0.00	53.85 ± 0.00	65.17	65.02	6.80			
	WiSE-FT	67.60 ± 5.67	69.51 ± 1.64	<u>72.25 ± 1.25</u>	63.65 ± 2.09	59.66 ± 0.93	62.93	63.92	3.00		67.34 ± 0.00	65.55 ± 0.00	78.66 ± 0.00	65.28 ± 0.00	55.17 ± 0.00	66.40	66.06	4.80			
	L2-SP	61.76 ± 1.22	59.53 ± 2.09	70.81 ± 0.79	<b>64.76 ± 2.40</b>	52.95 ± 0.45	61.96	62.02	3.60		<b>83.15 ± 0.03</b>	<u>66.76 ± 0.00</u>	<u>78.75 ± 0.74</u>	<u>68.22 ± 0.02</u>	58.86 ± 0.00	70.55	71.24	2.20			
	FEATURE-MAP	61.30 ± 1.94	55.91 ± 2.04	65.37 ± 0.59	61.18 ± 2.35	52.64 ± 1.03	59.28	59.36	5.60		77.49 ± 0.04	<b>67.13 ± 0.01</b>	78.57 ± 0.05	61.39 ± 0.01	<b>56.74 ± 0.00</b>	68.86	69.67	5.20			
	BSS	<b>67.94 ± 2.58</b>	60.40 ± 2.18	70.51 ± 1.82	60.39 ± 2.23	53.18 ± 0.46	62.48	62.91	3.00		69.74 ± 0.02	65.64 ± 0.00	<b>79.10 ± 0.00</b>	<b>68.47 ± 0.01</b>	54.97 ± 0.03	67.58	67.95	3.20			
	FULL-FT	66.75 ± 0.92	80.03 ± 0.54	43.23 ± 1.52	62.00 ± 3.04	47.81 ± 0.77	59.96	58.85	5.80		67.61 ± 0.01	<b>71.89 ± 5.76</b>	48.57 ± 0.01	52.54 ± 0.00	53.48 ± 0.00	58.82	57.88	5.20			
	LP	69.17 ± 0.41	78.19 ± 0.32	39.81 ± 0.34	48.97 ± 1.66	46.13 ± 0.24	56.45	54.76	7.00		71.21 ± 0.01	57.79 ± 0.00	40.44 ± 0.01	48.13 ± 0.00	<b>55.62 ± 0.00</b>	54.64	53.85	6.00			
	SURGICAL-FT	68.76 ± 0.63	82.19 ± 0.86	42.26 ± 2.37	56.73 ± 1.32	46.77 ± 0.14	59.34	57.71	5.60		71.70 ± 0.01	68.21 ± 0.00	46.06 ± 0.01	53.09 ± 0.00	54.86 ± 0.00	58.78	58.72	5.00			
	LP-FT	60.43 ± 0.30	82.00 ± 0.83	42.83 ± 1.39	61.12 ± 1.15	48.77 ± 0.32	60.83	59.77	4.20		68.90 ± 0.01	65.03 ± 0.01	47.57 ± 0.00	47.28 ± 0.00	54.15 ± 0.00	56.59	55.58	6.20			
SIZE	WiSE-FT	<b>70.76 ± 1.31</b>	81.92 ± 3.19	<b>65.58 ± 2.49</b>	56.58 ± 10.19	47.24 ± 0.57	64.42	64.31	4.00		72.03 ± 0.01	70.14 ± 5.65	45.24 ± 0.01	53.43 ± 0.00	53.59 ± 0.00	58.89	59.05	4.80			
	L2-SP	69.09 ± 1.06	<b>83.98 ± 1.98</b>	52.70 ± 4.51	63.68 ± 3.16	<b>50.80 ± 0.97</b>	64.05	61.82	2.60		72.95 ± 0.73	63.18 ± 0.27	63.46 ± 3.90	<b>66.83 ± 0.03</b>	54.89 ± 0.01	64.30	64.56	3.20			
	FEATURE-MAP	<b>67.57 ± 1.45</b>	82.52 ± 0.74	51.61 ± 1.25	<b>66.37 ± 3.56</b>	49.65 ± 0.57	63.54	61.85	3.00		<b>70.45 ± 0.06</b>	71.39 ± 0.05	65.20 ± 0.01	57.29 ± 0.43	53.01 ± 0.01	64.71	64.63	3.80			
	BSS	67.05 ± 1.32	80.29 ± 3.12	50.73 ± 6.35	62.56 ± 2.53	49.05 ± 0.64	62.06	60.31	4.00		72.26 ± 0.16	68.79 ± 6.08	<b>66.98 ± 0.01</b>	55.61 ± 0.00	<u>55.40 ± 0.01</u>	63.81	63.79	2.20			
	<b>FEWSHOT-100</b>																				
	FULL-FT	70.17 ± 2.25	86.87 ± 0.80	79.91 ± 0.70	<u>80.88 ± 1.37</u>	53.88 ± 0.69	72.05	73.16	4.20		69.31 ± 1.27	82.85 ± 0.00	83.76 ± 0.44	64.82 ± 2.36	56.88 ± 0.00	71.52	72.83	5.00			
	LP	70.45 ± 0.85	84.18 ± 0.82	73.16 ± 0.46	51.26 ± 1.30	52.78 ± 0.31	59.17	68.46	7.20		<u>81.55 ± 0.00</u>	80.80 ± 0.00	79.25 ± 0.00	51.60 ± 0.00	52.78 ± 0.00	70.37	72.63	5.00			
	SURGICAL-FT	81.54 ± 1.62	85.66 ± 0.52	77.00 ± 0.74	59.34 ± 0.42	53.63 ± 0.44	71.43	72.63	6.40		<u>75.51 ± 0.00</u>	<b>86.37 ± 0.00</b>	<b>84.51 ± 0.00</b>	<b>66.28 ± 0.00</b>	58.87 ± 0.00	74.31	75.43	2.00			
	LP-FT	<u>78.86 ± 1.12</u>	87.26 ± 0.81	78.86 ± 0.48	59.37 ± 0.51	54.31 ± 0.32	71.93	72.70	3.80		81.73 ± 0.32	83.54 ± 0.02	81.91 ± 0.04	<u>65.46 ± 0.62</u>	58.74 ± 0.00	74.28	76.37	5.20			
	WiSE-FT	<b>85.55 ± 1.43</b>	86.76 ± 0.42	74.53 ± 0.97	<b>61.90 ± 1.36</b>	<b>56.41 ± 0.69</b>	73.03	73.39	3.00		71.90 ± 1.49	83.15 ± 0.83	83.63 ± 0.95	57.66 ± 0.00	72.03	72.86	5.00				
L2-SP	79.13 ± 3.68	86.89 ± 0.40	79.66 ± 0.35	59.92 ± 1.04	54.64 ± 0.35	72.05	72.90	3.80		76.28 ± 0.02	81.15 ± 1.52	80.71 ± 1.44	64.00 ± 0.98	<b>59.02 ± 0.54</b>	72.23	73.66	4.40				
FEATURE-MAP	78.12 ± 2.01	<b>87.80 ± 0.62</b>	73.50 ± 0.69	59.97 ± 0.75	53.59 ± 0.24	70.58	70.53	5.40		<b>82.51 ± 0.15</b>	85.94 ± 0.56	82.09 ± 1.02	63.34 ± 0.11	57.52 ± 0.05	74.34	75.98	3.60				
BSS	79.00 ± 4.62	<u>87.38 ± 0.52</u>	<b>80.12 ± 0.33</b>	60.22 ± 1.07	53.88 ± 0.72	72.12	73.11	3.20		72.38 ± 1.42	<u>80.11 ± 0.78</u>	81.64 ± 0.64	63.65 ± 0.65	56.85 ± 0.81	70.93	72.05	6.80				
RANDOM	FULL-FT	70.51 ± 0.71	62.11 ± 1.32	68.39 ± 3.19	61.60 ± 1.74	52.20 ± 0.26	62.96	64.03	4.80		70.75 ± 0.00	65.39 ± 0.25	7								

Table 7: Robust fine-tuning performance on 4 **Regression** datasets (RMSE metrics) in the **Fewshot** setting (covering FEWSHOT-50, FEWSHOT-100, and FEWSHOT-500), evaluated across 3 dataset splits (RANDOM, SCAFFOLD, SIZE), over **MOLE-BERT** and **GRAPHIUM-TOY** models. AVG-R, AVG-R\* denote the average rank and the rank based on the average normalized performance over all the datasets for each evaluated method, respectively. Standard deviations across five replicates are shown in parentheses. We **bold** and underline the best and second-best performances in each scenario.

