

# Towards Effective and Generalizable Fine-tuning for Pre-trained Molecular Graph Models

Jun Xia<sup>1,2,3</sup>, Jiangbin Zheng<sup>2,3</sup>, Cheng Tan<sup>2,3</sup>, Ge Wang<sup>2,3</sup> and Stan Z. Li<sup>2,3</sup>

<sup>1</sup>Zhejiang University <sup>2</sup>School of Engineering, Westlake University

<sup>3</sup>Institute of Advanced Technology, Westlake Institute for Advanced Study

{xiajun, zhengjiangbin, tancheng, wangge, stan.zq.li}@westlake.edu.cn

## Abstract

Graph Neural Networks (GNNs) and Transformer have emerged as dominant tools for AI-driven drug discovery. Many state-of-the-art methods first pre-train GNNs or the hybrid of GNNs and Transformer on a large molecular database and then fine-tune on downstream tasks. However, different from other domains such as computer vision (CV) or natural language processing (NLP), getting labels for molecular data of downstream tasks often requires resource-intensive wet-lab experiments. Besides, the pre-trained models are often of extremely high complexity with huge parameters. These often cause the fine-tuned model to over-fit the training data of downstream tasks and significantly deteriorate the performance. To alleviate these critical yet under-explored issues, we propose two straightforward yet effective strategies to attain better generalization performance: 1. MolAug, which enriches the molecular datasets of down-stream tasks with chemical homologies and enantiomers; 2. WordReg, which controls the complexity of the pre-trained models with a smoothness-inducing regularization built on dropout. Extensive experiments demonstrate that our proposed strategies achieve notable and consistent improvements over vanilla fine-tuning and yield multiple state-of-the-art results. Also, these strategies are model-agnostic and readily pluggable into fine-tuning of various pre-trained molecular graph models. We will release the code and the fine-tuned models.

## 1 Introduction

Pre-trained language models (PLMs) have fundamentally changed the landscape of natural language processing (NLP) [Devlin *et al.*, 2019], which have established new state-of-the-art results for a large variety of NLP tasks. Inspired by their proliferation, tremendous efforts have been devoted to molecular graph pre-training which can exploit abundant knowledge of unlabelled molecular in the database [Hu *et al.*, 2020]. For the pre-training stage, existing works train the encoder with various pretext tasks in absence of labels [Rong *et al.*, 2020]. For the second fine-tuning stage, researchers adapt the pre-trained models to the downstream

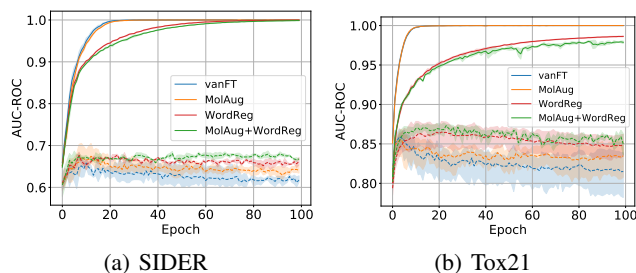


Figure 1: Training (solid lines) and testing (dashed lines) curves of various fine-tuning strategies on SIDER and Tox21 datasets. ‘vanFT’ refers to vanilla fine-tuning. The over-fitting issue of vanilla fine-tuning impede the performance improvements while our proposed MolAug, WordReg or their combination can alleviate this issue.

tasks via replacing the top layer of the pre-trained models by a task specific sub-network, and then continuing to train the new model with the limited data of the downstream task.

Despite the fruitful progress in the strategies for pre-training, the fine-tuning stage remains under-explored for pre-trained molecular graph models. There are two crucial issues impede the improvements of performance during fine-tuning: (1) insufficient labeled data for downstream tasks. Different from other domains that have abundant labeled data, getting high-quality labels for molecular data often requires resource-intensive wet-lab experiments [Xia *et al.*, 2021a]. (2) the pre-trained models are often of extremely high complexity with tens millions of parameters [Rong *et al.*, 2020; Li *et al.*, 2021b], which poses the capability of memorizing the limited samples and lead to poor generalization [Mohri *et al.*, 2012; Haoming *et al.*, 2020]. As shown in Figure 1, we conduct vanilla fine-tuning on one of the state-of-the-art pre-trained molecular models MPG [Li *et al.*, 2021b]. The over-fitting issue poses hurdle to the further improvements on various datasets. To mitigate this issue, existing fine-tuning methods in other domains often rely on hyper-parameter tuning heuristics. For example, Howard and Ruder [Howard and Ruder, 2018] follow a heuristic learning rate schedule and gradually unfreeze the layers of the pre-trained language model to improve the fine-tuning performance, which require significant tuning efforts. The other line of works propose various regularizations to control complexity of pre-trained models. Specifically, SMART [Haoming *et al.*, 2020] intro-

duces a regularization, which encourages the output of the model not to change much, when injecting an adversarial perturbation to the input. However, this strategy is not suitable for pre-trained graph models because synthesizing adversarial samples for molecular graphs online is time-consuming for its high complexity. Recently, Child-Tuning [Xu *et al.*, 2021b] is proposed to strategically mask out some gradients of network during the backward process. Albeit effective, it suffers a heavy computational overhead because it requires to search for suitable masks for each parameter in the pre-trained models. Now, we are naturally motivated to ask following question: *Can we devise more suitable strategies for effective and generalizable fine-tuning on pre-trained molecular graph models?*

To fully harness the power of pre-trained molecular graph models, we propose two strategies for better fine-tuning: MolAug and WordReg. MolAug is molecular graph augmentation with chemical enantiomers and homologies which share the similar (or identical) physical (permeability, solubility, etc.) or chemical (toxicity, etc.) properties with themselves. WordReg is a novel smoothness-inducing regularization built on dropout. To the best of our knowledge, we are the first to study the fine-tuning stage of the pre-trained molecular graph models, which is important while being neglected. We summarize our contributions here:

- We propose new data augmentations for molecular graph data, which introduce variations while not altering the physical or chemical properties of molecules too much.
- We propose a novel smoothness-inducing regularization built on dropout that can effectively control the high complexity of the pre-trained models. Also, it is domain-agnostic and we study its effects in fine-tuning pre-trained models of NLP and CV in the appendix.
- Through extensive experiments on 3 drug discovery sub-tasks (11 datasets in total) including molecular property prediction, drug-drug interaction prediction and drug-target interaction prediction, we observe consistent and notable improvements over vanilla fine-tuning and yield multiple new state-of-the-art results.

