

A MORE DETAILS OF PACIA

A.1 ENCODER

Encoder for MoleculeNet. As the main network of the encoder, to process a molecule with GNN, each node embedding \mathbf{h}_v represents an atom, and each edge e_{vu} represents a chemical bond. Here, we use GIN (Xu et al., 2019) as the main network in encoder, which is a powerful GNN structure. In GIN, the aggregation function in (1) is specified as adding all neighbors up: and for the update function is adding the aggregated embeddings and the target node, and feeding to a MLP:

$$\mathbf{h}_v^l = \text{MLP}_G^l \left((1 + \epsilon) \mathbf{h}_v^{l-1} + \sum_{u \in \mathcal{H}(v)} \mathbf{h}_u^{l-1} \right), \quad (14)$$

where ϵ is a scalar parameter to distinguish the target node. To obtain the molecular representation, the readout function in (2) is specified as

$$\mathbf{r} = \text{MLP}_R \left(\text{MEAN}(\{\mathbf{h}_v^L | v \in \mathcal{V}\}) \right). \quad (15)$$

Encoder for FS-Mol. Following existing works (Chen et al., 2022; Schimunek et al., 2023), we directly adopt the PNA (Corso et al., 2020) network provided in FS-Mol benchmark (Stanley et al., 2021) as the molecular encoder.

A.2 PREDICTOR

The classifier needs to make prediction of the query $\hat{\mathbf{h}}_{\tau,q}$, according to the N_τ labeled support samples $\{(\mathbf{h}_{\tau,s}, y_{\tau,s}) | \mathcal{X}_{\tau,s} \in \mathcal{S}_\tau\}$. We adopt an adaptive classifier (Requeima et al., 2019), which map the labeled samples in each class to the parameters of a linear classifier, i.e.,

$$\mathbf{w}_\pm = \text{MLP}_w \left(\frac{1}{|\mathcal{S}_\tau^\pm|} \sum_{\mathcal{X}_{\tau,s} \in \mathcal{S}_\tau^\pm} \mathbf{h}_{\tau,s} \right), \quad b_\pm = \text{MLP}_b \left(\frac{1}{|\mathcal{S}_\tau^\pm|} \sum_{\mathcal{X}_{\tau,s} \in \mathcal{S}_\tau^\pm} \mathbf{h}_{\tau,s} \right), \quad (16)$$

where \mathbf{w}_\pm has the same dimension with $\mathbf{h}_{\tau,q}$, and b_\pm is a scalar. Then the prediction is made by

$$\hat{y}_{\tau,q} = \text{softmax}([\mathbf{w}_-^\top \mathbf{h}_{\tau,q} + b_-, \mathbf{w}_+^\top \mathbf{h}_{\tau,q} + b_+]), \quad (17)$$

where $\text{softmax}(\mathbf{x}) = \exp(\mathbf{x}) / \sum_i \exp([\mathbf{x}]_i)$ and $[\mathbf{x}]_i$ means the i th element in \mathbf{x} .

A.3 UNIFIED GNN ADAPTER

The choice of e can be various (Wu et al., 2023), while in this work we adopt a simple feature-wise linear modulation (FiLM) (Perez et al., 2018) function.

A.4 HYPERPARAMETERS

Here we provide the detailed hyperparameter setting of PACIA.

Hyperparameters on MoleculeNet. The maximum layer number of the GNN $L_{\text{enc}} = 5$, the maximum depth of the relation graph $L_{\text{rel}} = 5$, During training, for each layer in GNN, we set dropout rate as 0.5 operated between the graph operation and FiLM layer. The dropout rate of MLP in (9) (10) and (11) is 0.1. For all baselines, we use Adam optimizer (Kingma & Ba, 2015) with learning rate 0.006 and the maximum episode number is 25000. In each episode, the meta-training tasks are learned one-by-one, query set size $M = 16$. The ROC-AUC is evaluated every 10 epochs on meta-testing tasks and the best performance is reported. Table 4 shows the details of the other parts. Experiments are conducted on a 24GB NVIDIA GeForce RTX 3090 GPU, with Python 3.8.13, CUDA version 11.7, Torch version 1.10.1.