SPLIT	METHODS	SELF-SUPERVISED PRE-TRAINING (MOLE-BERT)					SUPERVISED PRE-TRAINING (GRAPHIUM-TOY)				
		ESOL	LIPO	MALARIA	CEP	AVG-R   AVG-R*	ESOL	LIPO	MALARIA	CEP	AVG-R   AVG-R*
		FEWSHOT-50									
RANDOM	FULL-FT	1.390 ± 0.051	1.189 ± 0.016	1.276 ± 0.006	2.383 ± 0.046	3.50   4	1.223 ± 0.000	1.062 ± 0.000	1.284 ± 0.000	2.359 ± 0.000	6.25   7
	LP	2.654 ± 0.016	1.825 ± 0.011	1.296 ± 0.005	3.736 ± 0.020	8.00   8	1.085 ± 0.000	1.072 ± 0.000	1.272 ± 0.000	2.571 ± 0.000	4.00   3
	SURGICAL-FT	2.647 ± 0.022	1.618 ± 0.014	1.295 ± 0.004	3.596 ± 0.037	7.00   7	1.174 ± 0.000	1.009 ± 0.000	1.277 ± 0.000	2.355 ± 0.000	3.25   2
	LP-FT	1.422 ± 0.027	1.237 ± 0.027	1.291 ± 0.005	2.296 ± 0.012	5.25   6	1.386 ± 0.000	1.010 ± 0.000	1.286 ± 0.000	2.287 ± 0.000	5.25   8
	WISE-FT	1.384 ± 0.047	1.212 ± 0.020	1.276 ± 0.007	2.410 ± 0.051	4.25   5	1.219 ± 0.000	1.060 ± 0.000	1.280 ± 0.000	2.366 ± 0.000	5.25   4
	L2-SP	1.372 ± 0.029	1.196 ± 0.019	1.277 ± 0.006	2.280 ± 0.031	3.25   3	1.147 ± 0.026	1.092 ± 0.001	1.283 ± 0.000	2.312 ± 0.020	5.00   5
SCAFFOLD	FEATURE-MAP	1.329 ± 0.021	1.164 ± 0.010	1.271 ± 0.007	2.448 ± 0.010	2.25   1	1.089 ± 0.001	1.046 ± 0.000	1.276 ± 0.000	2.191 ± 0.017	2.00   1
	BSS	1.365 ± 0.028	1.186 ± 0.017	1.277 ± 0.006	2.275 ± 0.022	2.50   2	1.175 ± 0.011	1.128 ± 0.035	1.281 ± 0.000	2.262 ± 0.064	5.00   6
	FULL-FT	1.696 ± 0.058	1.124 ± 0.006	1.178 ± 0.005	2.356 ± 0.033	4.25   5	1.353 ± 0.000	1.071 ± 0.000	1.168 ± 0.000	2.001 ± 0.000	5.75   8
	LP	3.754 ± 0.020	1.858 ± 0.005	1.167 ± 0.002	3.849 ± 0.009	7.25   8	1.226 ± 0.000	1.013 ± 0.000	1.166 ± 0.000	2.450 ± 0.000	4.00   6
	SURGICAL-FT	3.599 ± 0.039	1.843 ± 0.006	1.167 ± 0.003	3.819 ± 0.017	6.75   7	1.239 ± 0.000	1.019 ± 0.000	1.162 ± 0.000	2.083 ± 0.000	3.00   2
	LP-FT	1.822 ± 0.014	1.134 ± 0.012	1.184 ± 0.004	2.292 ± 0.026	4.50   6	1.283 ± 0.000	1.033 ± 0.000	1.169 ± 0.000	1.949 ± 0.000	4.75   5
SIZE	WISE-FT	1.842 ± 0.056	1.177 ± 0.009	1.162 ± 0.004	2.454 ± 0.043	5.00   4	1.320 ± 0.000	1.071 ± 0.000	1.168 ± 0.000	1.992 ± 0.000	5.75   7
	L2-SP	1.699 ± 0.049	1.086 ± 0.009	1.162 ± 0.002	2.331 ± 0.024	2.75   2	1.273 ± 0.047	1.015 ± 0.007	1.166 ± 0.000	2.132 ± 0.048	6.00   4
	FEATURE-MAP	1.823 ± 0.028	1.036 ± 0.007	1.159 ± 0.000	2.425 ± 0.012	3.00   1	1.213 ± 0.001	0.991 ± 0.000	1.164 ± 0.000	2.128 ± 0.006	2.50   1
	BSS	1.680 ± 0.042	1.114 ± 0.008	1.165 ± 0.001	2.319 ± 0.025	2.50   3	1.222 ± 0.012	1.039 ± 0.000	1.166 ± 0.000	2.121 ± 0.029	4.25   3
	FULL-FT	2.382 ± 0.079	1.297 ± 0.040	0.929 ± 0.004	2.656 ± 0.039	2.75   4	1.441 ± 0.000	1.055 ± 0.000	0.914 ± 0.000	2.329 ± 0.000	5.00   7
	LP	4.534 ± 0.021	2.157 ± 0.012	0.941 ± 0.004	4.706 ± 0.022	7.75   8	1.443 ± 0.000	1.003 ± 0.000	0.936 ± 0.000	2.688 ± 0.000	6.50   8
RANDOM	SURGICAL-FT	4.344 ± 0.026	2.111 ± 0.021	0.943 ± 0.004	4.265 ± 0.028	7.25   7	1.469 ± 0.000	1.015 ± 0.000	0.914 ± 0.000	2.313 ± 0.000	5.25   5
	LP-FT	2.421 ± 0.060	1.395 ± 0.018	0.939 ± 0.007	2.525 ± 0.013	4.50   6	1.395 ± 0.000	0.999 ± 0.000	0.907 ± 0.000	2.410 ± 0.000	3.50   1
	WISE-FT	2.615 ± 0.072	1.391 ± 0.042	0.929 ± 0.004	2.762 ± 0.053	5.50   5	1.411 ± 0.000	1.071 ± 0.000	0.905 ± 0.000	2.324 ± 0.000	3.50   4
	L2-SP	2.393 ± 0.068	1.306 ± 0.037	0.915 ± 0.002	2.497 ± 0.019	2.00   2	1.446 ± 0.055	0.997 ± 0.000	0.908 ± 0.000	2.340 ± 0.020	4.25   3
	FEATURE-MAP	2.422 ± 0.021	1.327 ± 0.022	0.911 ± 0.002	2.659 ± 0.021	3.75   1	1.415 ± 0.005	0.989 ± 0.027	0.921 ± 0.002	2.254 ± 0.001	3.00   2
	BSS	2.369 ± 0.075	1.319 ± 0.050	0.925 ± 0.003	2.563 ± 0.022	2.50   3	1.499 ± 0.028	0.997 ± 0.000	0.907 ± 0.000	2.381 ± 0.006	5.00   6
FEWSHOT-100											
RANDOM	FULL-FT	1.141 ± 0.030	1.141 ± 0.023	1.256 ± 0.006	2.150 ± 0.021	2.00   1	1.191 ± 0.000	1.103 ± 0.000	1.258 ± 0.000	2.076 ± 0.118	5.25   4
	LP	2.273 ± 0.029	1.569 ± 0.008	1.280 ± 0.003	3.235 ± 0.019	8.00   8	1.066 ± 0.000	1.045 ± 0.000	1.267 ± 0.000	2.383 ± 0.000	4.75   5
	SURGICAL-FT	1.953 ± 0.039	1.281 ± 0.020	1.270 ± 0.006	3.019 ± 0.047	6.75   7	1.075 ± 0.000	1.030 ± 0.000	1.296 ± 0.000	1.935 ± 0.000	2.75   2
	LP-FT	1.244 ± 0.057	1.147 ± 0.018	1.277 ± 0.003	2.156 ± 0.019	5.25   6	1.689 ± 0.000	1.097 ± 0.000	1.273 ± 0.000	2.044 ± 0.015	6.25   8
	WISE-FT	1.189 ± 0.030	1.142 ± 0.025	1.256 ± 0.006	2.211 ± 0.028	3.50   2	1.131 ± 0.000	1.078 ± 0.000	1.256 ± 0.000	2.001 ± 0.071	3.75   3
	L2-SP	1.161 ± 0.016	1.149 ± 0.007	1.260 ± 0.004	2.131 ± 0.014	3.25   4	1.098 ± 0.012	1.077 ± 0.001	1.270 ± 0.001	2.261 ± 0.008	5.25   6
SCAFFOLD	FEATURE-MAP	1.120 ± 0.038	1.139 ± 0.017	1.266 ± 0.004	2.283 ± 0.011	3.25   5	0.995 ± 0.018	1.025 ± 0.000	1.258 ± 0.003	1.937 ± 0.023	1.75   1
	BSS	1.199 ± 0.033	1.149 ± 0.023	1.259 ± 0.006	2.132 ± 0.019	4.00   3	1.055 ± 0.009	1.136 ± 0.000	1.274 ± 0.000	2.269 ± 0.010	6.25   7
	FULL-FT	1.436 ± 0.054	1.026 ± 0.009	1.160 ± 0.011	2.198 ± 0.034	3.25   4	1.111 ± 0.000	1.037 ± 0.000	1.172 ± 0.000	1.965 ± 0.023	5.00   6
	LP	3.255 ± 0.025	1.503 ± 0.008	1.154 ± 0.003	3.350 ± 0.007	7.00   8	1.228 ± 0.000	0.960 ± 0.000	1.162 ± 0.000	2.423 ± 0.000	4.50   5
	SURGICAL-FT	2.587 ± 0.076	1.192 ± 0.015	1.156 ± 0.003	2.914 ± 0.066	6.50   7	1.087 ± 0.000	0.966 ± 0.000	1.156 ± 0.000	1.959 ± 0.000	1.25   1
	LP-FT	1.544 ± 0.042	1.010 ± 0.011	1.163 ± 0.004	2.187 ± 0.034	4.00   6	1.111 ± 0.000	0.984 ± 0.000	1.173 ± 0.000	2.149 ± 0.012	5.25   4
SIZE	WISE-FT	1.544 ± 0.063	1.041 ± 0.017	1.151 ± 0.007	2.301 ± 0.042	4.50   3	1.110 ± 0.000	1.027 ± 0.000	1.169 ± 0.000	2.013 ± 0.049	4.25   3
	L2-SP	1.473 ± 0.009	0.961 ± 0.003	1.153 ± 0.002	2.201 ± 0.038	2.75   2	1.252 ± 0.021	0.994 ± 0.013	1.163 ± 0.000	2.367 ± 0.052	5.75   7
	FEATURE-MAP	1.677 ± 0.020	0.937 ± 0.008	1.149 ± 0.003	2.356 ± 0.018	3.50   1	1.158 ± 0.020	0.966 ± 0.010	1.161 ± 0.000	2.024 ± 0.019	3.50   2
	BSS	1.463 ± 0.008	1.040 ± 0.018	1.160 ± 0.006	2.210 ± 0.018	4.50   5	1.253 ± 0.027	1.033 ± 0.015	1.167 ± 0.000	2.333 ± 0.022	6.50   8
	FULL-FT	1.889 ± 0.065	1.077 ± 0.028	0.918 ± 0.005	2.425 ± 0.024	4.00   3	1.411 ± 0.000	0.962 ± 0.000	0.921 ± 0.006	2.328 ± 0.015	4.75   5
	LP	3.851 ± 0.033	1.676 ± 0.025	0.911 ± 0.003	4.115 ± 0.038	6.75   8	1.253 ± 0.000	0.981 ± 0.000	0.924 ± 0.000	2.635 ± 0.000	6.00   8
RANDOM	SURGICAL-FT	3.237 ± 0.085	1.374 ± 0.031	0.912 ± 0.002	3.174 ± 0.048	6.25   7	1.329 ± 0.000	0.965 ± 0.000	0.910 ± 0.000	2.283 ± 0.000	3.25   2
	LP-FT	1.831 ± 0.066	1.085 ± 0.014	0.920 ± 0.008	2.468 ± 0.021	4.75   4	1.242 ± 0.000	0.962 ± 0.000	0.912 ± 0.000	2.375 ± 0.000	3.50   1
	WISE-FT	2.216 ± 0.056	1.124 ± 0.031	0.917 ± 0.004	2.543 ± 0.027	5.75   5	1.398 ± 0.000	0.963 ± 0.000	0.907 ± 0.002	2.319 ± 0.014	3.75   4
	L2-SP	1.731 ± 0.071	1.025 ± 0.028	0.905 ± 0.002	2.424 ± 0.024	1.25   1	1.418 ± 0.035	0.998 ± 0.038	0.906 ± 0.000	2.436 ± 0.072	5.50   6
	FEATURE-MAP	2.135 ± 0.077	1.049 ± 0.013	0.898 ± 0.003	2.500 ± 0.017	3.25   2	1.335 ± 0.005	0.967 ± 0.008	0.		

Table 8: Robust fine-tuning performance on 8 **Classification** datasets (AUC metrics) in the **Non-Fewshot** setting, evaluated across 3 dataset splits (RANDOM, SCAFFOLD, SIZE), over **MOLECULESTM** and **GRAPHIUM-LARGE** models. AVG, AVG-F, AVG-R denote the average AUC, average AUC without max and min values, and average rank over all the datasets for each method, respectively. Standard deviations across five replicates are shown in parentheses. We **bold** and underline the best and second-best performances in each scenario.

SPLIT	METHODS	CLINTOX	BBBP	BACE	HIV	MUV	SIDER	Tox21	ToxCast	AVG	AVG-F	AVG-R
SELF-SUPERVISED PRE-TRAINING (MOLECULESTM)												
RANDOM	FULL-FT	89.90 $\pm$ 1.49	93.43 $\pm$ 0.99	89.82 $\pm$ 1.08	84.72 $\pm$ 1.11	77.82 $\pm$ 3.46	62.12 $\pm$ 1.15	82.49 $\pm$ 0.41	72.95 $\pm$ 0.31	81.66	82.95	3.62
	LP	74.32 $\pm$ 1.90	84.76 $\pm$ 0.29	74.85 $\pm$ 0.27	74.15 $\pm$ 0.69	76.86 $\pm$ 1.07	59.69 $\pm$ 0.24	73.72 $\pm$ 0.20	66.19 $\pm$ 0.14	73.07	73.35	7.75
	SURGICAL-FT	86.04 $\pm$ 0.89	93.68 $\pm$ 0.51	89.99 $\pm$ 0.46	<b>85.68 <math>\pm</math> 0.84</b>	<u>79.59 <math>\pm</math> 2.47</u>	<b>63.64 <math>\pm</math> 0.78</b>	81.84 $\pm$ 0.66	71.83 $\pm$ 0.55	81.54	82.50	3.38
	LP-FT	86.39 $\pm$ 1.85	93.72 $\pm$ 0.93	89.82 $\pm$ 0.57	84.17 $\pm$ 1.41	76.87 $\pm$ 2.38	62.19 $\pm$ 1.00	82.54 $\pm$ 0.51	72.19 $\pm$ 0.52	80.99	82.00	3.75
	WiSE-FT	<b>90.35 <math>\pm</math> 1.26</b>	92.93 $\pm$ 0.80	<b>90.41 <math>\pm</math> 0.86</b>	84.38 $\pm$ 1.05	77.23 $\pm$ 3.08	62.17 $\pm$ 1.25	82.67 $\pm$ 0.32	<u>73.08 <math>\pm</math> 0.32</u>	81.65	83.02	2.88
	L <sup>2</sup> -SP	89.69 $\pm$ 1.39	<u>93.77 <math>\pm</math> 0.37</u>	89.21 $\pm$ 0.92	81.94 $\pm$ 1.20	50.21 $\pm$ 4.41	61.07 $\pm$ 1.22	<u>82.97 <math>\pm</math> 0.39</u>	71.02 $\pm$ 0.57	77.48	79.32	5.00
	FEATURE-MAP	79.93 $\pm$ 1.54	<b>90.59 <math>\pm</math> 0.39</b>	83.69 $\pm$ 0.24	77.66 $\pm$ 0.46	<b>80.03 <math>\pm</math> 1.01</b>	59.93 $\pm$ 0.14	75.32 $\pm$ 0.19	67.51 $\pm$ 0.30	76.83	77.36	6.25
BSS	<u>90.17 <math>\pm</math> 2.84</u>	<b>94.16 <math>\pm</math> 0.55</b>	<b>89.74 <math>\pm</math> 0.12</b>	83.96 $\pm$ 1.29	76.64 $\pm$ 1.29	61.87 $\pm$ 0.69	<b>83.26 <math>\pm</math> 0.57</b>	<b>74.55 <math>\pm</math> 0.31</b>	81.79	83.05	3.38	
SCAFFOLD	FULL-FT	74.94 $\pm$ 7.23	68.62 $\pm$ 0.80	75.35 $\pm$ 2.06	76.03 $\pm$ 0.91	73.43 $\pm$ 2.50	57.88 $\pm$ 1.18	76.67 $\pm$ 0.68	63.62 $\pm$ 0.27	70.82	72.00	4.25
	LP	65.07 $\pm$ 1.08	59.39 $\pm$ 0.35	69.24 $\pm$ 0.16	69.97 $\pm$ 0.57	71.81 $\pm$ 2.40	59.93 $\pm$ 0.37	69.87 $\pm$ 0.28	60.05 $\pm$ 0.25	65.67	65.69	7.00
	SURGICAL-FT	71.07 $\pm$ 4.16	67.78 $\pm$ 0.60	80.16 $\pm$ 2.36	<b>76.80 <math>\pm</math> 1.06</b>	<u>75.87 <math>\pm</math> 0.82</u>	59.24 $\pm$ 1.22	75.54 $\pm$ 0.64	63.27 $\pm$ 0.70	71.22	71.72	3.75
	LP-FT	<b>77.27 <math>\pm</math> 4.28</b>	67.78 $\pm$ 1.42	75.33 $\pm$ 1.14	76.68 $\pm$ 0.82	71.36 $\pm$ 1.39	58.51 $\pm$ 1.15	76.85 $\pm$ 0.63	62.98 $\pm$ 0.51	70.48	71.41	4.62
	WiSE-FT	<b>77.27 <math>\pm</math> 4.28</b>	68.72 $\pm$ 0.75	77.37 $\pm$ 1.44	75.91 $\pm$ 0.74	74.38 $\pm$ 2.20	58.19 $\pm$ 1.26	<b>76.89 <math>\pm</math> 0.69</b>	<b>64.05 <math>\pm</math> 0.34</b>	71.60	72.87	3.12
	L <sup>2</sup> -SP	74.62 $\pm$ 4.99	68.30 $\pm$ 1.19	79.91 $\pm$ 2.29	73.97 $\pm$ 0.78	61.62 $\pm$ 2.07	59.78 $\pm$ 0.33	75.39 $\pm$ 0.51	62.34 $\pm$ 0.82	69.49	69.37	5.25
	FEATURE-MAP	61.06 $\pm$ 2.00	65.12 $\pm$ 1.98	<b>82.66 <math>\pm</math> 0.62</b>	74.54 $\pm$ 1.00	72.81 $\pm$ 1.16	<b>60.47 <math>\pm</math> 0.45</b>	70.39 $\pm$ 0.11	60.10 $\pm$ 0.19	68.39	67.40	5.25
BSS	73.89 $\pm$ 6.04	<b>70.04 <math>\pm</math> 0.20</b>	77.94 $\pm$ 0.24	76.28 $\pm$ 1.28	<b>76.20 <math>\pm</math> 1.33</b>	<u>59.99 <math>\pm</math> 1.39</u>	75.86 $\pm$ 1.08	<u>63.62 <math>\pm</math> 0.50</u>	71.73	72.65	2.75	
SIZE	FULL-FT	61.94 $\pm$ 2.67	82.80 $\pm$ 2.31	63.62 $\pm$ 1.19	77.81 $\pm$ 2.99	72.05 $\pm$ 2.96	54.92 $\pm$ 0.79	71.08 $\pm$ 0.77	62.47 $\pm$ 0.83	68.34	68.16	5.12
	LP	55.54 $\pm$ 0.65	75.89 $\pm$ 0.90	42.31 $\pm$ 0.48	67.54 $\pm$ 1.27	69.87 $\pm$ 1.51	53.74 $\pm$ 0.43	68.10 $\pm$ 0.39	57.50 $\pm$ 0.19	61.31	62.05	7.75
	SURGICAL-FT	64.54 $\pm$ 8.03	<b>88.90 <math>\pm</math> 0.74</b>	61.99 $\pm$ 2.13	78.10 $\pm$ 0.96	<b>76.07 <math>\pm</math> 0.57</b>	<b>57.13 <math>\pm</math> 1.87</b>	72.24 $\pm$ 0.28	60.52 $\pm$ 0.95	69.94	68.91	2.50
	LP-FT	63.79 $\pm$ 3.29	83.12 $\pm$ 5.20	<b>65.48 <math>\pm</math> 0.70</b>	76.47 $\pm$ 3.53	72.24 $\pm$ 2.79	56.31 $\pm$ 0.72	<b>72.65 <math>\pm</math> 0.59</b>	61.71 $\pm$ 0.63	68.97	68.72	3.75
	WiSE-FT	63.85 $\pm$ 3.69	81.81 $\pm$ 2.80	62.71 $\pm$ 1.26	77.83 $\pm$ 2.02	73.40 $\pm$ 2.08	56.63 $\pm$ 0.63	71.27 $\pm$ 0.77	<u>62.70 <math>\pm</math> 0.87</u>	68.78	68.63	4.00
	L <sup>2</sup> -SP	63.67 $\pm$ 1.79	88.00 $\pm$ 1.00	63.98 $\pm$ 1.51	77.38 $\pm$ 1.25	58.29 $\pm$ 3.74	56.23 $\pm$ 1.70	71.93 $\pm$ 0.21	59.29 $\pm$ 0.72	67.35	65.76	4.50
	FEATURE-MAP	64.41 $\pm$ 1.38	86.82 $\pm$ 0.76	59.62 $\pm$ 1.17	70.71 $\pm$ 0.99	76.01 $\pm$ 0.60	55.03 $\pm$ 0.30	67.98 $\pm$ 0.41	57.91 $\pm$ 0.31	67.31	66.11	5.25
BSS	<b>67.80 <math>\pm</math> 4.60</b>	84.90 $\pm$ 2.20	62.77 $\pm$ 0.63	<b>78.13 <math>\pm</math> 2.21</b>	74.58 $\pm$ 1.13	54.91 $\pm$ 1.34	71.40 $\pm$ 0.44	<b>63.04 <math>\pm</math> 0.35</b>	69.69	69.62	3.12	
SUPERVISED PRE-TRAINING (GRAPHIUM-LARGE)												
RANDOM	FULL-FT	81.27 $\pm$ 3.88	69.17 $\pm$ 1.32	79.75 $\pm$ 1.07	76.42 $\pm$ 0.72	76.84 $\pm$ 1.80	63.63 $\pm$ 0.06	78.12 $\pm$ 0.46	66.37 $\pm$ 0.26	73.95	74.45	3.75
	LP	80.48 $\pm$ 0.00	66.90 $\pm$ 0.00	80.44 $\pm$ 0.00	75.83 $\pm$ 0.00	73.35 $\pm$ 0.00	62.03 $\pm$ 0.00	79.02 $\pm$ 0.00	66.09 $\pm$ 0.00	73.02	73.61	5.12
	SURGICAL-FT	86.17 $\pm$ 0.00	<b>73.71 <math>\pm</math> 0.00</b>	<b>84.16 <math>\pm</math> 0.00</b>	<b>77.47 <math>\pm</math> 0.00</b>	<b>78.87 <math>\pm</math> 0.00</b>	<b>64.02 <math>\pm</math> 0.00</b>	78.23 $\pm$ 0.00	<b>67.34 <math>\pm</math> 0.00</b>	76.25	76.63	1.38
	LP-FT	83.67 $\pm$ 3.53	69.98 $\pm$ 0.83	79.28 $\pm$ 0.32	76.17 $\pm$ 2.01	77.82 $\pm$ 1.15	61.20 $\pm$ 0.00	76.94 $\pm$ 0.00	66.28 $\pm$ 0.00	73.92	74.41	4.62
	WiSE-FT	85.40 $\pm$ 1.61	71.89 $\pm$ 1.79	78.13 $\pm$ 2.92	76.69 $\pm$ 1.76	74.37 $\pm$ 1.79	63.58 $\pm$ 0.00	77.98 $\pm$ 0.33	66.48 $\pm$ 0.43	74.31	74.26	3.62
	L <sup>2</sup> -SP	76.83 $\pm$ 8.87	67.35 $\pm$ 0.82	78.17 $\pm$ 0.02	73.69 $\pm$ 0.03	62.35 $\pm$ 0.15	62.21 $\pm$ 0.45	76.27 $\pm$ 0.32	62.75 $\pm$ 0.88	69.95	69.87	6.62
	FEATURE-MAP	<b>90.13 <math>\pm</math> 2.12</b>	70.99 $\pm$ 0.27	83.17 $\pm$ 0.49	73.61 $\pm$ 0.03	78.74 $\pm$ 0.76	62.12 $\pm$ 0.02	<b>79.99 <math>\pm</math> 0.12</b>	65.03 $\pm$ 0.08	75.47	75.25	3.50
BSS	79.99 $\pm$ 5.89	67.10 $\pm$ 0.93	78.12 $\pm$ 2.32	72.50 $\pm$ 0.51	61.20 $\pm$ 0.08	61.13 $\pm$ 0.95	69.71 $\pm$ 0.64	65.45 $\pm$ 0.89	70.27	70.18	3.78	
SIZE	FULL-FT	85.96 $\pm$ 4.28	87.62 $\pm$ 0.90	67.41 $\pm$ 2.44	81.47 $\pm$ 1.94	72.03 $\pm$ 2.55	54.72 $\pm$ 0.01	69.71 $\pm$ 0.37	61.31 $\pm$ 0.37	72.53	72.98	3.88
	LP	81.84 $\pm$ 0.02	78.09 $\pm$ 0.00	58.08 $\pm$ 0.01	77.48 $\pm$ 0.00	69.46 $\pm$ 0.00	53.59 $\pm$ 0.00	73.65 $\pm$ 0.00	61.25 $\pm$ 0.00	69.18	69.67	5.38
	SURGICAL-FT	86.59 $\pm$ 0.01	<b>89.07 <math>\pm</math> 0.00</b>	70.94 $\pm$ 0.01	<b>74.47 <math>\pm</math> 0.00</b>	<b>74.47 <math>\pm</math> 0.00</b>	<b>56.24 <math>\pm</math> 0.00</b>	72.30 $\pm$ 0.00	<b>62.74 <math>\pm</math> 0.00</b>	74.36	74.92	1.62
	LP-FT	<b>86.78 <math>\pm</math> 2.69</b>	88.02 $\pm$ 1.50	63.72 $\pm$ 1.85	<b>82.57 <math>\pm</math> 0.46</b>	73.51 $\pm$ 1.77	52.40 $\pm$ 0.00	68.23 $\pm$ 0.87	60.85 $\pm$ 0.00	72.01	72.61	4.00
	WiSE-FT	82.44 $\pm$ 3.02	87.76 $\pm$ 0.5	<b>72.89 <math>\pm</math> 0.66</b>	81.37 $\pm$ 1.07	73.67 $\pm$ 3.44						