## 2 Related work

### 2.1 Molecular Representation Learning

Molecular Representation Learning refers to represent molecules in the vector space. Initially, the traditional chemical fingerprints such as ECFP [Rogers and Hahn, 2010] encode the neighbors of atoms in the molecule into a fix-length vector. With the proliferation of deep learning, [Duvenaud *et al.*, 2015] first introduce convolutional neural networks to learn neural fingerprints of molecules. Subsequently, some researchers feed the SMILES (a line notation for describing the structure of chemical species using short ASCII strings) into RNN-based models to obtain molecular representations. In order to fulfill the topology information of molecular graphs, some recent works [Kearnes *et al.*, 2016; Xiong *et al.*, 2020] attempt to apply GNNs to molecular representation learning. Besides, MPNN [Gilmer and others., 2017] and its variants DMPNN [Yang *et al.*, 2019],

CMPNN [Song *et al.*, 2020], CoMPT [Chen *et al.*, 2021] utilize a message passing framework to better capture the interactions among atoms. However, they still require expensive annotations and barely generalize to unseen molecules, which pose hurdle to the practical applications.

### 2.2 Pre-training on Graphs

To address the fundamental challenges of extremely scarce labeled data and out-of-distribution generalization in graphs learning, tremendous efforts have been devoted to pre-training on graphs recently. One line of these works follow the contrastive paradigm [Zhu *et al.*, 2021b; Zhu *et al.*, 2021a; Xia *et al.*, 2021b]. For molecular pre-training, GraphCL [You *et al.*, 2020] and its variants [You *et al.*, 2021; Susheel *et al.*, 2021; Xia *et al.*, 2022] embeds augmented versions of the anchor molecular graph close to each other and pushes the embeddings of other molecules apart. The other line of works adopt generative pretext tasks. Prototypical examples are GPT-GNN [Hu and others., 2020] which introduces a self-supervised attributed graph generation task to pre-train GNNs so that they can capture the structural and semantic properties of the graph. For molecular graph pre-training, Hu *et al.* [Hu *et al.*, 2020] conduct attribute and structure prediction at the level of individual nodes as well as entire graphs. To capture the rich information in molecular graph motifs, GROVER [Rong *et al.*, 2020] and MGSSL [Zhang *et al.*, 2021] propose to predict or generate the motifs. Analogously, MPG [Li *et al.*, 2021b] learns to compare two half-graphs (each decomposed from a graph sample) and discriminate whether they come from the same source. Despite the progress in molecular graph pre-training, few efforts have been devoted to the fine-tuning except for a recent work [Han and others., 2021] that adaptively selects and combines various auxiliary tasks with the target task in the fine-tuning stage to improve performance, which is impractical because auxiliary tasks are often unavailable during fine-tuning.

## 3 Methodology

### 3.1 MolAug: Molecular Graph Augmentations with Chemical Enantiomers and Homologies

Initially, GraphCL [You *et al.*, 2020] augments molecular graph data in the form of naive random corruption (e.g., dropping bonds, dropping atoms and etc.). However, these augmentations may alter molecular graph semantics completely even if the perturbation is weak. For example, dropping a carbon atom in the phenyl ring will alter the aromatic system and result in an alkene chain, which will drastically change the molecular properties. Besides, MoCL [Sun *et al.*, 2021] attempts to incorporate domain knowledge into graph data augmentation via replacing valid substructures in molecular graph with bioisosteres that share similar properties. However, bioisosteres are used to modify some molecular properties as expected (e.g. reduce toxicity) in drug design [Mannhold *et al.*, 2012], which may introduce incorrect supervision for molecular properties prediction such as toxicity prediction. We compare MolAug with these augmentations in section 4.5.

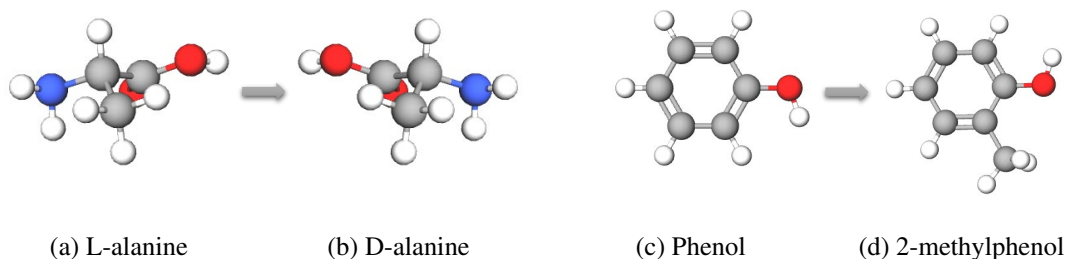


Figure 2: Illustration of MolAug. L-alanine (a) and D-alanine (b) is a pair of enantiomers; Phenol and 2-methylphenol is a pair of homologies. Gray, red, blue and white balls denote Carbon (C), Oxygen (O), Nitrogen (N) and Hydrogen (H) atoms respectively.

To alleviate these issues, we resort to chemical enantiomers and homologies. As shown in Figure 2 (a) and (b), an enantiomer is one of two stereoisomers that are mirror images of each other that are non-superposable, much as one’s left and right hands are mirror images of each other that cannot appear identical simply by reorientation. Despite the structural difference, enantiomers share the identical chemical and physical properties with each other in most cases and thus being a more suitable alternative for molecular graph augmentations in molecular properties prediction. In practice, we can obtain chemical enantiomers of a molecule by changing the chirality type which is provided in the atom (node) feature. Besides, as shown in Figure 2 (c) and (d), chemical homologies are a series of compounds differing from each other by a repeating unit, such as a methylene bridge, which share the same chemical properties with each other. Therefore, chemical homologies can serve as an ideal way of augmentation when predicting molecular chemical properties. For each molecule  $\mathbf{x}$  and its label  $y$  in the dataset  $\mathcal{D}$  of downstream task, we can obtain its various augmented versions with operation  $\text{MolAug}(\cdot)$  which can be enantiomers, homologies or their combinations. Now, we can formulate the loss of fine-tuning with MolAug as,

$$\mathcal{L}_M = \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} (l(y, f_\theta(\mathbf{x})) + \mu_1 l(y, f_\theta(\text{MolAug}(\mathbf{x})))), \quad (1)$$

where  $\mu_1$  is a trade-off parameter controlling the impact of MolAug and  $l$  is the loss function of downstream task. We can also obtain various augmentations for each molecule during each iterations with MolAug, which we study in section 4.6. Apart from pretrain-then-finetune paradigm, MolAug can also work well in training from scratch, which we validate in the appendix.

