# From Coarse to Fine: Enable Comprehensive Graph Self-supervised Learning with Multi-granular Semantic Ensemble

Author: Qianlong Wen, Mingxuan Ju, Zhongyu Ouyang, Chuxu Zhang, Yanfang Ye

# CONTENTS

**1** **Motivation**

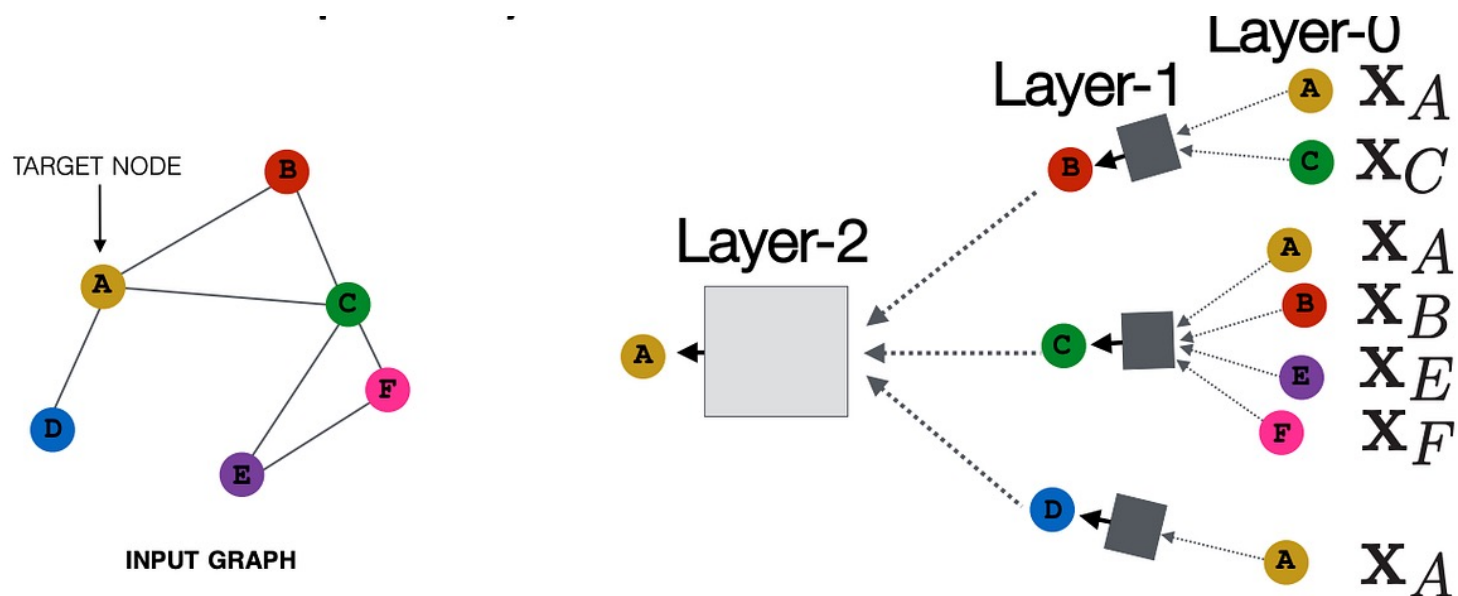**2** **Challenge**

**3** **Methodology**

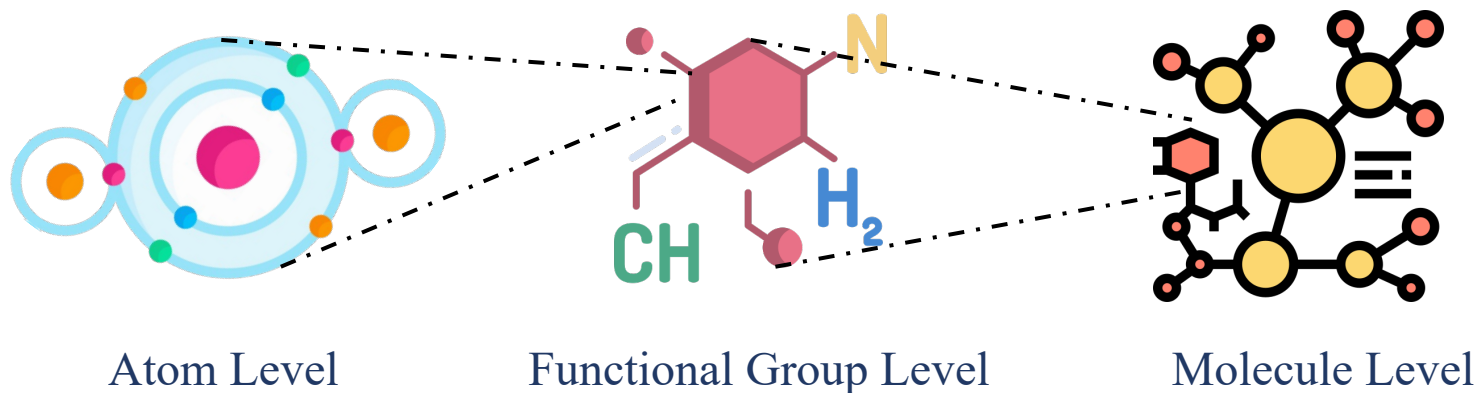**4** **Experiment**

## Graph Neural Networks (GNNs)

$$f_w(X, A) = \text{Softmax}(\hat{A} \ \text{ReLu}(\hat{A} \ X \ W_1) \ W_2),$$



The aggregated local neighborhood can usually represent the graph semantics.

## Graph Semantics (Substructures) in Different Granularities

There could exist different underlying patterns within