Hyperparameters on FS-Mol. The maximum layer number of the GNN $L_{\text{enc}} = 8$, the maximum depth of the relation graph $L_{\text{rel}} = 5$, During training, the dropout rate of MLP_L is 0.1. We use Adam optimizer (Kingma & Ba, 2015) with learning rate 0.0001 and the maximum episode number is 3000. In each episode, the meta-training tasks are learned with batch size 16, support set size $N_\tau = 64$, and the others are used as queries. The average precision s evaluated every 50 epochs on validation tasks and the the model with best validation performance is tested and reported. Table 5 shows the details of the other parts.

Table 4: Details of model structure for MoleculeNet.

	Layers	Output Dimension
MLP in (9)	input $\frac{1}{ \mathcal{V}_{\tau,s} } \sum_{v \in \mathcal{X}_{\tau,s}} [\mathbf{h}_v^l \mathbf{y}_{\tau,s}]$, fully connected, LeakyReLU	300
	fully connected with with residual skip connection, $\frac{1}{K} [\sum_{\mathcal{X}_{\tau,s} \in \mathcal{S}_\tau^+} (\cdot) \sum_{\mathcal{X}_{\tau,s} \in \mathcal{S}_\tau^-} (\cdot)]$	300
MLP in (10)	fully connected with residual skip connection	601
MLP in (11)	fully connected with residual skip connection	257
MLP _G ^l in (14)	input $(1 + \epsilon)\mathbf{h}_v^{l-1} + \mathbf{h}_{\text{agg}}^{l-1}$, fully connected, ReLU	600
	fully connected	300
MLP _L in (15)	input READOUT ($\{\mathbf{h}_v^T v \in \mathcal{V}_{t,i}\}$), fully connected, LeakyReLU	128
	fully connected	128
MLP in (4)	input $\exp(\mathbf{h}_{\tau,i}^{l-1} - \mathbf{h}_{\tau,i}^{l-1})$, fully connected, LeakyReLU	256
	fully connected, LeakyReLU	128
	fully connected	1
MLP in (5)	fully connected, LeakyReLU	256
	fully connected, LeakyReLU	128
MLP _w in (16)	input $\frac{1}{K} \sum_{y_{\tau,s}=c} \mathbf{h}_{\tau,s}$, fully connected with residual skip connection, LeakyReLU	128
	2× (fully connected with residual skip connection, LeakyReLU)	128
	fully connected	128
MLP _b in (16)	input $\frac{1}{K} \sum_{y_{\tau,s}=c} \mathbf{h}_{\tau,s}$, fully connected with residual skip connection, LeakyReLU	128
	2× (fully connected with residual skip connection, LeakyReLU)	128
	fully connected	1

Table 5: Details of model structure for FS-Mol.

	Layers	Output Dimension
MLP in (9)	input $\frac{1}{ \mathcal{V}_{\tau,s} } \sum_{v \in \mathcal{X}_{\tau,s}} [\mathbf{h}_v^l \mathbf{y}_{\tau,s}]$, fully connected, LeakyReLU	512
	fully connected with with residual skip connection, $[\frac{1}{ \mathcal{S}_\tau^+} \sum_{\mathcal{X}_{\tau,s} \in \mathcal{S}_\tau^+} (\cdot) \frac{1}{ \mathcal{S}_\tau^-} \sum_{\mathcal{X}_{\tau,s} \in \mathcal{S}_\tau^-} (\cdot)]$	512
MLP in (10)	fully connected with residual skip connection	1025
MLP in (11)	fully connected with residual skip connection	513
MLP in (4)	input $\exp(\mathbf{h}_{\tau,i}^{l-1} - \mathbf{h}_{\tau,i}^{l-1})$, fully connected, LeakyReLU	256
	fully connected, LeakyReLU	128
	fully connected	1
MLP in (5)	fully connected, LeakyReLU	256
	fully connected, LeakyReLU	256
MLP _w in (16)	input $\frac{1}{ \mathcal{S}_\tau^+} \sum_{y_{\tau,s}=c} \mathbf{h}_{\tau,s}$, fully connected with residual skip connection, LeakyReLU	256
	2× (fully connected with residual skip connection, LeakyReLU)	256
	fully connected	256
MLP _b in (16)	input $\frac{1}{ \mathcal{S}_\tau^+} \sum_{y_{\tau,s}=c} \mathbf{h}_{\tau,s}$, fully connected with residual skip connection, LeakyReLU	256
	2× (fully connected with residual skip connection, LeakyReLU)	256
	fully connected	1