### 3.2 WordReg: Smoothness-inducing Regularization Built on Dropout

Existing fine-tuning strategies [Haoming *et al.*, 2020; Xu *et al.*, 2021b; Aghajanyan *et al.*, 2021] in NLP enforce the output of the model not to change much, when injecting perturbations to the input or model parameters. In other words, they encourage the model to be insensitive to perturbations, and thus effectively control its capacity [Mohri *et al.*, 2012]. Compared with these strategies, WordReg imposes the perturbations on the neurons which is more efficient than perturbing the input or huge model parameters. Besides, as shown in

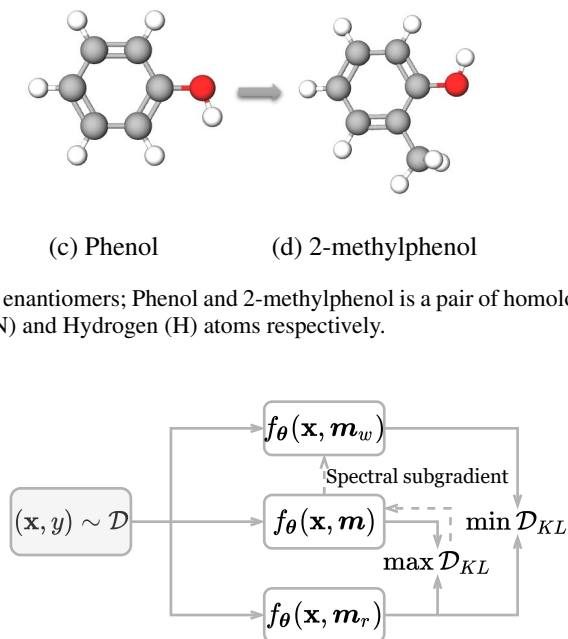


Figure 3: Illustration of WordReg. We obtain the worst-case dropout mask vector  $\mathbf{m}_w$  (corresponding to  $\mathbf{m}_r$ ) via solving a BQP problem with spectral subgradient-based method.

Figure 3, WordReg obtains the worst-case dropout via maximizing the divergence between the outputs with two different dropout which is task-dependent and can introduce stronger regularization than random dropout (validated in Figure 6). Formally, consider the pre-trained neural networks  $f_\theta(\mathbf{x}, \mathbf{m})$  that takes sample  $\mathbf{x}$  as input when fine-tuning on downstream task.  $\mathbf{m}$  is a mask vector denoting which neuron of the pre-trained model should be dropped. More specifically, for the  $i$ -th unit  $m_i$  of vector  $\mathbf{m}$ ,  $m_i = 1$  indicates that the neuron should be dropped while  $m_i = 0$  illustrates that the neuron should be preserved during dropout. To start, we introduce  $\mathbf{m}_r$  to denote the mask vector of random dropout of vanilla training. With the dropout ratio  $\sigma \in [0, 1]$  preset, we can define the constraint on  $\mathbf{m}$  as,

$$\mathcal{R}_m = \{\mathbf{m} \mid \mathbf{m} \in \{0, 1\}^N, \|\mathbf{m}\|_0 = \lfloor \sigma N \rfloor\}, \quad (2)$$

where  $\|\cdot\|_0$  is the  $\ell_0$  norm,  $N$  is the number of neurons in the model. Therefore, we can obtain the worst-case mask vector  $\mathbf{m}_w$  by solving an optimization problem,

$$\mathbf{m}_w = \arg \max_{\mathbf{m} \in \mathcal{R}_m} \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \mathcal{D}_{KL}(f_\theta(\mathbf{x}, \mathbf{m}_r) \| f_\theta(\mathbf{x}, \mathbf{m})), \quad (3)$$

where  $\mathcal{D}_{KL}$  is Kullback–Leibler divergence. With Taylor expansion, we can approximate the optimal solution as,

$$\mathbf{m}_w \approx \arg \max_{\mathbf{m} \in \mathcal{R}_m} \frac{1}{2} \mathbf{m}^T \mathbf{H} \mathbf{m}, \quad (4)$$

$$\mathbf{H} = \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \nabla^2 \mathcal{D}_{KL}(f_\theta(\mathbf{x}, \mathbf{m}_r) \| f_\theta(\mathbf{x}, \mathbf{m})) \Big|_{\mathbf{m}=\mathbf{0}}, \quad (5)$$

$\mathbf{H}$  is the Hessian matrix of the loss  $\mathcal{D}_{KL}$  at  $\mathbf{m} = \mathbf{0}$  which is a  $N \times N$  semi-positive definite matrix. Obviously, this is a Binary Quadratic Programming (BQP) problem, which is NP-hard but admits an approximate solution to

$\mathbf{m}_w$ . Here, we propose a novel method to solve this problem efficiently based on the spectral subgradient [Boyd *et al.*, 2003]. Firstly, we convert  $\{0, 1\}$ -constraint on  $\mathbf{m}$  of  $\mathcal{R}_m$  to  $\{-1, 1\}$ -constraint on  $\mathbf{n}$  via defining  $\mathbf{n} = 2\mathbf{m} - 1$ . Then we introduce a new variable  $\hat{\mathbf{n}}$  with  $N + 1$  dimensions,

$$\hat{n}_i = \begin{cases} 1, & \text{if } i = N + 1; \\ n_i, & \text{otherwise.} \end{cases} \quad (6)$$

Then we can rewrite constraint  $\mathcal{R}_m$  as a new constraint  $\mathcal{R}_{\hat{n}}$  on  $\hat{\mathbf{n}}$ ,

$$\mathcal{R}_{\hat{n}} = \{\hat{\mathbf{n}} \mid \hat{\mathbf{n}} \in \{\pm 1\}^{N+1}, \mathbf{e}^T \hat{\mathbf{n}} = c\}, \quad (7)$$

where  $c = 2\lfloor \sigma N \rfloor - N + 1$  and  $\mathbf{e} \in \mathbb{R}^{N+1}$  is an all-one vector. By these transformations, we can reformulate the BQP in Eq (5) as a new BQP in terms of  $\hat{\mathbf{n}}$ , where the constraint term  $\hat{\mathbf{n}} \in \{\pm 1\}^{N+1}$  can be rewritten as  $\hat{n}_i^2 = 1$ . We then introduce a Lagrange multiplier  $\lambda_i$  for each constraint  $\hat{n}_i^2 = 1$  and  $\lambda_0$  for the constraint  $\mathbf{e}^T \hat{\mathbf{n}} = c$ . Now, we can formulate the dual problem of the original BQP as,

$$\min_{\lambda, \lambda_0} d(\lambda, \lambda_0), \quad (8)$$

with

$$\begin{aligned} d(\lambda, \lambda_0) &= \max_{\|\hat{\mathbf{n}}\|^2 = N+1} \hat{\mathbf{n}}^T [\mathbf{L} + \text{diag}(\lambda)] \hat{\mathbf{n}} - \mathbf{e}^T \lambda - c\lambda_0 \\ &= (N + 1)\lambda_{\max} - \mathbf{e}^T \lambda - c\lambda_0 \end{aligned} \quad (9)$$

where

$$\mathbf{L} = \begin{pmatrix} \mathbf{H} & \mathbf{H}\mathbf{e} + \frac{1}{2}\lambda_0\mathbf{e} \\ \mathbf{e}^T\mathbf{H} + \frac{1}{2}\lambda_0\mathbf{e}^T & 0 \end{pmatrix} \in \mathbb{R}^{N+1 \times N+1}$$