a single graph, which makes it challenging to train a single model to fits well on the data. And those knowledge could be in diversifying granularities, which further increase the difficulties.

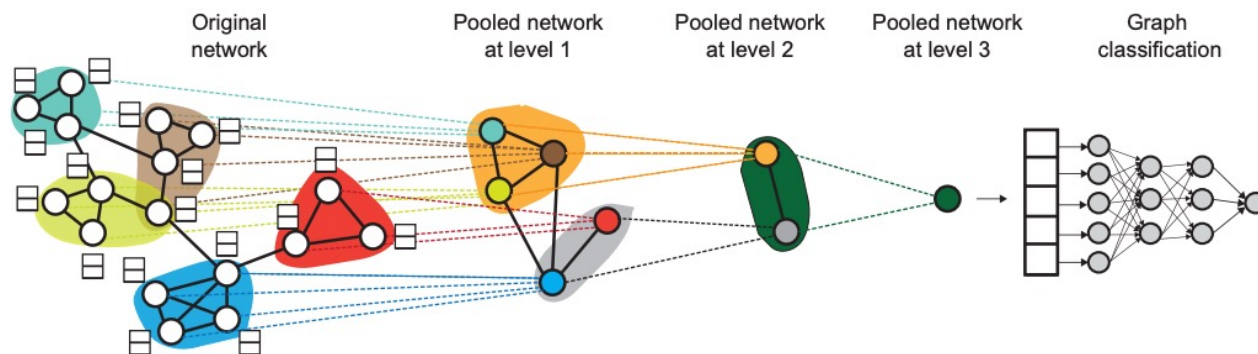

Atom Level      Functional Group Level      Molecule Level

Each level can used for the prediction tasks of different properties. Therefore, we need a model with larger capacity and higher expressiveness power on graph substructure learning so that it can extract the knowledge in different levels for better generalization ability.
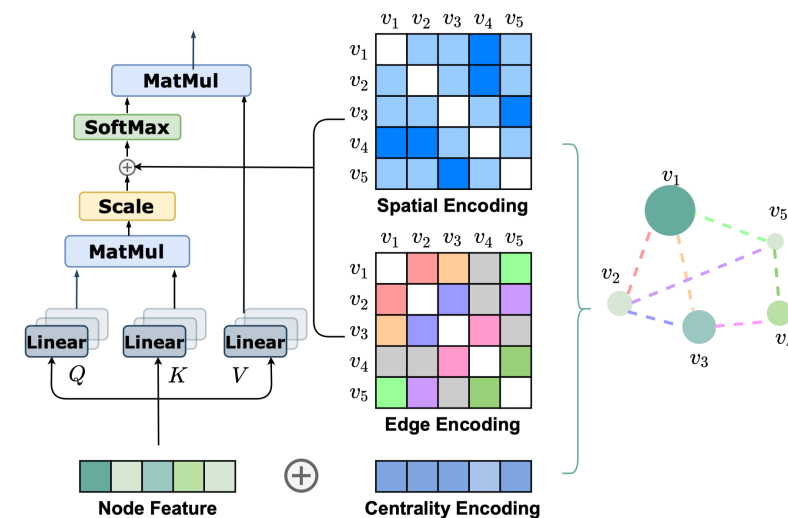
# How to Improve the Expressiveness Power on Graph Substructures?