Table 10: Robust fine-tuning performance on 5 **Classification** datasets (AUC metrics) in the **Fewshot** setting (covering FEWSHOT-50, FEWSHOT-100, FEWSHOT-500), evaluated across 3 dataset splits (RANDOM, SCAFFOLD, SIZE) over **MOLECULESTM** and **GRAPHIUM-LARGE** models. We **bold** and underline the best and second-best performances in each scenario.

SPLIT	METHODS	SELF-SUPERVISED PRE-TRAINING (MOLECULESTM)										SUPERVISED PRE-TRAINING (GRAPHIUM-LARGE)									
		CLINTox	BBBP	BACE	HIV	SIDER	AVG	AVG-F	AVG-R	CLINTox	BBBP	BACE	HIV	SIDER	AVG	AVG-F	AVG-R				
RANDOM	FULL-FT	49.60 ± 2.85	84.86 ± 1.30	74.74 ± 1.44	60.58 ± 1.47	49.47 ± 0.90	63.85	61.64	4.80	74.25 ± 0.00	82.09 ± 0.77	81.04 ± 0.00	62.83 ± 0.00	52.55 ± 0.00	70.55	72.71	6.00				
	LP	52.66 ± 3.14	78.85 ± 1.75	58.02 ± 3.19	52.39 ± 0.52	50.23 ± 0.47	58.43	54.36	6.40	64.37 ± 0.00	86.23 ± 0.00	81.47 ± 0.00	60.00 ± 0.00	54.28 ± 0.00	69.27	68.61	5.00				
	SURGICAL-FT	54.43 ± 4.39	<b>86.64 ± 0.06</b>	71.92 ± 0.05	<b>61.71 ± 0.64</b>	51.10 ± 0.82	65.76	63.69	2.00	70.03 ± 0.00	87.04 ± 0.00	<b>82.44 ± 0.00</b>	62.09 ± 0.00	53.09 ± 0.00	72.14	73.52	3.60				
	LP-FT	47.71 ± 2.16	84.36 ± 2.65	<u>71.92 ± 0.05</u>	55.82 ± 1.53	<b>51.62 ± 0.37</b>	62.89	60.79	4.60	<b>76.40 ± 0.00</b>	82.10 ± 0.00	73.86 ± 0.00	64.86 ± 0.00	54.11 ± 0.00	70.27	71.71	4.00				
	WISE-FT	<b>55.69 ± 5.37</b>	84.62 ± 1.45	74.02 ± 1.36	60.05 ± 1.26	49.41 ± 0.89	64.76	63.25	4.60	75.77 ± 0.00	84.05 ± 0.85	81.30 ± 0.00	62.46 ± 0.00	<b>55.49 ± 0.00</b>	71.81	73.18	3.60				
	L <sup>2</sup> -SP	50.07 ± 2.37	<u>85.69 ± 1.19</u>	<b>75.18 ± 1.16</b>	58.44 ± 1.98	50.58 ± 0.93	63.99	61.40	3.60	75.31 ± 2.24	84.45 ± 4.02	80.56 ± 0.00	<b>69.92 ± 0.00</b>	53.87 ± 0.00	72.82	75.26	4.60				
	FEATURE-MAP	54.09 ± 3.21	<u>78.77 ± 1.05</u>	67.88 ± 0.54	55.43 ± 1.21	50.12 ± 0.27	61.26	59.13	6.20	71.01 ± 0.00	<b>88.81 ± 0.00</b>	81.76 ± 0.03	61.15 ± 0.00	54.47 ± 0.15	71.44	71.31	3.80				
	BSS	52.06 ± 3.58	85.62 ± 1.18	74.31 ± 1.83	58.90 ± 0.76	<u>51.18 ± 0.69</u>	64.41	61.76	3.80	75.33 ± 0.00	81.30 ± 1.08	80.08 ± 0.00	64.67 ± 0.00	53.88 ± 0.51	71.23	73.66	5.20				
	FULL-FT	45.62 ± 5.48	<b>58.05 ± 2.70</b>	62.30 ± 1.27	<u>48.87 ± 6.91</u>	54.88 ± 0.29	53.94	53.93	2.60	74.79 ± 0.00	61.10 ± 0.00	74.43 ± 0.00	64.93 ± 0.00	54.35 ± 0.00	65.92	66.82	5.60				
	LP	30.76 ± 1.34	50.50 ± 1.35	56.94 ± 2.34	38.17 ± 1.21	53.17 ± 0.36	46.11	47.62	7.80	67.24 ± 0.00	63.31 ± 0.00	65.24 ± 0.00	50.89 ± 0.00	55.24 ± 0.00	60.38	61.60	5.60				
SCAFFOLD	SURGICAL-FT	45.90 ± 9.96	56.02 ± 1.54	63.07 ± 0.78	44.00 ± 3.78	<u>55.18 ± 0.47</u>	52.77	52.27	3.80	71.74 ± 0.00	62.43 ± 0.00	74.64 ± 0.00	65.60 ± 0.00	55.55 ± 0.00	65.99	66.99	4.00				
	LP-FT	33.97 ± 3.65	55.31 ± 2.06	61.87 ± 0.80	45.88 ± 1.92	<u>55.16 ± 0.46</u>	50.44	52.12	5.20	61.66 ± 0.00	63.39 ± 0.00	<b>76.82 ± 0.00</b>	56.11 ± 0.00	<u>56.50 ± 0.00</u>	62.90	60.52	4.60				
	WISE-FT	<b>47.69 ± 5.22</b>	<u>57.80 ± 2.92</u>	62.06 ± 1.03	47.33 ± 3.84	55.16 ± 0.57	54.01	53.55	2.60	73.93 ± 0.00	<b>65.16 ± 0.00</b>	<u>71.82 ± 0.00</u>	64.36 ± 0.00	<u>54.92 ± 0.00</u>	66.64	67.82	3.60				
	L <sup>2</sup> -SP	45.54 ± 5.40	56.06 ± 1.99	61.75 ± 1.66	45.56 ± 4.10	<b>55.29 ± 0.92</b>	52.84	52.30	4.20	68.43 ± 0.00	64.01 ± 0.93	74.63 ± 0.00	66.45 ± 0.00	<b>56.54 ± 0.00</b>	66.01	66.30	3.20				
	FEATURE-MAP	26.69 ± 2.38	56.71 ± 1.18	61.18 ± 5.30	43.71 ± 3.23	53.77 ± 0.39	48.41	51.40	6.60	65.60 ± 0.03	63.73 ± 0.00	70.32 ± 0.00	<b>70.97 ± 0.00</b>	54.72 ± 0.03	65.07	66.55	5.20				
	BSS	42.19 ± 1.78	57.09 ± 1.32	<b>63.74 ± 2.79</b>	<b>50.07 ± 8.79</b>	54.75 ± 0.37	53.57	53.97	3.20	<b>77.89 ± 0.04</b>	61.79 ± 0.00	74.27 ± 1.63	<u>62.56 ± 0.00</u>	55.03 ± 0.01	67.11	67.54	4.20				
	FULL-FT	58.52 ± 2.98	58.80 ± 9.95	36.17 ± 0.29	<u>52.04 ± 2.74</u>	<u>51.97 ± 1.34</u>	51.50	54.18	4.20	71.15 ± 0.00	80.00 ± 0.00	59.96 ± 3.09	48.05 ± 0.00	53.20 ± 0.00	62.47	61.44	4.60				
	LP	37.53 ± 4.82	45.54 ± 17.14	47.39 ± 1.62	48.21 ± 0.61	<u>50.89 ± 0.73</u>	49.91	48.83	6.60	62.05 ± 0.00	72.11 ± 0.00	56.89 ± 0.01	57.63 ± 0.00	49.15 ± 0.00	59.57	58.86	7.20				
	SURGICAL-FT	61.32 ± 8.19	54.19 ± 11.51	44.96 ± 7.70	51.79 ± 2.35	51.41 ± 0.98	52.73	52.46	4.80	<u>71.68 ± 0.00</u>	83.99 ± 0.00	<b>62.17 ± 0.00</b>	62.00 ± 0.00	54.99 ± 0.00	66.97	65.28	2.40				
	LP-FT	54.70 ± 9.04	55.56 ± 3.73	43.08 ± 1.91	47.90 ± 2.39	51.88 ± 0.55	50.62	51.49	5.80	70.52 ± 0.00	79.54 ± 0.00	59.30 ± 0.00	56.07 ± 0.00	52.10 ± 0.00	63.92	62.66	5.00				
RANDOM	WISE-FT	<u>61.69 ± 5.18</u>	56.83 ± 9.47	42.48 ± 6.40	50.61 ± 2.71	<b>52.28 ± 1.23</b>	52.76	53.24	3.80	70.51 ± 0.00	78.10 ± 0.00	59.48 ± 3.21	54.15 ± 0.00	53.24 ± 0.00	63.10	61.38	5.20				
	L <sup>2</sup> -SP	<u>60.54 ± 2.21</u>	<b>62.77 ± 6.52</b>	47.51 ± 8.30	<b>52.06 ± 2.80</b>	51.52 ± 1.67	54.88	54.71	2.60	65.70 ± 0.03	<b>85.88 ± 0.76</b>	56.81 ± 0.04	<u>62.79 ± 0.06</u>	<b>57.10 ± 0.00</b>	65.66	61.86	3.80				
	FEATURE-MAP	50.85 ± 1.96	50.21 ± 1.87	47.05 ± 3.15	44.09 ± 1.27	51.48 ± 0.50	50.66	49.78	5.40	69.15 ± 0.01	85.63 ± 0.31	61.95 ± 0.58	<b>64.82 ± 0.03</b>	50.81 ± 0.01	66.48	65.31	3.60				
	BSS	<b>62.26 ± 1.89</b>	60.79 ± 7.04	<b>49.70 ± 2.37</b>	51.85 ± 3.42	51.19 ± 1.56	55.16	54.61	2.80	<b>73.63 ± 0.01</b>	<u>72.93 ± 2.44</u>	56.91 ± 3.73	52.67 ± 1.33	<u>56.22 ± 0.72</u>	63.87	62.25	4.20				
	FULL-FT	<b>73.60 ± 7.53</b>	82.09 ± 2.90	80.72 ± 1.22	61.92 ± 2.62	51.58 ± 0.43	69.98	72.08	5.00	66.36 ± 0.01	86.40 ± 2.10	78.44 ± 0.00	63.35 ± 0.00	56.74 ± 0.00	70.26	69.38	6.20				
	LP	69.43 ± 1.40	73.63 ± 0.97	60.60 ± 3.89	54.74 ± 0.90	53.47 ± 0.21	62.37	61.59	6.60	65.67 ± 0.00	<b>90.26 ± 0.00</b>	81.88 ± 0.00	61.87 ± 0.00	57.00 ± 0.00	71.34	69.81	5.20				
	SURGICAL-FT	71.20 ± 2.70	83.50 ± 0.95	80.44 ± 0.62	62.65 ± 1.44	53.43 ± 0.50	70.24	71.45	4.20	71.48 ± 0.00	86.23 ± 0.00	85.03 ± 0.00	63.49 ± 0.00	<b>58.52 ± 0.00</b>	72.95	73.53	3.20				
	LP-FT	68.16 ± 1.86	<b>84.26 ± 1.37</b>	79.93 ± 2.67	60.14 ± 3.04	52.18 ± 0.81	68.93	64.41	5.20	70.77 ± 0.00	<u>89.94 ± 0.00</u>	<u>77.87 ± 2.04</u>	61.52 ± 0.00	57.76 ± 0.00	71.57	70.05	5.20				
	WISE-FT	72.72 ± 8.35	<u>83.52 ± 3.24</u>	<b>88.26 ± 1.45</b>	62.19 ± 2.74	51.66 ± 0.43	71.67	72.81	3.80	68.92 ± 0.01	86.48 ± 0.54	79.32 ± 0.00	63.14 ± 0.00	56.58 ± 0.00	70.89	70.46	6.00				
	L <sup>2</sup> -SP	73.05 ± 2.80	<u>82.49 ± 1.55</u>	81.60 ± 1.23	63.21 ± 2.21	<u>53.92 ± 0.82</u>	70.85	73.62	3.00	<u>74.74 ± 1.31</u>	86.30 ± 2.30	81.62 ± 1.01	63.65 ± 0.00	57.84 ± 0.00	72.81	73.44	4.00				
FEATURE-MAP	68.01 ± 2.06	78.35 ± 0.58	69.																		