, and  $\lambda_{\max}$  is the largest eigenvalue of  $\mathbf{L} + \text{diag}(\lambda)$ . The eigenvector of unit norm  $\mathbf{u}_{\max}$  corresponding to  $\lambda_{\max}$  can be derived via approximated by using a single-step power iteration instead of conducting naive eigenvalue decomposition. Then, the maximum  $\hat{\mathbf{n}}^*$  can be derived,

$$\hat{\mathbf{n}}^* = \sqrt{N + 1} \mathbf{u}_{\max}. \quad (10)$$

The dual problem Eq.(8) can be solved by the gradient descent method over iterations. We show the gradient of  $d$  w.r.t  $\lambda$  and  $\lambda_0$  as follows,

$$\nabla_{\lambda} d = (N + 1) \mathbf{u}_{\max}^2 - \mathbf{e}, \quad (11)$$

$$\frac{\partial d}{\partial \lambda_0} = \frac{1}{2}(N + 1) \mathbf{u}_{\max}^T \begin{pmatrix} 0 & \mathbf{e} \\ \mathbf{e} & 0 \end{pmatrix} \mathbf{u}_{\max} - c, \quad (12)$$

where  $\mathbf{u}_{\max}^2$  denotes an element-wise square of  $\mathbf{u}_{\max}$ . During fine-tuning the pre-trained models with WordReg, over each mini-batch, we compute the above gradient to make a one-step update of the Lagrange multipliers  $\lambda$  and  $\lambda_0$  along the descending direction, before the maximum  $\hat{\mathbf{n}}^*$  is taken with the updated multipliers. Finally, both  $\pm \hat{\mathbf{n}}^*$  are optimal and we should choose the one closer to  $\hat{n}_{N+1} = 1$  as required. Note that we seek the worst-case dropout mask layer-by-layer instead of applying it to an entire network as a whole. This can make the WordReg computationally efficient as well as prevent too many neurons from being dropped at a few layers.

With the worst-case dropout mask vector  $\mathbf{m}_w$  obtained, we can formulate the loss of fine-tuning with WordReg as,

$$\mathcal{L}_W = \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} (l(y, f_{\theta}(\mathbf{x})) + \mu_2 \mathcal{D}_{KL}(f_{\theta}(\mathbf{x}, \mathbf{m}_r), f_{\theta}(\mathbf{x}, \mathbf{m}_w))), \quad (13)$$

where  $\mu_2$  is a trade-off parameter controlling the impact of WordReg, which we study in section 4.6.

**Remark:** There are also some regularizations ELD [Ma *et al.*, 2017] and FD [Zolna *et al.*, 2018] built on dropout. WordReg differs from them in two aspects: (1) WordReg introduces the worst-case dropout, which is task-dependent and possess stronger regularization ability than their random drop; (2) ELD and FD are designed to reducing the gap between training (sub-model with dropout) and inference (full model without dropout) while WordReg controls the gap between sub-model with different dropouts.

### 3.3 Combining MolAug and WordReg

Additionally, we can conduct fine-tuning with MolAug and WordReg simultaneously. Formally, we specify the loss function of combining MolAug and WordReg in Eq (14).

$$\begin{aligned} \mathcal{L}_{M+W} &= \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} (l(y, f_{\theta}(\mathbf{x})) + \mu_1 l(y, f_{\theta}(\text{MolAug}(\mathbf{x}))) \\ &\quad + \mu_2 \mathcal{D}_{KL}(f_{\theta}(\mathbf{x}, \mathbf{m}_r), f_{\theta}(\mathbf{x}, \mathbf{m}_w))). \end{aligned} \quad (14)$$

The pseudo codes can be found in the appendix.

## 4 Experiments

### 4.1 Experimental Setup

**Base Pre-trained Molecular Graph Models.** We adopt recent proposed pre-trained molecular graph models MPG [Li *et al.*, 2021b] as the base model and fine-tune it with our proposed strategies. We also adopt MGSSL [Zhang *et al.*, 2021] and GraphLog [Xu *et al.*, 2021a] as our base models to validate that our strategies can be plugged into the fine-tuning of various pre-trained molecular graph models.

**Fine-tuning Tasks and Datasets.** Following previous works, we use 8 benchmark datasets from the MoleculeNet [Wu *et al.*, 2018] to perform molecular property prediction. To keep fair, we adopt the scaffold splitting method with a ratio for train/validation/test as 8:1:1 as previous works. We apply a grid search procedure to obtain the best hyper-parameters with validation set. More details of the hyper-parameters and datasets are deferred to the appendix. *The experiments on drug-drug interaction prediction and drug-target interaction prediction can be found in the appendix.*

**Baselines.** We adopt two types of baselines:

- *Training from scratch:* TF\_Robust [Rogers and Hahn, 2010] is a DNN-based multitask framework taking the molecular fingerprints as the input. GCN [Kipf. and Welling, 2017], Weave [Kearnes *et al.*, 2016] and SchNet [Schuett *et al.*, 2017] are three graph convolutional models. MPNN and its variants DMPNN, MGCN, CMPNN and CoMPT have been introduced in the related work. AttentiveFP [Xiong *et al.*, 2020] is an extension of the graph attention network.



Table 1: The performance comparison (*higher is better for classification task, lower is better for regression task*). *Relative Improvement* refers to the absolute improvement over MPG base model divided by the original results of MPG. We adopt it as a unified description of improvements for both classification and regression. The methods marked with ‘\*’ means the original papers follow a different data splitting with MPG, we fine-tune their public pre-trained models with the splitting as MPG and report the results here. The numbers in brackets are the standard deviation and the ones underlined are the previous best results. ‘—’ means the methods are too time-consuming or original implementations do not admit regression task without non-trivial modifications.