Basically, current research on improving the expressive power of deep graph models on graph substructure learning mainly focus on modifying model architectures to improve WL-Test,

- Hierarchical pooling operations: SubGNN, DiffPool and etc.

- Transformer backbone: Graphormer, SAN, GraphGPS and etc.



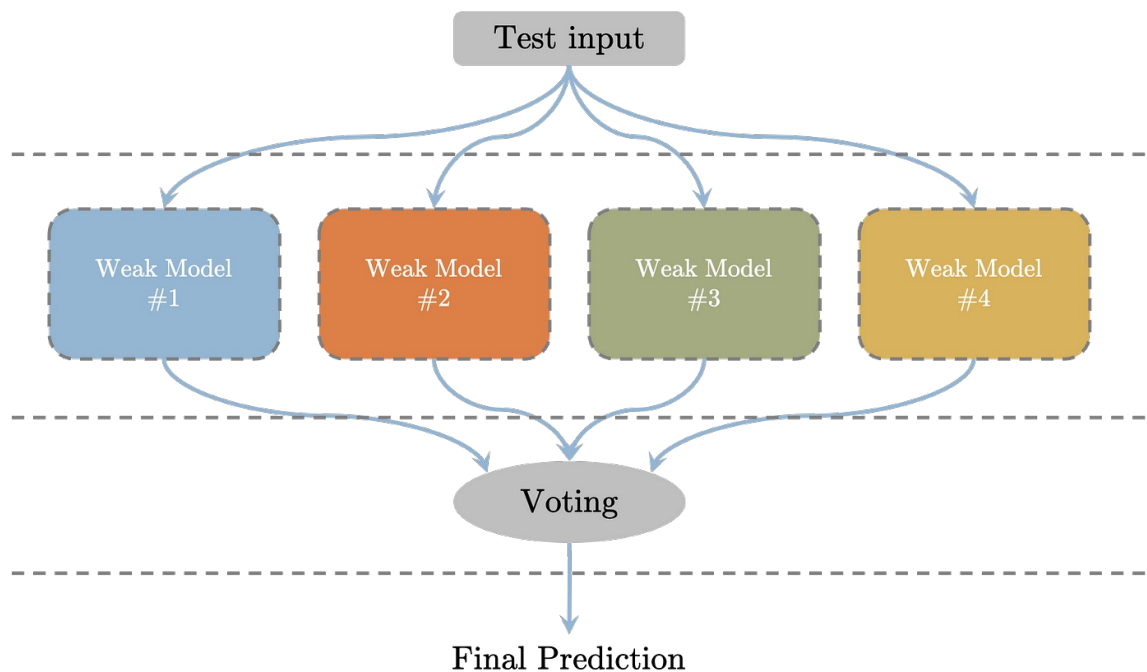*[NeurIPS 2018 ] Hierarchical Graph Representation Learning with Differentiable Pooling.*

*[NeurIPS 2021]Do Transformers Really Perform Bad for Graph Representation?*

## How to Improve the Expressiveness Power on Graph Substructures?

Meanwhile, another solution could be combining existing deep graph models into the ensemble learning and train each classifier to learn the graph semantics in different levels.

# Challenges on Incorporating Ensemble Learning

■ Higher Training and Inference Cost: Ensemble learning will naturally introduce extra computation cost, so the proposed framework should not be too expensive for training and inference.

■ Loading Balance: Ensuring that all classifiers are utilized effectively can be challenging, as some of them may become over-utilized while others remain under-utilized.