B ADOPTING MAML FOR PROPERTY-LEVEL ADAPTATION

Denote all model parameter as Θ . The model first predict samples in support set and get loss to do local-update. Denote the loss for local-update as $\mathcal{L}_\tau^S(\Theta) = \sum_{\mathcal{X}_{\tau,s} \in \mathcal{S}_\tau} \mathbf{y}_{\tau,s}^\top \log(\hat{\mathbf{y}}_{\tau,s})$, where $\hat{\mathbf{y}}_{\tau,s}$ is the prediction made by the main network with parameter Θ . The loss for global-update is

calculated with samples in query set, denoted as $\mathcal{L}_\tau^Q(\Theta'_\tau) = \sum_{\mathcal{X}_{\tau,q} \in \mathcal{Q}_\tau} \mathbf{y}_{\tau,q}^\top \log(\hat{\mathbf{y}}_{\tau,q})$, where $\hat{\mathbf{y}}_{\tau,q}$ is the prediction made by the main network with parameter Θ'_τ . Then Algorithm 3 can be adopted for meta-training, and Algorithm 4 for meta-testing.

Algorithm 3 Meta-training with MAML

Input: meta-training task set $\mathcal{T}_{\text{train}}$

- 1: initialize Θ randomly;
 - 2: **while** not done **do**
 - 3: **for** each task $\mathcal{T}_\tau \in \mathcal{T}_{\text{train}}$ **do**
 - 4: evaluate $\nabla_{\Theta} \mathcal{L}_\tau^S(\Theta)$ with respect to all samples in \mathcal{S}_τ ;
 - 5: compute adapted parameters with gradient descent: $\Theta'_\tau = \Theta - \nabla_{\Theta} \mathcal{L}_\tau^S(\Theta)$;
 - 6: **end for**
 - 7: update $\Theta \leftarrow \Theta - \nabla_{\Theta} \sum_{\mathcal{T}_\tau \in \mathcal{T}_{\text{train}}} \mathcal{L}_\tau^Q(\Theta'_\tau)$;
 - 8: **end while**
 - 9: **return** learned Θ^* .
-

Algorithm 4 Meta-testing with MAML

Input: learned Θ^* , a meta-testing task \mathcal{T}_τ ;

- 1: evaluate $\nabla_{\Theta} \mathcal{L}_\tau^S(\Theta)$ with respect to all samples in \mathcal{S}_τ ;
 - 2: compute adapted parameters with gradient descent: $\Theta'_\tau = \Theta - \nabla_{\Theta} \mathcal{L}_\tau^S(\Theta)$;
 - 3: make prediction $\mathbf{y}_{\tau,q}$ for $\mathcal{X}_{\tau,q} \in \mathcal{Q}_\tau$ with adapted parameter Θ'_τ ;
-

B.1 COMPARISON WITH EXISTING WORKS

We compare the proposed PACIA with existing few-shot MPP approaches in Table 6. As shown, we manage to compare in perspectives of support of pre-training, property-level adaptation, molecule-level adaptation, fast-adaptation and adaptation strategy. With the help of hypernetworks, our method not only introduces novel molecule-level adaptation, but also can adapt on property-level more effectively and efficiently.

Table 6: Comparison of the proposed PACIA with existing few-shot MPP methods.