Task Dataset # Molecules	Classification (ROC-AUC)						Regression (RMSE)	
	BBBP 2039	SIDER 1427	ClinTox 1478	BACE 1513	Tox21 7831	ToxCast 8575	FreeSolv 642	ESOL 1128
ECFP [Rogers and Hahn, 2010]	0.783 <sub>(0.050)</sub>	0.630 <sub>(0.019)</sub>	0.673 <sub>(0.031)</sub>	0.861 <sub>(0.024)</sub>	0.760 <sub>(0.009)</sub>	0.615 <sub>(0.017)</sub>	5.275 <sub>(0.751)</sub>	2.359 <sub>(0.454)</sub>
TF_Robust [Ramsundar and others., 2015]	0.860 <sub>(0.087)</sub>	0.607 <sub>(0.033)</sub>	0.765 <sub>(0.085)</sub>	0.824 <sub>(0.022)</sub>	0.698 <sub>(0.012)</sub>	0.585 <sub>(0.031)</sub>	4.122 <sub>(0.085)</sub>	1.722 <sub>(0.038)</sub>
GraphConv [Kipf. and Welling, 2017]	0.877 <sub>(0.036)</sub>	0.593 <sub>(0.035)</sub>	0.845 <sub>(0.051)</sub>	0.854 <sub>(0.011)</sub>	0.772 <sub>(0.041)</sub>	0.650 <sub>(0.025)</sub>	2.900 <sub>(0.135)</sub>	1.068 <sub>(0.050)</sub>
Weave [Kearnes et al., 2016]	0.837 <sub>(0.065)</sub>	0.543 <sub>(0.034)</sub>	0.823 <sub>(0.023)</sub>	0.791 <sub>(0.008)</sub>	0.741 <sub>(0.044)</sub>	0.678 <sub>(0.024)</sub>	2.398 <sub>(0.250)</sub>	1.158 <sub>(0.055)</sub>
SchNet [Schuett et al., 2017]	0.847 <sub>(0.024)</sub>	0.545 <sub>(0.038)</sub>	0.717 <sub>(0.042)</sub>	0.750 <sub>(0.033)</sub>	0.767 <sub>(0.025)</sub>	0.679 <sub>(0.021)</sub>	3.215 <sub>(0.755)</sub>	1.045 <sub>(0.064)</sub>
MPNN [Gilmer and others., 2017]	0.913 <sub>(0.041)</sub>	0.595 <sub>(0.030)</sub>	0.879 <sub>(0.054)</sub>	0.815 <sub>(0.044)</sub>	0.808 <sub>(0.024)</sub>	0.691 <sub>(0.013)</sub>	2.185 <sub>(0.952)</sub>	1.167 <sub>(0.430)</sub>
DMPNN [Yang et al., 2019]	0.919 <sub>(0.030)</sub>	0.632 <sub>(0.023)</sub>	0.897 <sub>(0.040)</sub>	0.852 <sub>(0.053)</sub>	0.826 <sub>(0.023)</sub>	0.718 <sub>(0.011)</sub>	2.177 <sub>(0.914)</sub>	0.980 <sub>(0.258)</sub>
MGCN [Lu et al., 2019]	0.850 <sub>(0.064)</sub>	0.552 <sub>(0.018)</sub>	0.634 <sub>(0.042)</sub>	0.734 <sub>(0.030)</sub>	0.707 <sub>(0.016)</sub>	0.663 <sub>(0.009)</sub>	3.349 <sub>(0.097)</sub>	1.266 <sub>(0.147)</sub>
AttentiveFP [Xiong et al., 2020]	0.908 <sub>(0.050)</sub>	0.605 <sub>(0.060)</sub>	0.933 <sub>(0.020)</sub>	0.863 <sub>(0.015)</sub>	0.807 <sub>(0.020)</sub>	0.579 <sub>(0.001)</sub>	2.030 <sub>(0.420)</sub>	0.853 <sub>(0.060)</sub>
TrimNet [Li et al., 2021a]	0.892 <sub>(0.025)</sub>	0.606 <sub>(0.006)</sub>	0.906 <sub>(0.017)</sub>	0.843 <sub>(0.025)</sub>	0.812 <sub>(0.019)</sub>	0.652 <sub>(0.032)</sub>	2.529 <sub>(0.111)</sub>	1.282 <sub>(0.029)</sub>
CMPNN [Song et al., 2020]	0.927 <sub>(0.017)</sub>	0.616 <sub>(0.003)</sub>	0.902 <sub>(0.008)</sub>	0.856 <sub>(0.012)</sub>	0.806 <sub>(0.016)</sub>	0.681 <sub>(0.008)</sub>	2.007 <sub>(0.442)</sub>	0.798 <sub>(0.112)</sub>
CoMPT [Chen et al., 2021]	0.938 <sub>(0.021)</sub>	0.634 <sub>(0.030)</sub>	0.934 <sub>(0.019)</sub>	0.859 <sub>(0.016)</sub>	0.817 <sub>(0.014)</sub>	0.693 <sub>(0.025)</sub>	1.855 <sub>(0.578)</sub>	0.774 <sub>(0.058)</sub>
Mol2Vec [Jaeger et al., 2018]	0.876 <sub>(0.030)</sub>	0.601 <sub>(0.023)</sub>	0.828 <sub>(0.023)</sub>	0.841 <sub>(0.052)</sub>	0.805 <sub>(0.015)</sub>	0.690 <sub>(0.014)</sub>	5.752 <sub>(1.245)</sub>	2.358 <sub>(0.452)</sub>
N-GRAM [Liu et al., 2019]	0.912 <sub>(0.013)</sub>	0.632 <sub>(0.005)</sub>	0.855 <sub>(0.037)</sub>	0.876 <sub>(0.035)</sub>	0.769 <sub>(0.027)</sub>	—	2.512 <sub>(0.190)</sub>	1.100 <sub>(0.160)</sub>
SMILES-BERT [Wang et al., 2019]	<u>0.959</u> <sub>(0.009)</sub>	0.568 <sub>(0.031)</sub>	<u>0.985</u> <sub>(0.014)</sub>	0.849 <sub>(0.021)</sub>	0.803 <sub>(0.010)</sub>	—	2.974 <sub>(0.510)</sub>	0.841 <sub>(0.096)</sub>
HU.et.al.* [Hu et al., 2020]	0.915 <sub>(0.040)</sub>	0.614 <sub>(0.006)</sub>	0.762 <sub>(0.058)</sub>	0.851 <sub>(0.027)</sub>	0.811 <sub>(0.015)</sub>	0.714 <sub>(0.019)</sub>	—	—
GraphCL* [You et al., 2020]	0.906 <sub>(0.028)</sub>	0.629 <sub>(0.037)</sub>	0.865 <sub>(0.017)</sub>	0.872 <sub>(0.023)</sub>	0.826 <sub>(0.037)</sub>	0.728 <sub>(0.025)</sub>	—	—
JOAO* [You et al., 2021]	0.895 <sub>(0.033)</sub>	0.634 <sub>(0.011)</sub>	0.869 <sub>(0.022)</sub>	0.863 <sub>(0.021)</sub>	0.818 <sub>(0.021)</sub>	0.733 <sub>(0.006)</sub>	—	—
AD-GCL* [Susheel et al., 2021]	0.918 <sub>(0.025)</sub>	0.651 <sub>(0.020)</sub>	0.884 <sub>(0.032)</sub>	0.881 <sub>(0.017)</sub>	0.826 <sub>(0.031)</sub>	0.742 <sub>(0.012)</sub>	—	—
GraphLog* [Xu et al., 2021a]	0.902 <sub>(0.007)</sub>	0.648 <sub>(0.015)</sub>	0.875 <sub>(0.009)</sub>	0.873 <sub>(0.022)</sub>	0.823 <sub>(0.017)</sub>	0.738 <sub>(0.0021)</sub>	—	—
MGSSL* [Zhang et al., 2021]	0.933 <sub>(0.009)</sub>	0.656 <sub>(0.011)</sub>	0.906 <sub>(0.019)</sub>	0.889 <sub>(0.006)</sub>	0.828 <sub>(0.015)</sub>	0.740 <sub>(0.018)</sub>	—	—
GROVER [Rong et al., 2020]	0.940 <sub>(0.019)</sub>	0.658 <sub>(0.023)</sub>	0.944 <sub>(0.021)</sub>	0.894 <sub>(0.028)</sub>	0.831 <sub>(0.025)</sub>	0.737 <sub>(0.010)</sub>	1.544 <sub>(0.397)</sub>	0.831 <sub>(0.120)</sub>
MPG [Li et al., 2021b]	0.922 <sub>(0.012)</sub>	<u>0.661</u> <sub>(0.007)</sub>	0.963 <sub>(0.028)</sub>	<u>0.920</u> <sub>(0.013)</sub>	<u>0.837</u> <sub>(0.019)</sub>	<u>0.748</u> <sub>(0.005)</sub>	<u>1.269</u> <sub>(0.192)</sub>	<u>0.741</u> <sub>(0.017)</sub>
MPG + MolAug	0.958 <sub>(0.022)</sub>	0.672 <sub>(0.025)</sub>	0.968 <sub>(0.017)</sub>	0.928 <sub>(0.07)</sub>	0.849 <sub>(0.022)</sub>	0.760 <sub>(0.010)</sub>	0.976 <sub>(0.028)</sub>	0.699 <sub>(0.052)</sub>
MPG + WordReg	0.975 <sub>(0.015)</sub>	0.676 <sub>(0.010)</sub>	<b>0.991</b> <sub>(0.012)</sub>	<b>0.946</b> <sub>(0.014)</sub>	0.857 <sub>(0.010)</sub>	0.762 <sub>(0.010)</sub>	0.988 <sub>(0.130)</sub>	<b>0.664</b> <sub>(0.048)</sub>
MPG + MolAug + WordReg	<b>0.978</b> <sub>(0.016)</sub>	<b>0.685</b> <sub>(0.012)</sub>	0.986 <sub>(0.021)</sub>	0.932 <sub>(0.016)</sub>	<b>0.862</b> <sub>(0.011)</sub>	<b>0.763</b> <sub>(0.008)</sub>	<b>0.934</b> <sub>(0.093)</sub>	0.695 <sub>(0.039)</sub>
Relative Improvement	+6.1%	+3.6%	+2.9%	+2.8%	+3.0%	+2.0%	+26.4%	+10.4%