■ Regularization on Multi-granularity: Each classifier is included to learn knowledge in different granularities, so regularization should be applied to achieve this goal

# MGSE Framework



$$\mathcal{L} = \frac{1}{KN}\sum_{k=1}^{K}\sum_{i=1}^{N}\left[\mathcal{D}\left(\mathbf{p}_{i,k}^{t},\mathbf{p}_{i,k}^{s}\right) - \lambda\mathcal{H}\left(\bar{\mathbf{p}}_{k}^{s}\right)\right] \quad \mathcal{D}\left(\mathbf{p}_{k}^{t},\mathbf{p}_{k}^{s}\right) = -\mathcal{H}\left(\mathbf{p}_{k}^{t},\mathbf{p}_{k}^{s}\right) = \frac{1}{ND_{k}}\sum_{i=1}^{N}\sum_{d=1}^{D_{k}} -\left[\mathbf{p}_{i,k}^{t}\right]_{d}\cdot\log\left[\mathbf{p}_{i,k}^{s}\right]_{d}.$$

Using knowledge distillation for quick knowledge transfer and model converge]
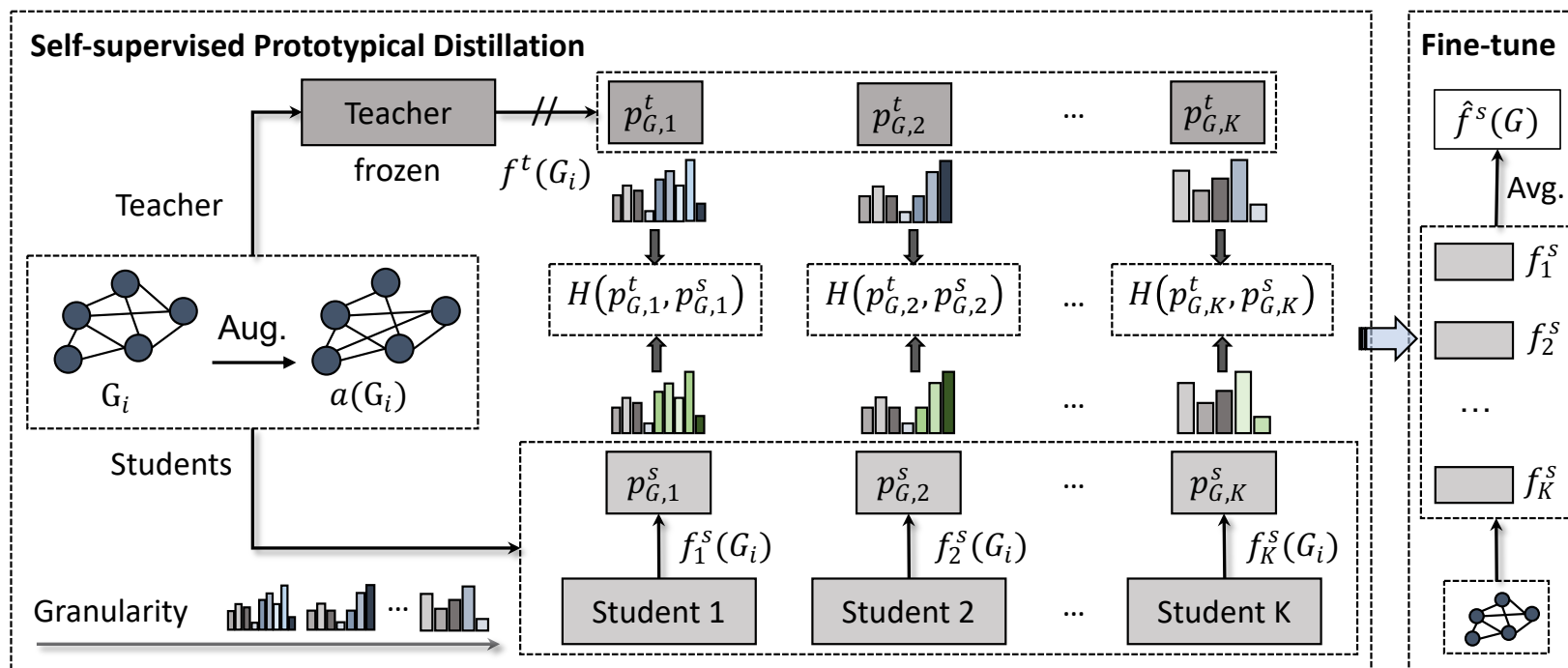
# MGSE Framework