Method	Support	Hierarchical adaptation		Fast adaptation	Adaptation Strategy
	Pre-training	Property-level	molecule-level		
IterRefLSTM	×	✓	×	✓	Pair-wise similarity
Meta-MGNN	✓	✓	×	×	Gradient
PAR	✓	✓	×	×	Attention+Gradient
ADKF-IFT	✓	✓	×	×	Gradient+statistical learning
MHNfs	✓	✓	✓	✓	Attention+pair-wise similarity
GS-META	✓	×	✓	✓	Message passing
PACIA	✓	✓	✓	✓	Hypernetwork

From the perspective of hypernetwork, the usage of the hypernetwork for encoder is related to GNN-FiLM (Brockschmidt, 2020), which considers a GNN as main network. It builds hypernetwork with target node as input to generate parameters of FiLM layers, to equip different nodes with different aggregation functions in the GNN. What and how to adapt are similar to ours, but it is different that the input of our hypernetwork for encoder is \mathcal{S}_τ and how we encode a set of labeled graphs.

In few-shot learning, some recent works (Requeima et al., 2019; Lin et al., 2021) use hypernetworks to process the task context to make the model task-adaptive. Their hypernetworks have similar functionality of our hypernetwork for encoder, that are used to map the support set to parameter to modulate the main network, but their main networks are convolutional neural network (CNN) and MLP respectively, there is significant difference about what to modulate. As for the usage of the hypernetwork for predictor, which is used to evaluate an unlabeled sample with a set of labeled ones to encode model architecture, did not appear in the literature.

C MORE DETAILS OF EXPERIMENTS

C.1 DATASETS

MoleculeNet. There are four sub-datasets for few-shot MPP: Tox21 (National Center for Advancing Translational Sciences, 2017), SIDER (Kuhn et al., 2016), MUV (Rohrer & Baumann, 2009) and ToxCast (Richard et al., 2016), which are included in MoleculeNet (Wu et al., 2018). We adopt the task splits provided by existing works (Altae-Tran et al., 2017; Wang et al., 2021). Tox21 is a collection of nuclear receptor assays related to human toxicity, containing 8014 compounds in 12 tasks, among which 9 are split for training and 3 are split for testing. SIDER collects information about side effects of marketed medicines, and it contains 1427 compounds in 21 tasks, among which 21 are split for training and 6 are split for testing. MUV contains compounds designed to be challenging for virtual screening for 17 assays, containing 93127 compounds in 17 tasks, among which 12 are split for training and 5 are split for testing. ToxCast collects compounds with toxicity labels, containing 8615 compounds in 617 tasks, among which 450 are split for training and 167 are split for testing.

FS-Mol. FS-Mol benchmark, which contains a set of few-shot learning tasks for molecular property prediction carefully collected from ChEMBL27 (Mendez et al., 2019) by Stanley et al. (2021). Following existing works (Chen et al., 2022; Schimunek et al., 2023), we use the same 10% of all tasks which contains 233,786 unique compounds, split into training (4,938 tasks), validation (40 tasks), and test (157 tasks) sets. Each task is associated with a protein target.

C.2 BASELINES

We compare our method with following baselines:

- **Siamese** (Koch et al., 2015): It learns two neural networks which are symmetric on structure to identify whether the input molecule pairs are from the same class. The performance is copied from (Altae-Tran et al., 2017) due to the lack of code.
- **ProtoNet**¹ (Snell et al., 2017): It makes classification according to inner-product similarity between the target and the prototype of each class. This method is incorporated as a classifier after the GNN encoder.
- **MAML**² (Finn et al., 2017): It learns a parameter initialization and the model is adapted to each task via few gradient steps on the support set. We adopt this method for all parameters in a model composed of a GNN encoder and a linear classifier.
- **EGNN**³ (Kim et al., 2019): It builds a relation graph that samples are refined, and it learns to predict edge-labels in the relation graph. This method is incorporated as the predictor after the GNN encoder.
- **GNN-FiLM** (Brockschmidt, 2020): GNN-FiLM has built hypernetwork with target node as input to generate parameters of FiLM layers, which equips different nodes with different aggregation functions in the GNN. We adopt this as encoder and a MLP as classifier and train on all samples in meta-training tasks and samples in support set of all meta-testing tasks;
- **IterRefLSTM** (Altae-Tran et al., 2017): It introduces matching networks combined with long short-term memory (LSTM) to refine the molecular representations according to the task context. The performance is copied from (Altae-Tran et al., 2017) due to the lack of code.
- **PAR**⁴ (Wang et al., 2021): It introduces an attention mechanism to capture task-dependent property and an inductive relation graph between samples, and incorporates MAML to train.
- **ADKF-IFT**⁵ (Wang et al., 2021): It adopt gradient-based strategy to learn the encoder where it proposes Implicit Function Theory to avoid computing the hyper-gradient. And a Gaussian Process is learned from scratch in each task as classifier.