- *Pretrain-then-finetune*: We also compare with several pre-trained models with vanilla fine-tuning: Mol2Vec [Jaeger et al., 2018], N-Gram [Liu et al., 2019], SMILES-BERT [Wang et al., 2019], GROVER, Hu et.al, MPG and several graph contrastive learning methods as introduced in the related work.

## 4.2 Comparisons with State-of-the-art Results

We first fine-tune the state-of-the-art pre-trained molecular graph model MPG with MolAug, WordReg and their combinations respectively. The results are reported in Table 1 which provides the following observations: (1) MolAug and WordReg consistently perform better than vanilla fine-tuning on MPG. The overall relative improvement is 7.2% (3.4% on classification tasks and 18.4% on regression tasks); (2) Equipped with MolAug, WordReg or their combinations, MPG yields new state-of-the-art results on molecular property prediction; (3) The combination of MolAug and WordReg brings a 26.4% relative improvement over MPG on a small dataset FreeSolv with only 642 labeled molecules. This confirms that our strategies can alleviate the over-fitting issues which is often severer in small-scaled datasets.

## 4.3 Comparisons with Other Fine-tuning Methods

Due to the fact that there is no fine-tuning methods specially designed for pre-trained molecular graph models, we consider some start-of-the-arts fine-tuning strategies in other domains. As can be observed in Table 2, generally, existing

Table 2: Comparison with other fine-tuning methods. The running time is evaluated on the same device (Tesla V100 GPU). We provide more details and results in the appendix.

Methods	BBBP		SIDER	
	ROC-AUC	Time (s)	ROC-AUC	Time (s)
Vanilla Fine-tuning (Base model)	0.922 <sub>(0.012)</sub>	654.6 <sub>(7.9)</sub>	0.661 <sub>(0.007)</sub>	327.6 <sub>(5.2)</sub>
SMART [Haoming et al., 2020]	0.959 <sub>(0.009)</sub>	9583.3 <sub>(26.2)</sub>	0.673 <sub>(0.005)</sub>	8927.5 <sub>(31.1)</sub>
R3F [Aghajanyan et al., 2021]	0.931 <sub>(0.008)</sub>	3876.5 <sub>(35.6)</sub>	0.659 <sub>(0.013)</sub>	2930.5 <sub>(16.7)</sub>
Adaptive Fine-tuning [Han and others., 2021]	0.944 <sub>(0.012)</sub>	8051.2 <sub>(8.2)</sub>	0.653 <sub>(0.012)</sub>	6982.3 <sub>(15.9)</sub>
ChildTuning [Xu et al., 2021b]	0.952 <sub>(0.005)</sub>	8541.4 <sub>(24.8)</sub>	0.668 <sub>(0.017)</sub>	7358.6 <sub>(22.8)</sub>
MolAug	0.958 <sub>(0.027)</sub>	789.8 <sub>(15.1)</sub>	0.672 <sub>(0.019)</sub>	405.6 <sub>(5.7)</sub>
WordReg	0.975 <sub>(0.015)</sub>	1659.8 <sub>(10.1)</sub>	0.676 <sub>(0.010)</sub>	1436.3 <sub>(20.3)</sub>
MolAug + WordReg	<b>0.978</b> <sub>(0.016)</sub>	2322.4 <sub>(16.8)</sub>	<b>0.685</b> <sub>(0.012)</sub>	1791.68 <sub>(14.9)</sub>

fine-tuning techniques can also bring improvements for Pre-trained molecular graph models. Compared with them, MolAug and WordReg are more suitable for pre-trained molecular graph models because they can avoid heavy tuning or heavy computational overhead of other fine-tuning methods and achieve the best performance among them.