$$\mathcal{L} = \frac{1}{KN} \sum_{k=1}^{K} \sum_{i=1}^{N} \left[ \mathcal{D}\left(\mathbf{p}_{i,k}^t, \mathbf{p}_{i,k}^s\right) - \lambda \mathcal{H}\left(\overline{\mathbf{p}}_k^s\right) \right] \qquad \mathcal{D}\left(\mathbf{p}_k^t, \mathbf{p}_k^s\right) = -\mathcal{H}\left(\mathbf{p}_k^t, \mathbf{p}_k^s\right) = \frac{1}{ND_k} \sum_{i=1}^{N} \sum_{d=1}^{D_k} -\left[\mathbf{p}_{i,k}^t\right]_d \cdot \log \left[\mathbf{p}_{i,k}^s\right]_d .$$

Manually setting different prototype number for $K$
student models to model the multi-granularity

# MGSE Framework



$$\mathcal{L} = \frac{1}{KN} \sum_{k=1}^{K} \sum_{i=1}^{N} \left[ \mathcal{D}\left(\mathbf{p}_{i,k}^{t}, \mathbf{p}_{i,k}^{s}\right) - \boxed{\lambda \mathcal{H}\left(\overline{\mathbf{p}}_{k}^{s}\right)} \right] \qquad \mathcal{H}\left(\overline{\mathbf{p}}_{k}^{s}\right) = \left[ \frac{1}{N} \sum_{i=1}^{N} \mathbf{p}_{i,k}^{s} \right] \cdot \log \left[ \frac{1}{N} \sum_{i=1}^{N} \mathbf{p}_{i,k}^{s} \right]$$

Adding ME-MAX regularization for loading balance

# Experiment - Dataset

Table 4. Statistics of MoleculeNet datasets and protein-protein interaction network datasets.

| Dataset | #Graphs | Avg #Nodes | Avg Degree | #Tasks (Class) | Metric | Category |
|---------|---------|------------|------------|----------------|--------|----------|
| ZINC15 | 2,000,000 | 26.62 | 57.72 | - | - | biochemical |
| PPI-306K | 306925 | 39.82 | 729.62 | - | - | Protein-Protein Intersection Networks |
| BBBP | 2,039 | 24.06 | 51.90 | 1 | ROC-AUC | biochemical |
| Tox21 | 7,813 | 18.57 | 38.58 | 12 | ROC-AUC | biochemical |
| ToxCast | 8,576 | 18.78 | 38.62 | 617 | ROC-AUC | biochemical |
| SIDER | 1,427 | 33.64 | 70.71 | 27 | ROC-AUC | biochemical |
| ClinTox | 1,477 | 26.15 | 55.76 | 2 | ROC-AUC | biochemical |
| MUV | 93,087 | 24.23 | 52.55 | 17 | ROC-AUC | biochemical |
| HIV | 41,127 | 25.51 | 54.93 | 1 | ROC-AUC | biochemical |
| BACE | 1,513 | 34.08 | 73.71 | 1 | ROC-AUC | biochemical |
| PPI | 88000 | 49.35 | 890.77 | 40 | ROC-AUC | Protein-Protein Intersection Networks |
| Cora | 1 | 2,708 | 5,429 | 7 | Accuracy | Citation Networks |
| Citeseer | 1 | 3,327 | 4,732 | 6 | Accuracy | Citation Networks |
| Pubmed | 1 | 19,717 | 44,338 | 3 | Accuracy | Citation Networks |
| ogbn-arxiv | 1 | 169,343 | 1,166,243 | 40 | Accuracy | Citation Networks |

# Experiment - Main Results

Table 1. Performance (i.e., AUC) of state-of-the-art SSL-based GNN frameworks in the transfer learning setting, and improvements after MGSE is applied. "-" means baseline results are not reported in the original papers. The percentage in the parentheses refers to the percentage of performance improvement brought by MGSE.