¹<https://github.com/jakesnell/prototypical-networks>

²<https://github.com/learnables/learn2learn>

³<https://github.com/khy0809/fewshot-egnn>

⁴<https://github.com/tatal661/PAR-NeurIPS21>

⁵<https://github.com/Wenlin-Chen/ADKF-IFT>

- **Pre-GNN**⁶ (Hu et al., 2019): It trains a GNN encoder on ZINC15 dataset with graph-level and node-level self-supervised tasks, and fine-tunes the pre-trained GNN on downstream tasks. We adopt the pre-trained GNN encoder and a linear classifier.
- **GraphLoG**⁷ (Xu et al., 2021): It introduces hierarchical prototypes to capture the global semantic clusters. And adopts an online expectation-maximization algorithm to learn. We adopt the pre-trained GNN encoder and a linear classifier.
- **MGSSL**⁸ (Hu et al., 2019): It trains a GNN encoder on ZINC15 dataset with graph-level, node-level and motif-level self-supervised tasks, and fine-tunes the pre-trained GNN on downstream tasks. We adopt the pre-trained GNN encoder and a linear classifier.
- **GraphMAE**⁹ (Hou et al., 2022): It presents a masked graph autoencoder for generative self-supervised graph pre-training and focus on feature reconstruction with both a masking strategy and scaled cosine error. We adopt the pre-trained GNN encoder and a linear classifier.
- **Meta-MGNN**¹⁰ (Guo et al., 2021): It incorporates self-supervised tasks such as bond reconstruction and atom type prediction to be jointly optimized via MAML. It uses the pre-trained GNN encoder provided by (Hu et al., 2019).
- **Pre-PAR**: The same as PAR but uses the pre-trained GNN encoder provided by (Hu et al., 2019).
- **Pre-ADKF-IFT**: The same as ADKF-IFT but uses the pre-trained GNN encoder provided by (Hu et al., 2019).

C.3 PERFORMANCE COMPARISON WITH PRE-TRAINING

Baselines with Pre-training. We compare with the following baselines with (w/) pre-training: (i) Methods which fine-tune pre-train GNN encoders, including **Pre-GNN** (Hu et al., 2019), **GraphLoG** (Xu et al., 2021), **MGSSL** (Zhang et al., 2021), **GraphMAE** (Hou et al., 2022); (ii) Few-shot MPP methods incorporating pre-trained encoders provided by (Hu et al., 2019), including **Meta-MGNN** (Guo et al., 2021), **Pre-PAR** (Wang et al., 2021) and **Pre-ADKF-IFT** (Chen et al., 2022). All encoders have the same structure (Hu et al., 2019) and are pre-trained on ZINC15 dataset (Sterling & Irwin, 2015). We equip our PACIA with the same pre-trained encoder, and name it as **Pre-PACIA**.

Table 7: Test ROC-AUC obtained with pre-trained GNN encoder.