Table 11: Robust fine-tuning performance on 4 **Regression** datasets (RMSE metrics) in the **Fewshot** setting (covering FEWSHOT-50, FEWSHOT-100, and FEWSHOT-500), evaluated across 3 dataset splits (RANDOM, SCAFFOLD, SIZE) over **MOLECULESTM** and **GRAPHIUM-LARGE** models. AVG-R, AVG-R\* denote the average rank and the rank based on the average normalized performance over all the datasets for each evaluated method, respectively. Standard deviations across five replicates are shown in parentheses. We **bold** and underline the best and second-best performances in each scenario.

SPLIT	METHODS	SELF-SUPERVISED PRE-TRAINING (MOLECULESTM)					SUPERVISED PRE-TRAINING (GRAPHIUM-LARGE)						
		ESOL	LIPO	MALARIA	CEP	AVG-R	AVG-R*	ESOL	LIPO	MALARIA	CEP	AVG-R*	
RANDOM	FEWSHOT-50												
	FULL-FT	2.128 ± 0.072	1.247 ± 0.031	1.310 ± 0.025	3.433 ± 0.226	5.00	6	1.125 ± 0.000	1.156 ± 0.019	1.277 ± 0.000	2.198 ± 0.001	7	
	LP	2.971 ± 0.017	1.638 ± 0.014	1.309 ± 0.012	3.519 ± 0.052	6.75	8	1.176 ± 0.000	1.131 ± 0.000	1.294 ± 0.000	2.113 ± 0.000	8	
	SURGICAL-FT	2.315 ± 0.081	1.327 ± 0.017	1.317 ± 0.024	3.272 ± 0.199	6.50	7	1.055 ± 0.000	1.076 ± 0.000	1.283 ± 0.000	2.192 ± 0.000	4	
	LP-FT	1.600 ± 0.129	1.181 ± 0.030	1.356 ± 0.011	2.358 ± 0.037	4.25	4	1.096 ± 0.000	1.032 ± 0.002	1.293 ± 0.000	2.092 ± 0.002	1	
	WISE-FT	2.135 ± 0.072	1.261 ± 0.035	1.298 ± 0.023	3.576 ± 0.235	5.50	5	1.116 ± 0.000	1.151 ± 0.024	1.278 ± 0.000	2.075 ± 0.004	3	
	L <sup>2</sup> -SP	1.472 ± 0.036	1.165 ± 0.037	1.297 ± 0.006	2.304 ± 0.055	1.50	1	1.161 ± 0.000	1.077 ± 0.019	1.276 ± 0.000	2.127 ± 0.015	5	
	FEATURE-MAP	1.632 ± 0.028	1.257 ± 0.025	1.301 ± 0.009	2.398 ± 0.037	4.00	3	1.133 ± 0.002	1.106 ± 0.003	1.277 ± 0.001	2.108 ± 0.002	2	
	BSS	1.450 ± 0.057	1.171 ± 0.021	1.314 ± 0.018	2.244 ± 0.036	2.50	2	1.188 ± 0.004	1.109 ± 0.021	1.276 ± 0.000	2.108 ± 0.029	6	
	SCAFFOLD	FULL-FT	2.790 ± 0.116	1.434 ± 0.072	1.195 ± 0.025	3.395 ± 0.191	5.75	6	1.237 ± 0.000	1.079 ± 0.000	1.175 ± 0.000	2.051 ± 0.000	7
LP		3.538 ± 0.075	1.755 ± 0.021	1.206 ± 0.012	3.870 ± 0.038	7.75	8	0.929 ± 0.000	1.096 ± 0.000	1.170 ± 0.000	2.053 ± 0.000	1	
SURGICAL-FT		3.018 ± 0.118	1.491 ± 0.085	1.191 ± 0.004	3.304 ± 0.347	5.75	7	1.240 ± 0.000	1.044 ± 0.000	1.180 ± 0.000	2.009 ± 0.000	2	
LP-FT		1.636 ± 0.021	1.181 ± 0.029	1.263 ± 0.009	2.294 ± 0.024	4.00	4	1.241 ± 0.000	1.085 ± 0.000	1.176 ± 0.000	2.044 ± 0.000	8	
WISE-FT		3.527 ± 0.112	1.392 ± 0.062	0.983 ± 0.053	3.435 ± 0.142	5.00	5	1.536 ± 0.000	1.149 ± 0.000	0.911 ± 0.000	2.321 ± 0.000	5	
L <sup>2</sup> -SP		1.654 ± 0.086	1.178 ± 0.022	1.185 ± 0.008	2.255 ± 0.026	2.25	2	1.280 ± 0.003	1.107 ± 0.002	1.175 ± 0.000	1.997 ± 0.016	6	
FEATURE-MAP		1.783 ± 0.034	1.252 ± 0.012	1.195 ± 0.008	2.401 ± 0.028	4.50	3	1.267 ± 0.110	1.037 ± 0.006	1.170 ± 0.143	2.073 ± 0.016	5	
BSS		1.632 ± 0.048	1.173 ± 0.022	1.182 ± 0.016	2.287 ± 0.028	1.50	1	1.159 ± 0.007	1.100 ± 0.002	1.162 ± 0.000	2.060 ± 0.009	3	
SIZE		FULL-FT	3.457 ± 0.086	1.407 ± 0.088	1.064 ± 0.067	3.311 ± 0.158	6.25	7	1.499 ± 0.000	1.108 ± 0.000	0.909 ± 0.000	2.321 ± 0.000	4
		LP	3.758 ± 0.010	1.773 ± 0.025	0.990 ± 0.056	4.114 ± 0.042	6.75	8	2.025 ± 0.000	1.325 ± 0.000	0.917 ± 0.000	2.358 ± 0.000	8
	SURGICAL-FT	3.429 ± 0.139	1.543 ± 0.083	0.990 ± 0.054	3.195 ± 0.306	5.25	6	1.675 ± 0.000	1.089 ± 0.000	0.916 ± 0.000	2.271 ± 0.000	1	
	LP-FT	2.035 ± 0.080	1.208 ± 0.078	1.102 ± 0.018	2.500 ± 0.045	4.00	4	1.540 ± 0.000	1.079 ± 0.001	0.994 ± 0.000	2.347 ± 0.001	7	
	WISE-FT	3.527 ± 0.112	1.392 ± 0.062	0.983 ± 0.053	3.386 ± 0.142	5.00	5	1.536 ± 0.000	1.149 ± 0.000	0.911 ± 0.000	2.321 ± 0.000	5	
	L <sup>2</sup> -SP	2.111 ± 0.091	1.159 ± 0.037	0.988 ± 0.032	2.421 ± 0.045	2.00	1	1.673 ± 0.030	1.072 ± 0.002	0.948 ± 0.007	2.304 ± 0.022	6	
	FEATURE-MAP	2.331 ± 0.050	1.225 ± 0.049	1.000 ± 0.034	2.439 ± 0.024	4.00	3	1.594 ± 0.010	1.070 ± 0.012	0.915 ± 0.001	2.306 ± 0.008	3	
	BSS	2.197 ± 0.084	1.106 ± 0.027	1.019 ± 0.033	2.419 ± 0.045	2.75	2	1.516 ± 0.008	1.076 ± 0.043	0.907 ± 0.000	2.313 ± 0.049	2	
	FEWSHOT-100												
	RANDOM	FULL-FT	1.842 ± 0.208	1.205 ± 0.059	1.289 ± 0.032	2.784 ± 0.110	5.75	6	1.121 ± 0.000	1.187 ± 0.020	1.259 ± 0.000	1.902 ± 0.011	6
LP		2.391 ± 0.044	1.623 ± 0.011	1.279 ± 0.007	3.176 ± 0.093	7.00	8	0.912 ± 0.000	1.068 ± 0.000	1.286 ± 0.000	1.920 ± 0.014	4	
SURGICAL-FT		1.650 ± 0.063	1.301 ± 0.037	1.277 ± 0.012	2.777 ± 0.181	5.00	4	0.952 ± 0.000	1.061 ± 0.000	1.269 ± 0.000	1.881 ± 0.000	2	
LP-FT		1.540 ± 0.123	1.234 ± 0.030	1.350 ± 0.016	2.203 ± 0.030	4.50	7	1.061 ± 0.005	1.126 ± 0.000	1.290 ± 0.011	1.918 ± 0.005	7	
WISE-FT		1.790 ± 0.147	1.207 ± 0.058	1.282 ± 0.017	2.842 ± 0.123	5.50	5	1.064 ± 0.000	1.121 ± 0.050	1.258 ± 0.000	1.905 ± 0.015	3	
L <sup>2</sup> -SP		1.486 ± 0.105	1.190 ± 0.038	1.267 ± 0.007	2.207 ± 0.046	1.75	1	1.109 ± 0.082	1.094 ± 0.007	1.276 ± 0.000	1.916 ± 0.022	5	
FEATURE-MAP		1.557 ± 0.034	1.252 ± 0.007	1.269 ± 0.002	2.130 ± 0.020	3.25	2	0.897 ± 0.009	1.053 ± 0.007	1.273 ± 0.000	1.881 ± 0.011	1	
BSS		1.543 ± 0.044	1.190 ± 0.031	1.285 ± 0.011	2.170 ± 0.028	3.25	3	1.159 ± 0.012	1.129 ± 0.022	1.276 ± 0.004	2.036 ± 0.139	8	
SCAFFOLD		FULL-FT	2.036 ± 0.119	1.108 ± 0.017	1.205 ± 0.050	2.942 ± 0.208	5.75	6	1.238 ± 0.000	1.027 ± 0.000	1.187 ± 0.000	1.986 ± 0.019	7
		LP	2.906 ± 0.093	1.389 ± 0.008	1.180 ± 0.017	3.635 ± 0.051	6.75	8	1.184 ± 0.013	0.998 ± 0.000	1.163 ± 0.000	1.935 ± 0.000	3
	SURGICAL-FT	1.956 ± 0.170	1.190 ± 0.027	1.183 ± 0.016	2.848 ± 0.120	5.50	5	1.121 ± 0.000	0.977 ± 0.000	1.172 ± 0.000	1.914 ± 0.000	1	
	LP-FT	1.775 ± 0.178	1.103 ± 0.024	1.288 ± 0.012	2.310 ± 0.034	4.75	7	1.210 ± 0.001	1.062 ± 0.003	1.206 ± 0.000	1.918 ± 0.002	8	
	WISE-FT	2.052 ± 0.082	1.112 ± 0.023	1.188 ± 0.027	3.049 ± 0.246	6.25	4	1.199 ± 0.000	1.002 ± 0.000	1.160 ± 0.000	1.988 ± 0.028	5	
	L <sup>2</sup> -SP	1.559 ± 0.047	1.069 ± 0.044	1.166 ± 0.004	2.227 ± 0.036	1.75	1	1.210 ± 0.030	0.999 ± 0.035	1.176 ± 0.015	2.000 ± 0.000	6	
	FEATURE-MAP	1.576 ± 0.028	1.123 ± 0.009	1.181 ± 0.005	2.216 ± 0.014	3.50	3	1.106 ± 0.025	0.957 ± 0.008	1.159 ± 0.003	2.047 ± 0.008	2	
	BSS	1.680 ± 0.098	1.081 ± 0.019	1.163 ± 0.004	2.212 ± 0.018	1.75	2	1.169 ± 0.035	1.025 ± 0.000	1.170 ± 0.014	1.938 ± 0.030	4	
	SIZE	FULL-FT	2.527 ± 0.152	1.113 ± 0.054	1.022 ± 0.046	2.587 ± 0.100	6.25	7	1.675 ± 0.003	1.132 ± 0.000	0.909 ± 0.000	2.317 ± 0.000	6
		LP	3.020 ± 0.061	1.492 ± 0.039	0.951 ± 0.011	3.408 ± 0.041	6.75	8	1.740 ± 0.000	1.245 ± 0.000	0.934 ± 0.000	2.355 ± 0.000	8
SURGICAL-FT		2.435 ± 0.119	1.119 ± 0.037	0.970 ± 0.020	2.607 ± 0.040	6.25							