## 4.4 Improving the Performance of Various Pre-trained Molecular Graph Models

To validate that MolAug and WordReg are model-agnostic, we fine-tune other pre-trained molecular graph models including MGSSL and GraphLog (see in the appendix for the limited space) with our MolAug and WordReg. Figure 4 shows the comparisons between our strategies and vanilla fine-tuning on the pre-trained model MGSSL. Obviously, our strategies outperform vanilla fine-tuning across various datasets, which validates that MolAug & WordReg

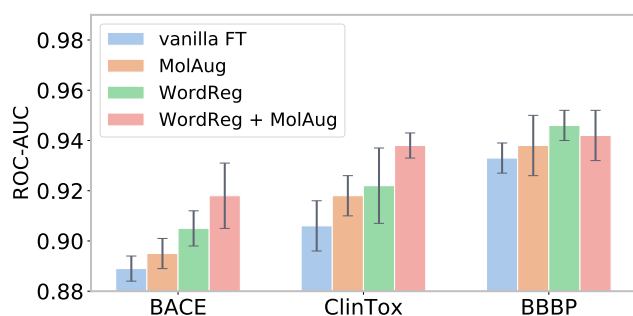


Figure 4: MolAug and WordReg’s performance (compared with vanilla fine-tuning) on another pre-trained model MGSSL.

can also improve the performance when fine-tuning various pre-trained molecular graph models.

## 4.5 Ablation Study

**Ablation Study for MolAug.** In this section, we substitute MolAug with general graph augmentations including nodes drop (ND), edges perturbation (EP) in GraphCL [You *et al.*, 2020] and domain knowledge-enriched molecular graph augmentations (DK) in MoCL [Sun *et al.*, 2021]. Besides, we also consider chemical enantiomers (Enan.) and homologies (Homo.) of MolAug in isolation. The results demonstrated in Figure 5 illustrate that general graph augmentations and DK bring negative or limited improvements over vanilla fine-tuning, which may alter molecular property completely. Besides, using Enan. and Homo. in isolation is not the optimal alternative. Instead, the combinations of these two ingredients of MolAug achieve the best performance.

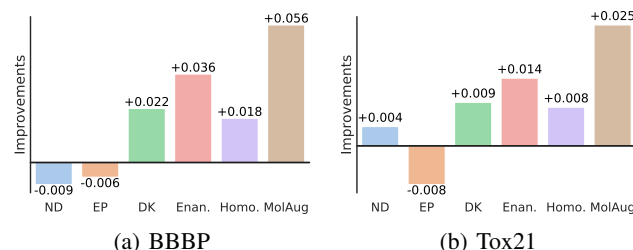


Figure 5: Comparisons between MolAug and existing molecular graph augmentations. The improvements refer to the improvements over vanilla fine-tuning.

**Ablation Study for WordReg.** We substitute the worst-case dropout with random dropout to study its influence. As shown in Figure 6, WordReg outperforms the random dropout across various drop ratios, which validates that looking for the worst-case dropout is necessary and effective. Figure 6(a) is symmetrical because both two dropouts are random.

## 4.6 Hyper-parameters Analysis

**Worst-case dropout ratio.** We have study dropout ratio in Figure 6(b). Obviously, WordReg can achieve brilliant performance when the two dropout ratios are in a reasonable

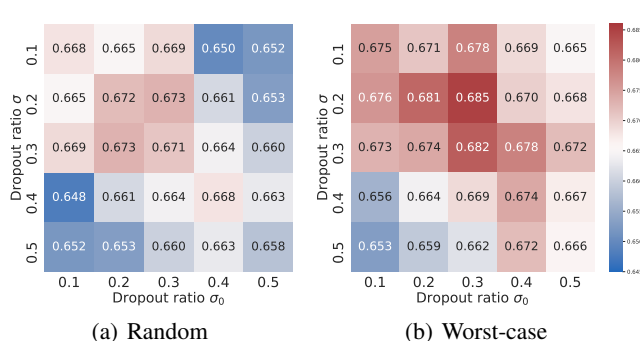


Figure 6: Comparisons between WordReg and random dropout regularization. The experiments are conducted with MPG base model on SIDER dataset. Here  $\sigma_0$  and  $\sigma$  are dropout ratio of vanilla training (fine-tuning) and worst-case drop, respectively.

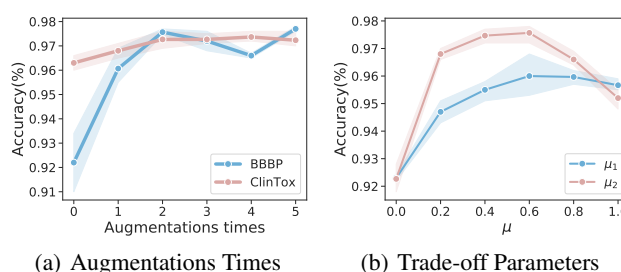


Figure 7: Hyper-parameters Analysis.

range (0.1-0.3). Besides, WordReg tends to perform better when the two dropout ratios are close or identical.

**Augmentations times of MolAug.** As shown in Figure 7(a), more augmentations for each molecule at each iteration bring improvements. However, the performance is prone to be steady when the augmentations times are up to 3. On the other hand, more augmentations times will lead to unnecessary computational overhead.

**Trade-off parameters  $\mu_1, \mu_2$ .** Figure 7(b) shows a different influence between  $\mu_1$  and  $\mu_2$ . For  $\mu_1$ , the accuracy will be improved and then tend to converge with its value increasing, which illustrates that MolAug can also work well even if its trade-off parameter is oversized. In contrast, for  $\mu_2$ , the accuracy sees a dramatical drop, which can be imputed to its over-powerful regularization.

## 5 Conclusions and Future Work

We propose two effective strategies, MolAug & WordReg, for fine-tuning pre-trained molecular graph models to alleviate the over-fitting issues. The empirical results suggest that our strategies can advance the performance of vanilla fine-tuning on various pre-trained models and outperform state-of-the-arts results by a large margin. Also, we validate that our strategies are superior to existing fine-tuning techniques. Interesting direction for future work is applying our strategies to other molecule-related tasks such as molecule generation.