Method	Tox21		SIDER		MUV		ToxCast	
	10-shot	1-shot	10-shot	1-shot	10-shot	1-shot	10-shot	1-shot
Pre-GNN	83.02 _(0.13)	82.75 _(0.09)	77.55 _(0.14)	67.34 _(0.30)	67.22 _(2.16)	65.79 _(1.68)	73.03 _(0.67)	71.26 _(0.85)
GraphLoG	81.61 _(0.35)	79.23 _(0.93)	75.18 _(0.27)	67.52 _(1.40)	67.83 _(1.65)	66.56 _(1.46)	73.92 _(0.15)	73.10 _(0.39)
MGSSL	83.24 _(0.09)	<u>83.21</u> _(0.12)	77.87 _(0.18)	69.66 _(0.21)	68.58 _(1.32)	66.93 _(1.74)	73.51 _(0.45)	72.89 _(0.63)
GraphMAE	84.01 _(0.27)	81.54 _(0.18)	76.07 _(0.15)	67.60 _(0.38)	67.99 _(1.28)	<u>67.50</u> _(2.12)	74.15 _(0.33)	72.67 _(0.71)
Meta-MGNN	83.44 _(0.14)	82.67 _(0.20)	77.84 _(0.34)	74.62 _(0.41)	68.31 _(3.06)	66.10 _(3.98)	74.69 _(0.57)	73.29 _(0.85)
Pre-PAR	84.95 _(0.24)	<u>83.01</u> _(0.28)	<u>78.05</u> _(0.15)	<u>75.29</u> _(0.32)	69.88 _(1.57)	66.96 _(2.63)	75.48 _(0.99)	<u>73.90</u> _(1.21)
Pre-ADKF-IFT	<u>86.06</u> _(0.35)	80.97 _(0.48)	70.95 _(0.60)	62.16 _(1.03)	95.74 _(0.37)	67.25 _(3.87)	76.22 _(0.13)	71.13 _(1.15)
Pre-PACIA	86.40 _(0.27)	84.35 _(0.14)	83.97 _(0.22)	80.70 _(0.28)	<u>73.43</u> _(1.96)	69.26 _(2.35)	76.22 _(0.73)	75.09 _(0.95)

Performance with Pre-training. Table 7 shows the results. We can see that Pre-PACIA obtains significantly better performance except the 10-shot case on MUV, surpassing the second-best method Pre-ADKF-IFT by 3.10%. MGSSL defeats the other methods which fine-tune pre-trained GNN encoders, i.e., Pre-GNN, GraphLoG, and GraphMAE. However, it still performs worse than Pre-PACIA equipped with Pre-GNN, which validates the necessity of designing a few-shot MPP method

⁶<http://snap.stanford.edu/gnn-pretrain>

⁷<http://proceedings.mlr.press/v139/xu21g/xu21g-suppl.zip>

⁸<https://github.com/zaixizhang/MGSSL>

⁹<https://github.com/THUDM/GraphMAE>

¹⁰<https://github.com/zhichunguo/Meta-Meta-MGNN>

instead of simply fine-tuning a pre-trained GNN encoder. Moreover, comparing Pre-PACIA and PACIA in Table 1, the pre-trained encoder brings 3.05% improvement in average performance due to a better starting point of learning.

C.4 A CLOSER LOOK AT MOLECULE-LEVEL ADAPTATION

In this section, we pay a closer look at our molecule-level adaptation mechanism, proving evidence of its effectiveness.

C.4.1 PERFORMANCE UNDER DIFFERENT PROPAGATION DEPTH

Figure 5 compares Pre-PACIA with “w/o M” (introduced in Section 5.3) using different fixed layers of relation graph refinement on Tox21, where the maximum Depth $L = 5$. As can be seen, Pre-PACIA equipped performs much better than “w/o M” which takes the same depth of relation graph refinement as in PAR. This validates the necessity of molecule-level adaptation.