Table 12: Robust fine-tuning performance on 8 **Classification** datasets (AUC metrics) in the **Non-Fewshot** setting, evaluated across 3 dataset splits (RANDOM, SCAFFOLD, SIZE), over **GRAPHMAE** and **GRAPHGPS** models. AVG, AVG-F, AVG-R denote the average AUC, average AUC without max and min values, and average rank over all the datasets for each method, respectively. Standard deviations across five replicates are shown in parentheses. We **bold** and underline the best and second-best performances in each scenario.

SPLIT	METHODS	CLINTOX	BBBP	BACE	HIV	MUV	SIDER	Tox21	ToxCast	AVG	AVG-F	AVG-R
SELF-SUPERVISED PRE-TRAINING (GRAPHMAE)												
RANDOM	FULL-FT	83.22 ± 2.07	94.70 ± 0.32	89.26 ± 0.40	85.31 ± 0.29	80.71 ± 0.58	61.53 ± 0.48	82.35 ± 0.15	73.01 ± 0.16	81.26	82.31	4.00
	LP	78.82 ± 1.55	83.16 ± 0.58	77.65 ± 1.27	74.45 ± 0.31	78.54 ± 1.16	61.51 ± 0.35	73.57 ± 0.16	66.96 ± 0.16	74.33	75.00	7.50
	SURGICAL-FT	83.85 ± 1.52	92.11 ± 0.35	86.77 ± 0.09	84.56 ± 0.30	<b>82.71 ± 0.81</b>	61.79 ± 0.19	79.90 ± 0.14	71.51 ± 0.21	80.40	81.55	4.50
	LP-FT	<b>88.09 ± 1.04</b>	94.68 ± 0.19	89.58 ± 0.23	<b>86.06 ± 0.43</b>	80.75 ± 1.53	61.69 ± 0.26	<b>82.50 ± 0.21</b>	<b>73.66 ± 0.07</b>	82.13	83.44	2.25
	WISE-FT	80.01 ± 4.00	93.04 ± 0.46	<b>90.15 ± 0.50</b>	85.42 ± 0.52	82.07 ± 2.10	62.18 ± 0.49	81.55 ± 0.43	72.48 ± 0.26	80.86	81.95	3.38
	L2-SP	83.39 ± 1.88	93.89 ± 0.28	88.70 ± 0.10	80.22 ± 0.17	73.35 ± 1.54	<b>62.36 ± 0.43</b>	77.45 ± 0.47	68.71 ± 0.31	78.51	78.64	5.00
SCAFFOLD	FEATURE-MAP	73.08 ± 0.89	85.36 ± 0.46	75.88 ± 0.75	77.04 ± 0.26	79.53 ± 1.25	62.06 ± 0.32	75.36 ± 0.13	65.69 ± 0.24	74.25	74.43	6.75
	BSS	<u>83.98 ± 3.00</u>	<b>94.85 ± 0.31</b>	89.31 ± 0.21	<u>86.05 ± 0.40</u>	80.55 ± 0.75	61.92 ± 0.21	<u>82.48 ± 0.28</u>	<u>73.22 ± 0.07</u>	81.54	82.60	2.62
	FULL-FT	74.74 ± 0.93	66.35 ± 0.65	80.33 ± 0.37	<u>77.22 ± 0.38</u>	77.47 ± 1.33	60.98 ± 0.19	<u>76.18 ± 0.31</u>	<u>64.27 ± 0.36</u>	72.19	72.70	3.88
	LP	71.34 ± 1.48	64.36 ± 0.24	61.70 ± 7.34	70.62 ± 0.64	<u>79.13 ± 1.20</u>	58.23 ± 1.29	70.89 ± 0.10	60.03 ± 0.13	67.04	66.49	6.75
	SURGICAL-FT	71.88 ± 1.07	66.81 ± 0.29	80.24 ± 0.90	76.90 ± 0.30	<b>79.20 ± 0.50</b>	<b>64.00 ± 0.09</b>	74.18 ± 0.40	62.60 ± 0.27	71.98	72.16	4.12
	LP-FT	74.88 ± 2.00	67.39 ± 0.55	80.67 ± 0.57	<b>77.97 ± 0.38</b>	75.13 ± 1.06	60.76 ± 0.45	<u>76.18 ± 0.20</u>	<b>64.29 ± 0.23</b>	72.16	72.64	3.25
SIZE	WISE-FT	<b>77.30 ± 5.30</b>	<b>69.29 ± 2.34</b>	<b>82.16 ± 1.50</b>	76.75 ± 0.69	77.76 ± 1.55	59.76 ± 0.86	74.99 ± 0.44	63.61 ± 0.34	72.70	73.28	3.25
	L2-SP	73.40 ± 0.45	67.39 ± 0.90	80.36 ± 0.92	74.63 ± 0.44	73.20 ± 0.90	63.40 ± 0.29	73.16 ± 0.14	61.29 ± 0.38	70.85	70.86	5.00
	FEATURE-MAP	64.74 ± 0.62	62.46 ± 0.26	69.22 ± 2.06	72.34 ± 0.58	75.63 ± 0.54	57.13 ± 1.08	71.25 ± 0.13	57.78 ± 0.26	66.32	66.30	7.38
	BSS	<u>75.80 ± 1.11</u>	<u>67.46 ± 1.35</u>	<u>80.82 ± 0.62</u>	77.10 ± 0.77	78.53 ± 1.03	62.29 ± 0.51	<b>76.45 ± 0.24</b>	64.03 ± 0.09	72.81	73.23	2.38
	FULL-FT	56.52 ± 0.81	80.05 ± 2.01	59.94 ± 1.83	77.21 ± 0.94	74.64 ± 1.72	53.04 ± 0.74	70.87 ± 0.24	60.80 ± 0.50	66.63	66.66	4.62
	LP	57.44 ± 0.94	73.52 ± 1.68	51.46 ± 0.97	73.91 ± 0.89	65.97 ± 3.36	51.84 ± 0.31	67.56 ± 0.10	57.49 ± 0.11	62.40	62.30	7.38
RANDOM	SURGICAL-FT	57.47 ± 1.16	81.96 ± 0.78	55.85 ± 2.81	<b>80.48 ± 0.18</b>	75.86 ± 2.96	54.32 ± 0.43	71.19 ± 0.30	59.45 ± 0.18	67.07	66.72	3.12
	LP-FT	56.35 ± 0.62	76.80 ± 2.24	61.61 ± 1.01	77.14 ± 0.69	<b>79.10 ± 0.89</b>	52.69 ± 0.35	<b>71.33 ± 0.26</b>	60.98 ± 0.27	67.00	67.37	4.00
	WISE-FT	<b>59.25 ± 3.49</b>	<b>82.99 ± 1.91</b>	61.16 ± 2.31	75.90 ± 1.94	75.09 ± 3.95	<b>55.74 ± 1.28</b>	70.94 ± 0.42	<b>61.53 ± 0.56</b>	67.83	67.31	2.50
	L2-SP	59.11 ± 0.88	80.40 ± 1.50	61.10 ± 1.52	70.67 ± 1.61	65.11 ± 0.75	53.81 ± 0.72	68.96 ± 0.47	57.85 ± 0.36	65.38	64.80	4.88
	FEATURE-MAP	59.02 ± 0.89	77.60 ± 0.45	43.17 ± 0.32	79.17 ± 0.23	73.54 ± 0.29	52.23 ± 0.32	68.74 ± 0.09	53.39 ± 0.51	63.36	64.09	5.75
	BSS	58.58 ± 1.31	80.86 ± 1.92	<b>61.96 ± 2.00</b>	79.14 ± 0.79	73.35 ± 1.27	53.14 ± 0.63	70.76 ± 0.37	60.62 ± 0.35	67.30	67.40	3.75
SUPERVISED PRE-TRAINING (GRAPHGPS)												
RANDOM	FULL-FT	99.77 ± 0.01	99.99 ± 0.01	<b>100.00 ± 0.00</b>	84.80 ± 0.33	57.06 ± 0.00	87.13 ± 0.39	87.17 ± 0.48	86.90 ± 0.17	87.85	90.96	4.00
	LP	99.48 ± 0.04	86.96 ± 0.40	80.94 ± 0.45	86.70 ± 0.42	<u>63.97 ± 0.80</u>	84.77 ± 0.08	82.70 ± 0.14	83.93 ± 0.04	83.68	84.33	5.50
	SURGICAL-FT	99.65 ± 0.05	99.16 ± 0.00	98.14 ± 0.04	86.58 ± 0.03	60.52 ± 0.64	47.74 ± 0.95	51.53 ± 0.00	51.71 ± 0.00	74.38	74.61	5.88
	LP-FT	99.54 ± 0.14	89.67 ± 5.14	84.88 ± 7.57	85.71 ± 1.11	63.96 ± 0.80	85.97 ± 2.43	83.98 ± 2.45	84.48 ± 1.09	84.77	85.78	5.12
	WISE-FT	97.04 ± 1.00	58.14 ± 3.98	68.29 ± 2.24	67.14 ± 4.36	49.94 ± 0.01	80.52 ± 0.07	67.81 ± 0.12	77.50 ± 0.03	70.80	69.90	7.62
	L2-SP	<b>99.84 ± 0.03</b>	<b>100.00 ± 0.00</b>	<b>100.00 ± 0.00</b>	<b>97.75 ± 0.07</b>	<b>74.51 ± 1.12</b>	<b>92.16 ± 0.44</b>	<b>92.28 ± 0.46</b>	<b>89.79 ± 0.07</b>	93.29	95.30	1.25
SCAFFOLD	FEATURE-MAP	99.79 ± 0.09	100.00 ± 0.00	100.00 ± 0.00	<b>99.42 ± 0.01</b>	53.07 ± 0.82	91.64 ± 0.06	91.61 ± 0.16	89.39 ± 0.06	90.62	95.31	2.62
	BSS	99.77 ± 0.00	100.00 ± 0.00	100.00 ± 0.00	84.87 ± 0.02	58.93 ± 3.25	87.61 ± 0.05	87.52 ± 0.10	86.75 ± 0.05	88.18	91.09	4.00
	FULL-FT	99.76 ± 0.04	99.99 ± 0.01	<b>100.00 ± 0.00</b>	83.67 ± 1.61	57.08 ± 1.77	87.26 ± 0.15	87.16 ± 0.21	86.71 ± 0.12	87.70	90.76	4.12
	LP	99.47 ± 0.04	86.84 ± 0.49	81.04 ± 0.53	86.66 ± 0.44	<u>63.98 ± 0.82</u>	84.74 ± 0.08	82.70 ± 0.14	83.93 ± 0.04	83.67	84.32	5.75
	SURGICAL-FT	99.64 ± 0.08	99.33 ± 0.14	98.14 ± 0.06	87.61 ± 0.63	61.75 ± 0.39	76.46 ± 1.75	72.53 ± 1.99	55.58 ± 0.35	81.38	82.64	5.75
	LP-FT	99.54 ± 0.15	89.53 ± 5.24	84.35 ± 6.50	84.81 ± 2.36	62.46 ± 1.48	85.96 ± 2.47	83.96 ± 2.42	84.52 ± 1.17	84.98	85.52	5.38
SIZE	WISE-FT	97.32 ± 0.16	64.50 ± 3.69	100.00 ± 0.00	67.98 ± 4.58	49.84 ± 0.72	80.53 ± 0.07	68.04 ± 0.17	77.53 ± 0.02	75.73	76.00	7.00
	L2-SP	99.83 ± 0.03	<b>100.00 ± 0.00</b>	100.00 ± 0.00	98.35 ± 0.43	<b>74.63 ± 0.95</b>	<b>92.33 ± 0.21</b>	<b>92.43 ± 0.34</b>	<b>89.85 ± 0.17</b>	93.43	95.46	1.50
	FEATURE-MAP	<b>99.85 ± 0.01</b>	<b>100.00 ± 0.00</b>	100.00 ± 0.00	<b>99.26 ± 0.13</b>	55.32 ± 0.31	91.63 ± 0.04	91.61 ± 0.11	89.30 ± 0.06	90.87	95.27	2.62
	BSS	99.81 ± 0.04	99.99 ± 0.01	100.00 ± 0.00	85.03 ± 0.57	60.82 ± 4.94	89.80 ± 3.20	87.36 ± 0.09	86.85 ± 0.12	88.71	91.47	3.88
	FULL-FT	99.76 ± 0.03	99.99 ± 0.01	<b>100.00 ± 0.00</b>	83.42 ± 1.75	56.61 ± 1.51	87.41 ± 0.51	87.06 ± 0.10	86.90 ± 0.13	87.64	90.76	4.12
	LP	99.47 ± 0.05	86.56 ± 0.34	80.81 ± 0.52	86.66 ± 0.44	64.02 ± 0.78	84.74 ± 0.08	82.38 ± 0				

Table 13: Robust fine-tuning performance on 4 **Regression** datasets (RMSE metrics) in the **Non-Fewshot** setting, evaluated across 3 dataset splits (RANDOM, SCAFFOLD, SIZE) over **GRAPHMAE** and **GRAPHGPS** models. AVG-R, AVG-R\* denote the average rank and the rank based on the average normalized performance over all the datasets for each method, respectively. Standard deviations across five replicates are shown in parentheses. We **bold** and underline the best and second-best performances in each scenario.