## References

- [Aghajanyan *et al.*, 2021] Armen Aghajanyan, Akshat Shrivastava, Ankit Gupta, Naman Goyal, Luke Zettlemoyer, and Sonal Gupta. Better fine-tuning by reducing representational collapse. In *ICLR*, 2021.
- [Boyd *et al.*, 2003] Stephen Boyd, Lin Xiao, and others. Subgradient methods. 2003.
- [Chen *et al.*, 2021] Jianwen Chen, Shuangjia Zheng, and others. Learning attributed graph representation with communicative message passing transformer. *IJCAI*, 2021.
- [Devlin *et al.*, 2019] Jacob Devlin, Ming-Wei Chang, and others. Bert: Pre-training of deep bidirectional transformers for language understanding. *NAACL*, 2019.
- [Duvenaud *et al.*, 2015] David Duvenaud, Dougal Maclaurin, and others. Convolutional networks on graphs for learning molecular fingerprints. *NIPS*, 2015.
- [Gilmer and others., 2017] Justin Gilmer and others. Neural message passing for quantum chemistry. *ICML*, 2017.
- [Han and others., 2021] Xueting Han and others. Adaptive transfer learning on graph neural networks. In *KDD*, 2021.
- [Haoming *et al.*, 2020] Jiang Haoming, He Pengcheng, and others. Smart: Robust and efficient fine-tuning for pre-trained natural language models through principled regularized optimization. *ACL*, 2020.
- [Howard and Ruder, 2018] J. Howard and S. Ruder. Universal language model fine-tuning for text classification. In *ACL*, 2018.
- [Hu and others., 2020] Ziniu Hu and others. Gpt-gnn: Generative pre-training of graph neural networks. *KDD*, 2020.
- [Hu *et al.*, 2020] Weihua Hu, Bowen Liu, and others. Strategies for pre-training graph neural networks. *ICLR*, 2020.
- [Jaeger *et al.*, 2018] Sabrina Jaeger, Simone Fulle, and others. Mol2vec: Unsupervised machine learning approach with chemical intuition. *J CHEM INF MODEL*, 2018.
- [Kearnes *et al.*, 2016] Steven Kearnes, Kevin McCloskey, and others. Molecular graph convolutions: Moving beyond fingerprints. *J Comput Aid Mol Des*, 2016.
- [Kipf. and Welling, 2017] N. Thomas Kipf. and M. Welling. Semi-supervised classification with graph convolutional networks. *ICLR*, 2017.
- [Li *et al.*, 2021a] Pengyong Li, Yuquan Li, and others. Trimnet: learning molecular representation from triplet messages for biomedicine. *BIB*, 2021.
- [Li *et al.*, 2021b] Pengyong Li, Jun Wang, and others. An effective self-supervised framework for learning expressive molecular global representations to drug discovery. *BIB*, 2021.
- [Liu *et al.*, 2019] Shengchao Liu, Mehmet Demirel, and others. N-gram graph - simple unsupervised representation for graphs, with applications to molecules. *NeurIPS*, 2019.
- [Lu *et al.*, 2019] Chengqiang Lu, Qi Liu, and others. Molecular property prediction: A multilevel quantum interactions modeling perspective. *AAAI*, 2019.
- [Ma *et al.*, 2017] Xuezhe Ma, Yingkai Gao, and others. Dropout with expectation-linear regularization. 2017.
- [Mannhold *et al.*, 2012] Raimund Mannhold, Hugo Kubinyi, and others. *Bioisosteres in medicinal chemistry*. 2012.
- [Mohri *et al.*, 2012] Mehryar Mohri, Afshin Rostamizadeh, and others. Foundations of machine learning. 2012.
- [Ramsundar and others., 2015] B. Ramsundar and others. Massively multitask networks for drug discovery. 2015.
- [Rogers and Hahn, 2010] David Rogers and Mathew. Hahn. Extended-connectivity fingerprints. *J chem inf*, 2010.
- [Rong *et al.*, 2020] Yu Rong, Yatao Bian, and others. Self-supervised graph transformer on large-scale molecular data. *NIPS*, 2020.
- [Schuett *et al.*, 2017] T. K. Schuett, Kindermans, and others. Schnet: A continuous-filter convolutional neural network for modeling quantum interactions. *NIPS*, 2017.
- [Song *et al.*, 2020] Ying Song, Shuangjia Zheng, and others. Communicative representation learning on attributed molecular graphs. In *IJCAI*, 2020.
- [Sun *et al.*, 2021] Mengying Sun, Jing Xing, and others. Mocl: Contrastive learning on molecular graphs with multi-level domain knowledge. *KDD*, 2021.
- [Susheel *et al.*, 2021] Susheel, Pan Li, and others. Adversarial graph augmentation to improve graph contrastive learning. 2021.
- [Wang *et al.*, 2019] Sheng Wang, Yuzhi Guo, and others. Smiles-bert - large scale unsupervised pre-training for molecular property prediction. *BCB*, 2019.
- [Wu *et al.*, 2018] Zhenqin Wu, Bharath Ramsundar, and others. Moleculenet: a benchmark for molecular machine learning. *Chemical science*, 2018.
- [Xia *et al.*, 2021a] Jun Xia, Haitao Lin, Yongjie Xu, Lirong Wu, Zhangyang Gao, Siyuan Li, and Stan Z. Li. Towards robust graph neural networks against label noise, 2021.
- [Xia *et al.*, 2021b] Jun Xia, Lirong Wu, Jintao Chen, Ge Wang, and Stan Z Li. Debiased graph contrastive learning. *arXiv preprint arXiv:2110.02027*, 2021.
- [Xia *et al.*, 2022] Jun Xia, Lirong Wu, Jintao Chen, Bozhen Hu, and Stan Z. Li. SimGRACE: A Simple Framework for Graph Contrastive Learning without Data Augmentation. In *Proceedings of The Web Conference 2022*. Association for Computing Machinery, 2022.
- [Xiong *et al.*, 2020] Zhaoping Xiong, Dingyan Wang, and others. Pushing the boundaries of molecular representation for drug discovery with graph attention mechanism. *J Med Chem*, 2020.
- [Xu *et al.*, 2021a] Minghao Xu, Hang Wang, and others. Self-supervised graph-level representation learning with local and global structure. *ICML*, 2021.
- [Xu *et al.*, 2021b] Runxin Xu, Fuli Luo, and others. Raise a child in large language model: Towards effective and generalizable fine-tuning. *EMNLP*, 2021.

- [Yang *et al.*, 2019] Kevin Yang, Kyle Swanson, and others. Analyzing learned molecular representations for property prediction. *J CHEM INF MODEL*, 2019.
- [You *et al.*, 2020] Y. You, T. Chen, and others. Graph contrastive learning with augmentations. In *NeurIPS*, 2020.
- [You *et al.*, 2021] Yuning You, Tianlong Chen, and others. Graph contrastive learning automated. *ICML*, 2021.
- [Zhang *et al.*, 2021] Zaixi Zhang, Qi Liu, and others. Motif-based graph self-supervised learning for molecular property prediction. *NeurIPS*, 2021.
- [Zhu *et al.*, 2021a] Yanqiao Zhu, Yichen Xu, Qiang Liu, and Shu Wu. An empirical study of graph contrastive learning. 2021.
- [Zhu *et al.*, 2021b] Yanqiao Zhu, Yichen Xu, Feng Yu, Qiang Liu, Shu Wu, and Liang Wang. Graph contrastive learning with adaptive augmentation. In *Proceedings of the Web Conference 2021*, pages 2069–2080, 2021.
- [Zolna *et al.*, 2018] Konrad Zolna, Devansh Arpit, and others. Fraternal dropout. In *ICLR*, 2018.