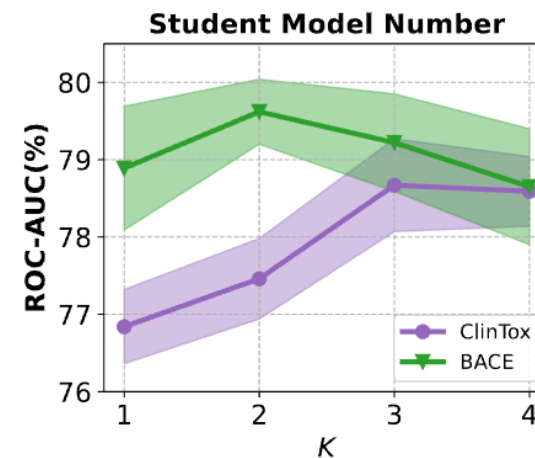| Model | BBBP | Tox 21 | ToxCast | SIDER | ClinTox | HIV | BACE | MUV | PPI |
|---|---|---|---|---|---|---|---|---|---|
| No Pre-train | $65.8_{\pm4.5}$ | $74.0_{\pm0.8}$ | $63.4_{\pm0.6}$ | $57.3_{\pm1.6}$ | $58.0_{\pm4.4}$ | $75.3_{\pm1.9}$ | $70.1_{\pm5.4}$ | $71.8_{\pm2.5}$ | $64.8_{\pm1.0}$ |
| GraphCL | $69.68_{\pm0.67}$ | $73.87_{\pm0.66}$ | $62.40_{\pm0.57}$ | $60.53_{\pm0.88}$ | $75.99_{\pm2.65}$ | $78.47_{\pm1.22}$ | $75.38_{\pm1.44}$ | $69.80_{\pm2.66}$ | $67.88_{\pm0.85}$ |
| +MGSE | $72.26_{\pm0.65}$ | $75.89_{\pm0.33}$ | $64.57_{\pm0.34}$ | $61.44_{\pm0.68}$ | $78.67_{\pm2.89}$ | $79.07_{\pm0.72}$ | $79.22_{\pm0.93}$ | $71.46_{\pm1.45}$ | $69.11_{\pm0.70}$ |
| Perf. (↑) | +2.58 (3.7%) | +2.02 (2.7%) | +2.17 (3.5%) | +0.91 (1.5%) | +2.68 (3.5%) | +0.60 (0.8%) | +3.84 (5.1%) | +1.66 (2.4%) | +1.23 (1.8%) |
| RGCL | $71.42_{\pm0.66}$ | $75.20_{\pm0.34}$ | $63.33_{\pm0.17}$ | $61.38_{\pm0.61}$ | $83.38_{\pm0.90}$ | $77.90_{\pm0.80}$ | $76.03_{\pm0.77}$ | $76.66_{\pm0.99}$ | - |
| +MGSE | $71.65_{\pm0.78}$ | $76.82_{\pm0.62}$ | $64.85_{\pm0.20}$ | $63.72_{\pm0.63}$ | $84.88_{\pm2.01}$ | $78.33_{\pm0.85}$ | $77.40_{\pm1.27}$ | $77.18_{\pm0.81}$ | - |
| Perf. (↑) | +0.23 (0.3%) | +1.62 (2.2%) | +1.52 (2.4%) | +2.34 (3.8%) | +1.50 (1.8%) | +0.43 (0.6%) | +1.37 (1.8%) | +0.52 (0.7%) | - |
| AD-GCL | $70.00_{\pm1.07}$ | $76.54_{\pm0.82}$ | $63.07_{\pm0.72}$ | $63.28_{\pm0.79}$ | $79.78_{\pm3.52}$ | $78.28_{\pm0.97}$ | $78.51_{\pm0.80}$ | $72.30_{\pm1.61}$ | $68.83_{\pm1.26}$ |
| +MGSE | $70.44_{\pm0.70}$ | $76.80_{\pm0.80}$ | $64.60_{\pm0.59}$ | $63.50_{\pm0.92}$ | $83.05_{\pm2.64}$ | $78.91_{\pm0.57}$ | $79.65_{\pm1.07}$ | $74.32_{\pm0.85}$ | $68.95_{\pm0.83}$ |