SPLIT	METHODS	SELF-SUPERVISED PRE-TRAINING (GRAPHMAE)					SUPERVISED PRE-TRAINING (GRAPHGPS)						
		ESOL	LIPO	MALARIA	CEP	AVG-R	AVG-R*	ESOL	LIPO	MALARIA	CEP	AVG-R	AVG-R*
RANDOM	FULL-FT	0.987 ± 0.013	0.734 ± 0.007	1.109 ± 0.015	1.342 ± 0.015	3.00	3	0.191 ± 0.019	0.211 ± 0.012	0.955 ± 0.008	0.587 ± 0.000	4.50	4
	LP	1.394 ± 0.012	1.156 ± 0.001	1.263 ± 0.002	3.079 ± 0.105	8.00	8	0.737 ± 0.005	0.877 ± 0.004	1.031 ± 0.003	1.602 ± 0.006	6.50	6
	SURGICAL-FT	1.088 ± 0.011	0.883 ± 0.007	1.120 ± 0.012	1.697 ± 0.012	6.25	6	1.565 ± 0.313	2.284 ± 0.179	0.800 ± 0.022	0.881 ± 0.000	6.00	7
	LP-FT	<b>0.953 ± 0.009</b>	0.743 ± 0.006	<u>1.006 ± 0.009</u>	<b>1.322 ± 0.025</b>	1.75	1	<b>0.139 ± 0.016</b>	0.197 ± 0.003	0.925 ± 0.007	0.646 ± 0.087	3.25	3
	WISE-FT	1.210 ± 0.032	0.846 ± 0.023	<b>1.060 ± 0.008</b>	1.531 ± 0.030	4.50	5	2.488 ± 0.137	1.224 ± 0.007	1.187 ± 0.001	2.574 ± 0.015	7.75	8
	L2-SP	0.995 ± 0.024	0.787 ± 0.008	1.115 ± 0.006	1.363 ± 0.040	4.25	4	0.169 ± 0.009	0.194 ± 0.010	0.559 ± 0.022	0.451 ± 0.036	2.00	2
SCAFFOLD	FEATURE-MAP	1.297 ± 0.007	1.080 ± 0.002	1.115 ± 0.016	1.473 ± 0.018	6.25	7	<b>0.187 ± 0.026</b>	<b>0.134 ± 0.008</b>	<b>0.243 ± 0.009</b>	<b>0.215 ± 0.026</b>	1.75	1
	BSS	<u>0.975 ± 0.019</u>	<b>0.725 ± 0.011</b>	1.100 ± 0.004	<u>1.334 ± 0.004</u>	2.00	2	<u>0.177 ± 0.013</u>	0.213 ± 0.005	0.921 ± 0.013	0.651 ± 0.079	4.25	5
	FULL-FT	1.332 ± 0.015	0.808 ± 0.008	1.104 ± 0.007	1.327 ± 0.017	3.50	3	0.218 ± 0.054	0.202 ± 0.022	0.929 ± 0.011	0.528 ± 0.123	4.25	4
	LP	1.793 ± 0.016	1.043 ± 0.006	1.150 ± 0.003	3.102 ± 0.136	7.50	8	0.752 ± 0.006	0.849 ± 0.005	1.008 ± 0.000	1.539 ± 0.009	7.75	7
	SURGICAL-FT	1.335 ± 0.005	0.884 ± 0.007	1.111 ± 0.013	1.669 ± 0.022	5.50	5	1.574 ± 0.314	0.362 ± 0.013	0.818 ± 0.007	0.917 ± 0.000	5.50	6
	LP-FT	<b>1.312 ± 0.024</b>	<b>0.788 ± 0.005</b>	1.104 ± 0.006	<u>1.318 ± 0.017</u>	1.75	1	<b>0.145 ± 0.020</b>	<u>0.181 ± 0.012</u>	0.944 ± 0.015	0.585 ± 0.036	3.25	3
SIZE	WISE-FT	1.617 ± 0.031	0.891 ± 0.009	<b>1.077 ± 0.004</b>	1.498 ± 0.034	5.50	7	2.338 ± 0.519	1.262 ± 0.015	1.220 ± 0.017	2.610 ± 0.082	8.00	8
	L2-SP	1.329 ± 0.030	0.835 ± 0.011	1.108 ± 0.011	1.325 ± 0.021	3.50	4	0.208 ± 0.037	0.183 ± 0.004	0.733 ± 0.151	0.460 ± 0.500	2.75	2
	FEATURE-MAP	1.551 ± 0.013	0.994 ± 0.004	<u>1.092 ± 0.008</u>	1.415 ± 0.030	5.00	6	0.194 ± 0.009	<b>0.142 ± 0.004</b>	<b>0.327 ± 0.034</b>	<b>0.252 ± 0.026</b>	1.50	1
	BSS	<u>0.926 ± 0.028</u>	<u>0.704 ± 0.010</u>	<u>1.092 ± 0.008</u>	<b>1.302 ± 0.012</b>	2.00	2	<u>0.172 ± 0.012</u>	0.213 ± 0.005	0.929 ± 0.013	0.651 ± 0.079	4.25	5
	FULL-FT	1.332 ± 0.009	0.814 ± 0.013	0.908 ± 0.005	1.472 ± 0.016	3.25	3	0.192 ± 0.022	0.221 ± 0.013	0.836 ± 0.044	0.474 ± 0.042	3.75	3
	LP	2.309 ± 0.030	1.024 ± 0.014	0.927 ± 0.010	3.814 ± 0.175	7.75	8	0.752 ± 0.006	0.881 ± 0.004	0.996 ± 0.005	1.540 ± 0.015	6.75	7
SIZE	SURGICAL-FT	1.915 ± 0.036	0.886 ± 0.013	0.925 ± 0.003	2.135 ± 0.038	6.00	5	1.589 ± 0.314	0.353 ± 0.005	0.787 ± 0.018	0.943 ± 0.000	5.25	6
	LP-FT	<b>1.754 ± 0.075</b>	<b>0.775 ± 0.005</b>	0.907 ± 0.020	<b>1.710 ± 0.010</b>	1.75	1	<b>0.145 ± 0.007</b>	<b>0.195 ± 0.007</b>	0.902 ± 0.067	0.575 ± 0.078	3.25	4
	WISE-FT	2.233 ± 0.041	0.974 ± 0.016	0.895 ± 0.011	1.982 ± 0.039	5.50	7	2.264 ± 0.336	1.226 ± 0.096	1.180 ± 0.022	2.683 ± 0.151	8.00	8
	L2-SP	1.849 ± 0.025	0.911 ± 0.009	0.911 ± 0.009	1.748 ± 0.014	4.50	4	0.192 ± 0.014	0.196 ± 0.009	0.759 ± 0.029	0.456 ± 0.109	3.00	3
	FEATURE-MAP	2.136 ± 0.030	1.007 ± 0.010	<b>0.891 ± 0.012</b>	1.947 ± 0.013	4.75	6	0.209 ± 0.014	<b>0.153 ± 0.009</b>	<b>0.354 ± 0.007</b>	<b>0.227 ± 0.048</b>	1.50	1
	BSS	<b>1.188 ± 0.039</b>	0.818 ± 0.025	0.899 ± 0.006	1.712 ± 0.021	2.50	2	<b>0.188 ± 0.019</b>	0.211 ± 0.006	0.946 ± 0.006	0.550 ± 0.040	4.00	5



Table 14: Robust fine-tuning performance on 5 **Classification** datasets (AUC metrics) in the **Fewshot** setting (covering FEWSHOT-50, FEWSHOT-100, FEWSHOT-500), evaluated across 3 dataset splits (RANDOM, SCAFFOLD, SIZE) over **GRAPHMAE** and **GRAPHGPS** models. We **bold** and **underline** the best and second-best performances in each scenario.

SPLIT	METHODS	SELF-SUPERVISED PRE-TRAINING (GRAPHMAE)										SUPERVISED PRE-TRAINING (GRAPHGPS)									
		CLINTOX		BBBP		HIV		SIDER		AVG		CLINTOX		BBBP		HIV		SIDER		AVG	
		BACK	FOR	BACK	FOR	BACK	FOR	BACK	FOR	AVG-F	AVG-R	BACK	FOR	BACK	FOR	BACK	FOR	AVG-F	AVG-R		
FEWSHOT-50																					
RANDOM	FULL-FT	59.67 ± 3.35	83.04 ± 0.39	74.97 ± 1.30	<b>62.63 ± 0.92</b>	52.52 ± 0.19	66.57	65.76	4.20	98.37 ± 0.28	56.80 ± 1.97	59.17 ± 2.03	72.16 ± 3.04	85.04 ± 0.42	74.31	72.12	3.80				
	LP	57.56 ± 4.09	71.69 ± 0.89	72.96 ± 0.91	48.27 ± 4.06	55.09 ± 0.22	61.11	61.45	6.20	97.63 ± 0.56	56.07 ± 1.71	56.41 ± 0.20	63.26 ± 0.78	82.77 ± 2.22	71.23	67.48	6.40				
	SURGICAL-FT	59.83 ± 2.64	78.37 ± 1.06	<u>75.25 ± 0.92</u>	53.35 ± 0.81	54.97 ± 0.43	64.35	63.35	4.40	97.73 ± 0.49	<u>60.73 ± 3.76</u>	58.30 ± 0.86	72.63 ± 0.13	<b>86.59 ± 0.75</b>	75.20	73.32	3.20				
	LP-FT	60.20 ± 2.14	<b>84.54 ± 0.41</b>	<b>76.92 ± 0.34</b>	62.24 ± 0.58	54.41 ± 0.32	67.64	66.32	2.60	97.37 ± 0.89	55.13 ± 0.23	56.28 ± 0.08	61.68 ± 3.35	83.52 ± 1.15	70.80	67.16	7.20				
	WISE-FT	<b>63.50 ± 7.72</b>	70.77 ± 1.42	70.57 ± 1.13	58.10 ± 2.35	51.23 ± 2.01	62.83	64.06	6.00	97.59 ± 0.27	53.55 ± 1.95	53.17 ± 3.12	64.90 ± 3.22	83.69 ± 0.21	70.58	67.38	6.80				
	L <sup>2</sup> -SP	<u>61.02 ± 2.03</u>	83.79 ± 0.60	74.24 ± 0.96	61.58 ± 0.81	<u>55.34 ± 0.44</u>	67.19	65.61	3.20	<b>98.74 ± 0.40</b>	58.95 ± 2.37	<u>61.20 ± 3.46</u>	<u>72.90 ± 2.19</u>	<u>85.15 ± 1.21</u>	75.39	73.08	2.20				
SCAFFOLD	FEATURE-MAP	59.96 ± 3.80	73.57 ± 1.12	71.18 ± 2.60	48.24 ± 4.14	<b>55.85 ± 0.10</b>	61.77	62.34	5.20	98.10 ± 0.33	59.51 ± 0.56	<b>61.65 ± 0.88</b>	68.77 ± 2.81	82.13 ± 0.22	74.15	71.05	4.20				
	BSS	58.86 ± 3.63	83.81 ± 0.57	74.38 ± 1.20	62.06 ± 0.80	54.46 ± 0.56	66.71	65.10	4.20	98.43 ± 0.09	<b>63.68 ± 3.86</b>	59.82 ± 3.70	<b>73.10 ± 1.05</b>	85.03 ± 0.39	76.01	73.94	2.20				
	FULL-FT	55.61 ± 2.60	58.53 ± 0.58	58.21 ± 7.54	45.89 ± 4.20	<u>54.86 ± 0.67</u>	54.62	56.23	5.60	98.29 ± 0.28	52.89 ± 0.45	<b>64.90 ± 1.55</b>	72.07 ± 2.45	84.83 ± 0.05	74.60	73.93	3.80				
	LP	62.76 ± 3.66	56.21 ± 1.38	56.67 ± 6.74	52.12 ± 3.82	53.39 ± 0.50	56.23	55.42	6.20	97.98 ± 0.48	56.21 ± 2.18	56.28 ± 0.18	63.27 ± 0.78	82.52 ± 0.30	71.25	67.36	6.40				
	SURGICAL-FT	63.53 ± 3.11	59.33 ± 0.82	60.97 ± 3.53	52.62 ± 1.46	<b>54.94 ± 0.39</b>	58.28	58.41	3.00	97.72 ± 0.49	<b>61.37 ± 2.90</b>	58.30 ± 0.86	72.63 ± 0.13	<b>86.59 ± 0.75</b>	75.20	73.32	3.20				
	LP-FT	60.20 ± 2.14	84.54 ± 0.41	76.92 ± 0.34	62.24 ± 0.58	54.41 ± 0.32	67.64	66.32	2.60	97.37 ± 0.89	55.14 ± 0.44	56.41 ± 0.20	63.27 ± 0.78	83.49 ± 1.19	70.83	67.19	6.60				
SIZE	WISE-FT	55.45 ± 5.80	59.33 ± 0.74	<b>67.39 ± 0.29</b>	<u>58.03 ± 0.66</u>	53.77 ± 4.09	58.79	57.60	5.20	98.23 ± 0.05	50.43 ± 0.95	54.67 ± 0.12	66.17 ± 5.35	83.75 ± 0.09	60.65	68.19	6.20				
	L <sup>2</sup> -SP	<u>64.57 ± 0.27</u>	59.09 ± 0.16	<u>75.91 ± 0.12</u>	54.91 ± 0.32	54.37 ± 0.42	65.84	67.40	5.20	<b>98.72 ± 0.42</b>	57.64 ± 2.35	59.55 ± 0.38	72.49 ± 0.17	85.89 ± 0.16	78.32	72.35	2.80				
	FEATURE-MAP	<b>68.84 ± 1.77</b>	<u>56.59 ± 1.37</u>	<u>64.71 ± 2.65</u>	43.90 ± 0.98	<u>50.07 ± 0.75</u>	58.62	57.12	5.20	98.33 ± 0.07	58.93 ± 0.76	59.64 ± 0.10	68.71 ± 3.16	82.73 ± 0.19	73.67	70.37	3.80				
	BSS	<u>60.27 ± 3.40</u>	<b>60.16 ± 0.57</b>	<b>61.83 ± 0.07</b>	<b>62.17 ± 1.89</b>	<b>54.35 ± 0.96</b>	67.96	60.75	3.60	<u>98.53 ± 0.13</u>	<u>59.09 ± 2.98</u>	<u>59.36 ± 2.79</u>	<b>72.34 ± 1.36</b>	<b>85.12 ± 0.23</b>	75.07	72.57	2.20				
	FULL-FT	55.86 ± 4.15	58.45 ± 1.93	<u>64.53 ± 8.42</u>	51.39 ± 0.87	52.27 ± 0.60	52.36	56.51	5.40	98.34 ± 0.26	55.58 ± 1.28	56.71 ± 0.33	<b>73.23 ± 1.85</b>	83.13 ± 0.22	74.60	73.02	3.20				
	LP	52.46 ± 1.91	58.50 ± 0.51	58.50 ± 0.51	52.62 ± 1.46	52.62 ± 1.46	50.67	50.67	5.40	97.90 ± 0.47	50.67 ± 0.47	50.67 ± 0.47	63.27 ± 0.78	82.52 ± 0.30	71.25	67.36	6.40				
RANDOM	SURGICAL-FT	53.27 ± 3.82	48.97 ± 1.11	<b>62.03 ± 0.95</b>	52.11 ± 0.15	53.37 ± 0.34	51.95	54.27	4.40	97.70 ± 0.51	61.72 ± 5.14	59.07 ± 0.65	72.65 ± 0.15	<b>86.48 ± 0.70</b>	75.52	73.02	6.80				
	LP-FT	54.33 ± 3.19	59.46 ± 1.82	60.76 ± 2.04	<b>67.05 ± 1.85</b>	<b>63.41 ± 0.19</b>	53.02	54.96	3.40	97.37 ± 0.85	53.39 ± 0.19	56.32 ± 0.16	71.35 ± 3.29	83.17 ± 0.75	70.48	67.15	3.20				
	WISE-FT	55.45 ± 5.80	59.33 ± 0.74	<b>67.39 ± 0.29</b>	<u>58.03 ± 0.66</u>	53.77 ± 4.09	58.79	57.60	5.20	98.23 ± 0.05	50.43 ± 0.95	54.67 ± 0.12	66.17 ± 5.35	83.75 ± 0.09	60.65	68.19	6.20				
	L <sup>2</sup> -SP	<u>64.57 ± 0.27</u>	59.09 ± 0.16	<u>75.91 ± 0.12</u>	54.91 ± 0.32	54.37 ± 0.42	65.84	67.40	5.20	<b>98.72 ± 0.42</b>	57.64 ± 2.35	59.55 ± 0.38	72.49 ± 0.17	85.89 ± 0.16	78.32	72.35	2.80				
	FEATURE-MAP	53.75 ± 1.04	60.21 ± 2.22	60.45 ± 1.64	54.23 ± 4.82	51.88 ± 0.54	53.18	53.02	4.20	98.04 ± 0.40	59.55 ± 0.79	<b>61.03 ± 0.46</b>	68.77 ± 3.28	82.85 ± 0.21	74.13	71.03	4.40				
	BSS	<b>58.40 ± 1.49</b>	59.13 ± 1.42	66.02 ± 0.89	59.34 ± 1.41																