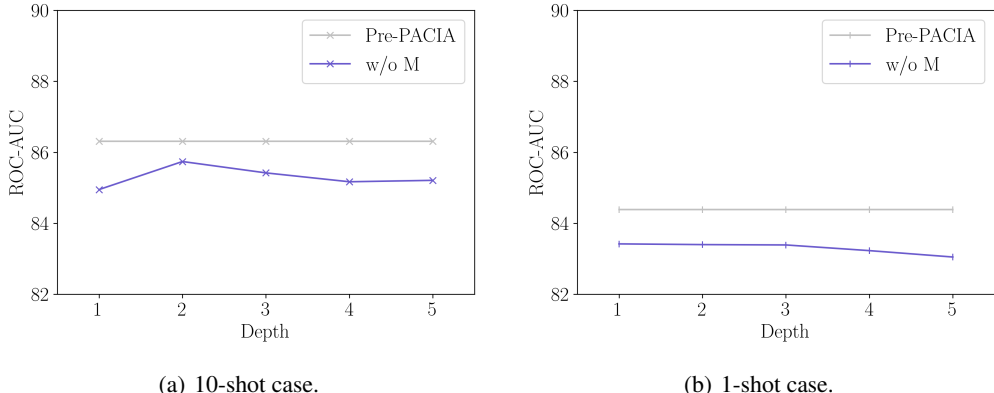


Figure 5: Comparing Pre-PACIA with “w/o M” using different fixed propagation depth of relation graph on Tox21.

C.4.2 DISTRIBUTION OF PROPAGATION DEPTH

Figure 6(a) plots the distribution of learned l' for query molecules in meta-testing tasks for 10-shot case of Tox21. The three meta-testing tasks contain different number of query molecules in scale: 6447 in task SR-HSE, 5790 in task SR-MPP, and 6754 in task SR-p53. We can see that Pre-PACIA choose different l' for query molecules in the same task. Besides, the distribution of learned l' varies across different meta-testing tasks: molecules in task SR-MPP mainly choose smaller depth while molecules in the other two tasks tend to choose greater depth. This can be explained as most molecules in task SR-MPP are relatively easy to classify, which is consistent with the fact that Pre-PACIA obtains the highest ROC-AUC on SR-MPP among the three meta-testing tasks (83.75 for SR-HSE, 88.79 for SR-MPP and 86.39 for SR-p53).

Further, we pick out molecules with $l' = 1$ (denote as **Group A**) and $l' = 4$ (denote as **Group B**) as they are more extreme cases. We then apply “w/o M” with different fixed depth for Group A and Group B, and compare them with Pre-PACIA. Figure 6(b) shows the results. Different observations can be made for these two groups. Molecules in Group A have good performance with smaller depth relation graph, they can achieve higher ROC-AUC score than the average of all molecules using Pre-PACIA. These indicate they are easier to classify and it is reasonable that Pre-PACIA choose $l' = 1$ for them. While molecules in Group B are harder to classify and requires $l' = 4$.

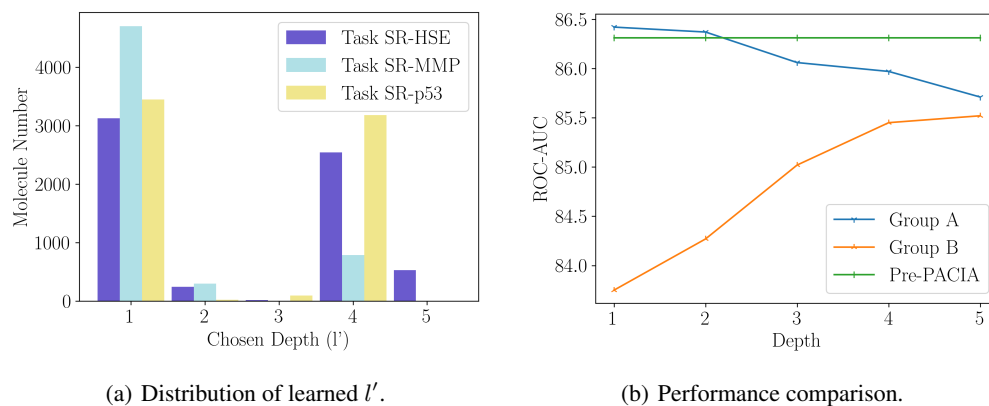


Figure 6: Examine molecule-level adaptation of Pre-PACIA on 10-shot tasks of Tox21.