| Perf. (↑) | +0.44 (0.6%) | +0.26 (0.3%) | +1.53 (2.4%) | +0.22 (0.3%) | +3.27 (4.1%) | +0.63 (0.8%) | +1.14 (1.5%) | +2.02 (2.8%) | +0.12 (0.2%) |
| JOAO | $70.22_{\pm0.98}$ | $74.98_{\pm0.29}$ | $62.94_{\pm0.48}$ | $59.97_{\pm0.79}$ | $81.32_{\pm2.49}$ | $76.73_{\pm1.23}$ | $77.34_{\pm0.48}$ | $71.66_{\pm1.43}$ | $64.43_{\pm1.38}$ |
| +MGSE | $71.93_{\pm0.50}$ | $76.20_{\pm0.33}$ | $64.26_{\pm0.27}$ | $61.02_{\pm0.86}$ | $83.30_{\pm2.44}$ | $77.50_{\pm0.67}$ | $79.82_{\pm0.71}$ | $73.52_{\pm0.62}$ | $65.37_{\pm0.96}$ |
| Perf. (↑) | +1.71 (2.4%) | 1.22 (1.6%) | +1.32 (2.1%) | +1.05 (1.8%) | +1.98 (2.4%) | +0.77 (1.0%) | +2.48 (3.2%) | +1.86 (2.6%) | +0.96 (1.5%) |
| GraphMAE | $72.0_{\pm0.6}$ | $75.5_{\pm0.6}$ | $64.1_{\pm0.3}$ | $60.3_{\pm1.1}$ | $82.3_{\pm1.2}$ | $77.20_{\pm1.0}$ | $83.1_{\pm0.9}$ | $76.3_{\pm2.4}$ | - |
| +MGSE | $71.62_{\pm0.51}$ | $76.52_{\pm0.48}$ | $65.31_{\pm0.38}$ | $62.46_{\pm0.52}$ | $84.41_{\pm2.20}$ | $78.03_{\pm0.70}$ | $82.92_{\pm0.75}$ | $77.15_{\pm0.75}$ | - |
| Perf. (↑) | -0.38 (-0.5%) | +1.02 (1.4%) | +1.21 (1.9%) | +2.16 (3.6%) | +2.11 (2.6%) | +0.83 (1.1%) | -0.18 (-0.2%) | +0.85 (1.1%) | - |
| GraphLoG | $72.5_{\pm0.8}$ | $75.7_{\pm0.5}$ | $63.5_{\pm0.7}$ | $61.2_{\pm1.1}$ | $76.7_{\pm3.3}$ | $77.8_{\pm0.8}$ | $83.5_{\pm1.2}$ | $76.0_{\pm1.1}$ | $66.95_{\pm1.32}$ |
| +MGSE | $72.57_{\pm1.13}$ | $76.84_{\pm0.58}$ | $64.88_{\pm0.39}$ | $63.08_{\pm0.86}$ | $83.72_{\pm2.02}$ | $78.64_{\pm0.80}$ | $83.18_{\pm1.24}$ | $77.22_{\pm0.94}$ | $68.26_{\pm1.06}$ |
| Perf. (↑) | +0.07 (0.1%) | +1.14 (1.5%) | +1.38 (2.2%) | +1.88 (3.1%) | +7.02 (9.2%) | +0.84 (1.1%) | -0.32 (-0.4%) | +1.22 (1.6%) | +1.31 (2.0%) |
| Avg. Perf. (↑) | +0.78 (1.1%) | +1.21 (1.6%) | +1.52 (2.4%) | +1.43 (2.3%) | +3.09 (3.9%) | +0.68 (0.9%) | +1.39 (1.8%) | +1.36 (1.9%) | +0.91 (1.1%) |