Table 15: Robust fine-tuning performance on 4 **Regression** datasets (RMSE metrics) in the **Fewshot** setting (covering FEWSHOT-50, FEWSHOT-100, and FEWSHOT-500), evaluated across 3 dataset splits (RANDOM, SCAFFOLD, SIZE) over **GRAPHMAE** and **GRAPHGPS** models. AVG-R, AVG-R\* denote the average rank and the rank based on the average normalized performance over all the datasets for each evaluated method, respectively. Standard deviations across five replicates are shown in parentheses. We **bold** and underline the best and second-best performances in each scenario.

SPLIT	METHODS	SELF-SUPERVISED PRE-TRAINING (GRAPHMAE)					SUPERVISED PRE-TRAINING (GRAPHGPS)							
		ESOL	LIPO	MALARIA	CEP	AVG-R	AVG-R*	ESOL	LIPO	MALARIA	CEP	AVG-R	AVG-R*	
RANDOM	FEWSHOT-50													
	FULL-FT	1.432 ± 0.019	1.328 ± 0.051	1.297 ± 0.015	2.927 ± 0.226	4.25	6	0.896 ± 0.015	1.221 ± 0.016	1.192 ± 0.017	2.072 ± 0.050	4.00	4	
	LP	1.646 ± 0.027	1.395 ± 0.076	1.334 ± 0.009	4.133 ± 0.372	7.50	8	1.183 ± 0.012	1.223 ± 0.007	1.193 ± 0.009	2.219 ± 0.016	6.00	6	
	SURGICAL-FT	1.497 ± 0.017	1.303 ± 0.051	1.309 ± 0.017	3.300 ± 0.406	5.00	7	3.573 ± 0.101	2.168 ± 0.089	1.203 ± 0.010	4.263 ± 0.096	7.75	8	
	LP-FT	1.886 ± 0.022	1.217 ± 0.021	1.399 ± 0.033	2.840 ± 0.226	3.75	5	1.037 ± 0.209	1.199 ± 0.041	1.178 ± 0.014	2.156 ± 0.145	3.50	5	
	WISE-FT	1.622 ± 0.053	1.343 ± 0.010	1.248 ± 0.008	2.385 ± 0.026	3.75	2	2.488 ± 0.137	1.221 ± 0.007	1.180 ± 0.019	2.574 ± 0.053	6.00	7	
	L2-SP	1.444 ± 0.027	1.354 ± 0.052	1.294 ± 0.005	2.315 ± 0.106	3.75	1	0.881 ± 0.037	1.303 ± 0.022	1.184 ± 0.013	2.091 ± 0.049	2.50	3	
	FEATURE-MAP	1.655 ± 0.027	1.312 ± 0.020	1.278 ± 0.003	2.363 ± 0.127	3.75	3	0.882 ± 0.059	1.173 ± 0.013	1.193 ± 0.006	2.050 ± 0.029	2.75	2	
	BSS	1.439 ± 0.029	1.351 ± 0.051	1.294 ± 0.005	2.682 ± 0.115	4.00	4	0.822 ± 0.024	1.204 ± 0.021	1.189 ± 0.011	2.109 ± 0.036	3.25	1	
	SCAFFOLD	FULL-FT	1.717 ± 0.028	1.214 ± 0.051	1.169 ± 0.005	2.612 ± 0.178	5.25	6	0.859 ± 0.065	1.219 ± 0.025	1.426 ± 0.243	2.100 ± 0.031	4.00	3
LP		2.209 ± 0.039	1.183 ± 0.045	1.170 ± 0.004	2.656 ± 0.048	6.00	8	1.213 ± 0.015	1.223 ± 0.006	1.194 ± 0.012	2.261 ± 0.019	5.50	6	
SURGICAL-FT		1.834 ± 0.031	1.198 ± 0.049	1.166 ± 0.001	3.142 ± 0.589	5.25	7	3.589 ± 0.101	2.168 ± 0.089	1.204 ± 0.010	4.261 ± 0.096	7.00	8	
LP-FT		1.642 ± 0.026	1.147 ± 0.038	1.300 ± 0.061	2.879 ± 0.264	4.25	4	1.053 ± 0.180	1.198 ± 0.043	1.174 ± 0.019	2.633 ± 0.025	4.00	4	
WISE-FT		2.221 ± 0.047	1.175 ± 0.016	1.166 ± 0.002	2.326 ± 0.031	4.00	3	1.020 ± 0.045	1.259 ± 0.027	1.238 ± 0.012	2.123 ± 0.035	6.00	7	
L2-SP		1.718 ± 0.053	1.200 ± 0.053	1.202 ± 0.062	2.366 ± 0.059	5.00	5	0.897 ± 0.058	1.196 ± 0.030	1.205 ± 0.032	2.105 ± 0.032	3.00	2	
FEATURE-MAP		2.197 ± 0.075	1.148 ± 0.023	1.163 ± 0.003	2.400 ± 0.175	3.00	1	0.898 ± 0.040	1.200 ± 0.013	1.229 ± 0.014	2.115 ± 0.031	4.25	5	
BSS		1.715 ± 0.056	1.168 ± 0.050	1.168 ± 0.002	2.551 ± 0.121	3.25	2	0.863 ± 0.024	1.308 ± 0.016	1.186 ± 0.010	2.081 ± 0.037	2.25	1	
SIZE		FULL-FT	2.654 ± 0.075	1.557 ± 0.093	0.943 ± 0.026	2.550 ± 0.053	4.25	5	0.886 ± 0.054	1.209 ± 0.011	1.179 ± 0.017	2.098 ± 0.048	3.75	4
		LP	2.818 ± 0.087	1.676 ± 0.115	0.963 ± 0.030	5.414 ± 0.036	7.00	8	1.176 ± 0.011	1.232 ± 0.007	1.181 ± 0.009	2.248 ± 0.016	6.75	7
	SURGICAL-FT	2.658 ± 0.088	1.641 ± 0.114	0.929 ± 0.027	3.423 ± 0.550	5.75	6	3.589 ± 0.101	2.168 ± 0.089	1.192 ± 0.010	4.258 ± 0.095	8.00	8	
	LP-FT	2.440 ± 0.056	1.422 ± 0.111	1.166 ± 0.053	2.339 ± 0.049	2.75	1	1.049 ± 0.186	1.204 ± 0.047	1.174 ± 0.016	2.167 ± 0.129	5.25	6	
	WISE-FT	3.050 ± 0.087	1.513 ± 0.049	0.969 ± 0.001	3.223 ± 0.224	5.75	7	1.045 ± 0.054	1.230 ± 0.015	1.171 ± 0.025	2.126 ± 0.039	4.50	5	
	L2-SP	2.606 ± 0.085	1.614 ± 0.112	0.914 ± 0.016	2.466 ± 0.079	3.00	2	0.851 ± 0.036	1.194 ± 0.015	1.169 ± 0.005	2.101 ± 0.022	2.00	1	
	FEATURE-MAP	2.630 ± 0.036	1.697 ± 0.080	0.920 ± 0.007	2.408 ± 0.057	4.00	4	0.867 ± 0.035	1.180 ± 0.014	1.183 ± 0.007	2.039 ± 0.022	3.00	3	
	BSS	2.579 ± 0.066	1.613 ± 0.110	0.926 ± 0.018	2.580 ± 0.157	3.50	3	0.844 ± 0.037	1.200 ± 0.028	1.171 ± 0.033	2.104 ± 0.032	2.75	2	
	FEWSHOT-100													
	RANDOM	FULL-FT	1.304 ± 0.041	1.239 ± 0.032	1.289 ± 0.003	3.028 ± 0.310	3.25	1	0.412 ± 0.033	1.061 ± 0.017	1.140 ± 0.016	1.976 ± 0.031	3.00	3
LP		1.609 ± 0.032	1.285 ± 0.043	1.334 ± 0.009	4.562 ± 0.047	7.50	8	0.902 ± 0.037	1.185 ± 0.007	1.174 ± 0.004	2.239 ± 0.010	7.25	7	
SURGICAL-FT		1.356 ± 0.022	1.219 ± 0.010	1.288 ± 0.008	3.100 ± 0.805	4.50	5	3.371 ± 0.120	1.925 ± 0.045	1.162 ± 0.013	4.076 ± 0.046	7.50	8	
LP-FT		1.310 ± 0.021	1.228 ± 0.021	1.274 ± 0.045	2.241 ± 0.438	4.75	6	0.785 ± 0.230	1.144 ± 0.049	1.153 ± 0.030	2.158 ± 0.100	5.50	6	
WISE-FT		1.600 ± 0.051	1.324 ± 0.013	1.245 ± 0.017	2.294 ± 0.024	4.75	7	0.671 ± 0.104	1.068 ± 0.049	1.159 ± 0.036	2.017 ± 0.095	5.00	5	
L2-SP		1.323 ± 0.034	1.253 ± 0.029	1.276 ± 0.014	2.271 ± 0.065	3.25	2	0.405 ± 0.034	1.055 ± 0.022	1.129 ± 0.016	1.951 ± 0.045	1.75	1	
FEATURE-MAP		1.526 ± 0.030	1.243 ± 0.027	1.276 ± 0.004	2.271 ± 0.116	3.75	3	0.422 ± 0.021	1.014 ± 0.006	1.170 ± 0.013	1.883 ± 0.012	3.25	4	
BSS		1.322 ± 0.053	1.251 ± 0.028	1.293 ± 0.006	2.541 ± 0.128	4.25	4	0.405 ± 0.060	1.050 ± 0.003	1.147 ± 0.014	1.986 ± 0.018	2.75	2	
SCAFFOLD		FULL-FT	1.695 ± 0.045	1.168 ± 0.030	1.167 ± 0.003	3.087 ± 0.765	4.50	2	0.497 ± 0.045	1.125 ± 0.034	1.215 ± 0.015	2.036 ± 0.073	4.75	6
		LP	2.045 ± 0.044	1.211 ± 0.064	1.173 ± 0.004	4.579 ± 0.037	7.50	8	0.971 ± 0.036	1.185 ± 0.008	1.174 ± 0.004	2.247 ± 0.005	6.25	5
	SURGICAL-FT	1.693 ± 0.019	1.146 ± 0.017	1.169 ± 0.003	3.226 ± 0.563	4.50	1	3.386 ± 0.120	1.927 ± 0.041	1.162 ± 0.013	4.073 ± 0.048	7.00	8	
	LP-FT	1.626 ± 0.016	1.123 ± 0.011	1.312 ± 0.023	2.782 ± 0.364	3.75	5	0.730 ± 0.236	1.136 ± 0.029	1.154 ± 0.029	2.167 ± 0.117	4.75	3	
	WISE-FT	2.069 ± 0.066	1.205 ± 0.014	1.158 ± 0.008	2.244 ± 0.068	4.25	7	1.069 ± 0.332	1.124 ± 0.023	1.228 ± 0.016	2.143 ± 0.115	6.00	7	
	L2-SP	1.679 ± 0.045	1.169 ± 0.048	1.168 ± 0.003	3.237 ± 0.030	3.50	4	0.497 ± 0.060	1.098 ± 0.015	1.155 ± 0.022	2.031 ± 0.061	3.25	2	
	FEATURE-MAP	1.964 ± 0.034	1.164 ± 0.029	1.164 ± 0.001	2.341 ± 0.095	3.50	6	0.489 ± 0.040	1.039 ± 0.014	1.185 ± 0.010	2.008 ± 0.022	2.75	4	
	BSS	1.681 ± 0.043	1.191 ± 0.046	1.169 ± 0.004	2.566 ± 0.149	4.50	3	0.396 ± 0.010	1.054 ± 0.033	1.139 ± 0.005	1.972 ± 0.010	1.25	1	
	SIZE	FULL-FT	2.414 ± 0.081	1.283 ± 0.070	0.911 ± 0.008	2.677 ± 0.139	3.00	1	0.431 ± 0.059	1.039 ± 0.026	1.118 ± 0.014	1.968 ± 0.056	3.25	4
		LP	2.859 ± 0.078	1.493 ± 0.115	0.951 ± 0.030	5.420 ± 0.033	7.50	8	0.901 ± 0.037	1.192 ± 0.007	1.163 ± 0.004	2.236 ± 0.011	7.25	7
SURGICAL-FT		2.537 ± 0.059	1.301 ± 0.074	0.909 ± 0.002	3.707 ± 0.589	4.75	6	3.386 ± 0.120	1.933 ± 0.038	1.151 ± 0.012	4.077 ± 0.040	7.50	8	
LP-FT		2.217 ± 0.047	1.146 ± 0.022	1.065 ± 0.020	2.562 ± 0.076	3.50	2	0.733 ± 0.232	1.166 ± 0.029	1.147 ± 0.022	2.138 ± 0.123	5.50	6	
WISE-FT		2.507 ± 0.098	1.297 ± 0.038	0.904 ± 0.002	2.823 ± 0.031	3.75	3	0.708 ± 0.099	1.079 ± 0.040	1.147 ± 0.040	1.987 ± 0.050	4.75	5	
L2-SP		2.442 ± 0.047	1.362 ± 0.082	0.916 ± 0.009	2.451 ± 0.093	4.50	5	0.409 ± 0.024	1.037 ± 0.030	1.125 ± 0.016	1.942 ± 0.032	2.00	1	
FEATURE-MAP		2.716 ± 0.026	1.531 ± 0.085	0.912 ± 0.003	2.424 ± 0.039	5.00	7	0.419 ± 0.016	1.009 ± 0.013	1.160 ± 0.010	1.886 ± 0.031	3.00	3	
BSS		2.434 ± 0.046	1.358 ± 0.084	0.912 ± 0.005	2.533 ± 0.103	3.75	3	0.387 ± 0.020	1.038 ± 0.021	1.136 ± 0.013	1.967 ± 0.023	2.50	2	
FEWSHOT-500														
RANDOM		FULL-FT	1.042 ± 0.017	1.023 ± 0.022	1.290 ± 0.004	1.958 ± 0.038	4.00	5	0.135 ± 0.019	0.070 ± 0.005	0.787 ± 0.009	1.554 ± 0.044	3.25	3
	LP	1.487 ± 0.011	1.233 ± 0.019	1.331 ± 0.012	4.602 ± 0.019	8.00	8	0.769 ± 0.108	0.854 ± 0.008	1.035 ± 0.001	1.941 ± 0.004	6.00	6	
	SURGICAL-FT	1.164 ± 0.010	1.127 ± 0.007	1.240 ± 0.011	3.577 ± 0.498	5.00	7	2.376 ± 0.207	0.806 ± 0.037	0.803 ± 0.010	3.058 ± 0.054	6.75	7	
	LP-FT	0.995 ± 0.010	0.975 ± 0.007	1.310 ± 0.019	2.004 ± 0.056	3.75	4	0.545 ± 0.293	0.605 ± 0.352	0.793 ± 0.018	1.566 ± 0.027	5.50	5	
	WISE-FT	1.251 ± 0.029	0.979 ± 0.010	1.231 ± 0.016	1.975 ± 0.017	3.25	1	2.512 ± 0.245	1.563 ± 0.200	1.197 ± 0.017	2.177 ± 0.063	7.75	8	
	L2-SP	1.048 ± 0.014	1.036 ± 0.009	1.241 ± 0.007	1.886 ± 0.032	3.25	1	0.141 ± 0.043	0.080 ± 0.026	0.781 ± 0.010	1.549 ± 0.022	2.75	2	
	FEATURE-MAP	1.340 ± 0.007	1.202 ± 0.014	1.241 ± 0.007	1.992 ± 0.013	5.75	6	0.155 ± 0.021	0.104 ± 0.005	0.778 ± 0.004	1.565 ± 0.028	3.25	4	
	BSS	1.031 ± 0.013	1.020 ± 0.006	1.272 ± 0.007	1.896 ± 0.034	3.00	3	0.129 ± 0.018	0.018 ± 0.004	0.779 ± 0.007	1.543 ± 0.028	1.25	1	
	SCAFFOLD	FULL-FT	1.406 ± 0.016	0.945 ± 0.021	1.199 ± 0.025	2.057 ± 0.072	4.75	5	0.145 ± 0.023	0.072 ± 0.005	0.776 ± 0.006	1.564 ± 0.033	2.50	3
		LP	1.849 ± 0.028	1.102 ± 0.019	1.182 ± 0.007	4.607 ± 0.020	7.00	8	0.771 ± 0.018	0.854 ± 0.008	1.035 ± 0.001	1.941 ± 0.004	6.50	6
SURGICAL-FT		1.436 ± 0.010	1.029 ± 0.006	1.156 ± 0.010	2.874 ± 0.632	5.00	9	2.377 ± 0.207	0.806 ± 0.041	0.803 ± 0.011	3.053 ± 0.051	6.50	6	
LP-FT		1.054 ± 0.011	0.940 ± 0.012	1.275 ± 0.014	2.052 ± 0.036	3.75	4	0.446 ± 0.030	0.605 ± 0.352	0.793 ± 0.018	1.566 ± 0.027	5.50	5	
WISE-FT		1.707 ± 0.029	1.028 ± 0.025	1.125 ± 0.008	1.906 ± 0.020	3.50	3	2.476 ± 0.266	1.459 ± 0.258	1.207 ± 0.030	2.173 ± 0.061	7.75	8	
L2-SP		1.413 ± 0.045	0.943 ± 0.022	1.156 ± 0.012	1.931 ± 0.054	3.25	2	0.137 ± 0.017	0.070 ± 0.009	0.782 ± 0.005	1.524 ± 0.014	2.25	2	
FEATURE-MAP		1.880 ± 0.021	1.081 ± 0.006	1.129 ± 0.006	1.992 ± 0.008	5.25	7	0.163 ± 0.010	0.111 ± 0.002	0.786 ± 0.005	1.592 ± 0.013	4.00	4	
BSS		0.941 ± 0.013	1.199 ± 0.013	1.272 ± 0.007	1.896 ± 0.034	3.00	3	0.107 ± 0.013	0.068 ± 0.004	0.779 ± 0.007	1.543 ± 0.028	1.25	1	
SIZE		FULL-FT	2.102 ± 0.080	0.968 ± 0.032	0.955 ± 0.031	2.283 ± 0.090	3.50	4	0.142 ± 0.049	0.070 ± 0.003	0.723 ± 0.008	1.548 ± 0.011	3.00	3
		LP	2.486 ± 0.040	1.140 ± 0.046	0.968 ± 0.027	5.452 ± 0.018	7.50	8	0.771 ± 0.018	0.855 ± 0.009	1.008 ± 0.004	1.938 ± 0.014	6.50	6
	SURGICAL-FT	2.142 ± 0.062	0.982 ± 0.014	0.949 ± 0.032	3.765 ± 0.499	4.50	7	2.384 ± 0.212	0.812 ± 0.042	0.745 ± 0.011				

Table 17: DWiSE-FT performance on 4 **Regression** datasets (RMSE metrics) in the **Fewshot** setting with 50,100, 500 samples, evaluated across 3 dataset splits (RANDOM, SCAFFOLD, SIZE) given **MOLE-BERT** model. AVG-R denote the average rank. Standard deviations across five replicates are shown in parentheses. We **bold** and underline the best and second-best performances in each scenario.

SPLIT	METHODS	FEWSHOT 50					FEWSHOT 100					FEWSHOT 500				
		ESOL	LIPO	MALARIA	CEP	AVG	ESOL	LIPO	MALARIA	CEP	AVG	ESOL	LIPO	MALARIA	CEP	AVG
RANDOM	WISE-FT	1.384 ± 0.047	1.212 ± 0.020	1.270 ± 0.007	2.410 ± 0.051	3.75	1.189 ± 0.030	1.142 ± 0.025	<b>1.256 ± 0.006</b>	2.211 ± 0.028	3.00	0.995 ± 0.010	0.855 ± 0.011	1.193 ± 0.003	1.893 ± 0.021	3.75
	L <sup>2</sup> -SP	<u>1.372 ± 0.029</u>	1.196 ± 0.019	1.277 ± 0.006	2.280 ± 0.031	3.00	1.161 ± 0.016	1.149 ± 0.007	1.260 ± 0.004	<u>2.131 ± 0.014</u>	3.25	<b>0.878 ± 0.020</b>	<b>0.806 ± 0.007</b>	<b>1.192 ± 0.004</b>	1.893 ± 0.018	1.50
	Top	<b>1.329 ± 0.021</b>	<b>1.164 ± 0.010</b>	<b>1.271 ± 0.007</b>	<u>2.273 ± 0.022</u>	1.25	<b>1.120 ± 0.038</b>	<u>1.139 ± 0.017</u>	<b>1.256 ± 0.006</b>	<u>2.131 ± 0.014</u>	1.50	<b>0.878 ± 0.020</b>	<b>0.806 ± 0.007</b>	<b>1.192 ± 0.004</b>	<b>1.862 ± 0.010</b>	1.00
	DWiSE-FT	1.378 ± 0.055	1.189 ± 0.020	1.273 ± 0.009	<b>2.222 ± 0.050</b>	2.00	<u>1.132 ± 0.025</u>	<b>1.138 ± 0.028</b>	<b>1.256 ± 0.004</b>	<b>2.129 ± 0.020</b>	1.25	0.918 ± 0.012	0.818 ± 0.013	<b>1.192 ± 0.004</b>	1.865 ± 0.030	2.25
SCAFFOLD	WISE-FT	1.842 ± 0.056	1.177 ± 0.009	<u>1.162 ± 0.004</u>	2.454 ± 0.043	3.50	1.544 ± 0.063	1.041 ± 0.017	<u>1.151 ± 0.007</u>	2.301 ± 0.042	3.50	1.388 ± 0.023	0.834 ± 0.012	<b>1.114 ± 0.002</b>	1.936 ± 0.037	3.25
	L <sup>2</sup> -SP	1.699 ± 0.049	1.066 ± 0.009	<u>1.162 ± 0.002</u>	2.331 ± 0.024	2.50	1.473 ± 0.069	0.961 ± 0.003	1.153 ± 0.007	2.201 ± 0.038	2.50	1.163 ± 0.026	0.813 ± 0.010	1.126 ± 0.011	1.885 ± 0.011	2.50
	Top	1.680 ± 0.042	<b>1.036 ± 0.007</b>	<b>1.159 ± 0.000</b>	<b>2.292 ± 0.026</b>	1.25	<b>1.436 ± 0.054</b>	<b>0.937 ± 0.008</b>	1.149 ± 0.003	2.187 ± 0.034	1.25	<b>1.112 ± 0.015</b>	<b>0.802 ± 0.003</b>	<b>1.114 ± 0.002</b>	<b>1.881 ± 0.010</b>	1.00
	DWiSE-FT	<b>1.616 ± 0.047</b>	1.110 ± 0.013	1.173 ± 0.005	<u>2.306 ± 0.030</u>	2.50	1.485 ± 0.041	0.979 ± 0.014	1.158 ± 0.009	<b>2.149 ± 0.040</b>	2.75	1.266 ± 0.021	0.823 ± 0.010	1.121 ± 0.004	1.900 ± 0.019	3.00
SIZE	WISE-FT	2.615 ± 0.072	1.301 ± 0.042	0.929 ± 0.004	2.762 ± 0.053	4.00	2.216 ± 0.056	1.124 ± 0.031	0.917 ± 0.004	2.543 ± 0.027	3.75	2.071 ± 0.078	0.902 ± 0.016	0.912 ± 0.003	2.379 ± 0.086	3.75
	L <sup>2</sup> -SP	2.393 ± 0.068	1.306 ± 0.037	0.915 ± 0.002	<b>2.497 ± 0.019</b>	2.50	<u>1.731 ± 0.071</u>	<b>1.025 ± 0.028</b>	0.905 ± 0.002	<u>2.424 ± 0.024</u>	1.75	1.629 ± 0.084	0.821 ± 0.011	0.904 ± 0.003	2.368 ± 0.013	2.50
	Top	<u>2.369 ± 0.075</u>	<u>1.297 ± 0.040</u>	<b>0.911 ± 0.002</b>	<b>2.497 ± 0.019</b>	1.50	<u>1.731 ± 0.071</u>	<b>1.025 ± 0.028</b>	<b>0.898 ± 0.003</b>	<u>2.424 ± 0.024</u>	1.50	1.629 ± 0.084	<b>0.803 ± 0.006</b>	<b>0.895 ± 0.002</b>	<u>2.328 ± 0.017</u>	1.50
	DWiSE-FT	<b>1.488 ± 0.101</b>	<b>1.113 ± 0.021</b>	<u>0.913 ± 0.007</u>	2.539 ± 0.023	1.75	<b>1.469 ± 0.062</b>	1.031 ± 0.022	0.920 ± 0.006	<b>2.390 ± 0.025</b>	2.25	<b>1.466 ± 0.040</b>	<u>0.816 ± 0.022</u>	0.915 ± 0.003	<b>2.322 ± 0.031</b>	2.00