Our proposed method can generally further improve the performance of existing graph SSL methods.

## Experiment – Multi-Granularity Design
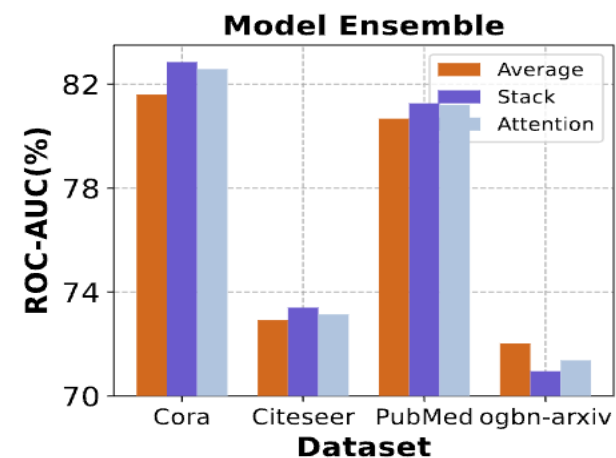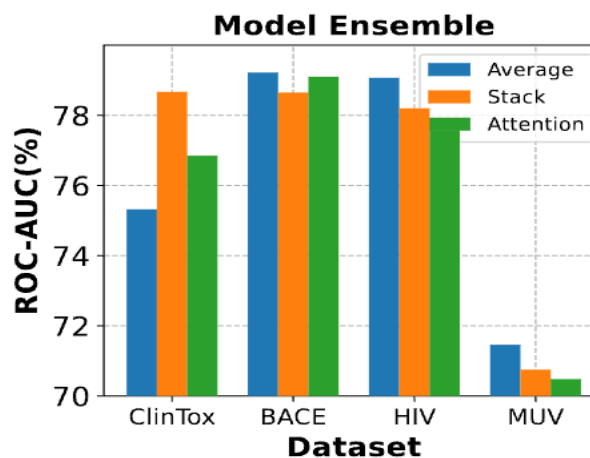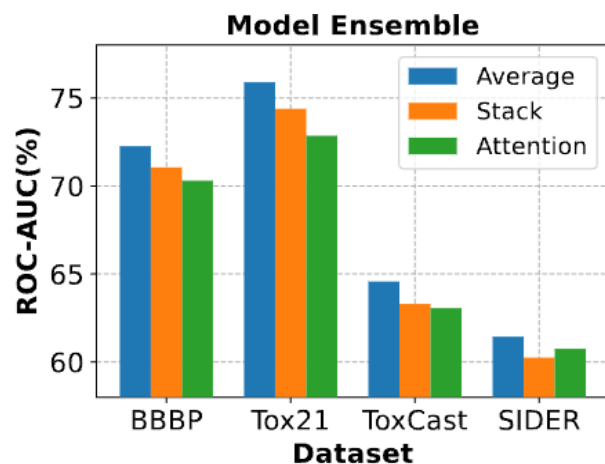
Table 2. Impacts of the multi-granularity design.

| Dataset | BBBP | Tox21 | ToxCast | SIDER |
|---|---|---|---|---|
| Teacher | $69.68_{\pm0.67}$ | $73.87_{\pm0.66}$ | $62.40_{\pm0.57}$ | $60.53_{\pm0.88}$ |
| Multi-teacher | $70.45_{\pm1.04}$ | $73.85_{\pm0.74}$ | $63.10_{\pm0.37}$ | $60.68_{\pm0.70}$ |
| Granularity #1 ($D_K = 2$) | $70.80_{\pm0.70}$ | $74.64_{\pm0.70}$ | $61.75_{\pm0.71}$ | $59.26_{\pm0.52}$ |
| Granularity #2 ($D_K = 21$) | $71.68_{\pm0.77}$ | $75.48_{\pm0.50}$ | $63.10_{\pm0.33}$ | $60.88_{\pm0.65}$ |
| Granularity #3 ($D_K = 50$) | $71.20_{\pm0.92}$ | $75.33_{\pm0.70}$ | $63.85_{\pm0.35}$ | $60.83_{\pm0.65}$ |
| MGSE (Multi-granularity) | $\mathbf{72.26}_{\pm0.65}$ | $\mathbf{75.89}_{\pm0.33}$ | $\mathbf{64.57}_{\pm0.34}$ | $\mathbf{61.44}_{\pm0.68}$ |



Student Model Number

- Multi-granularity has more significant advantage on multi-label prediction tasks which require more complex knowledge.

- Our design provides more flexibility to combine different granularities for various tasks.

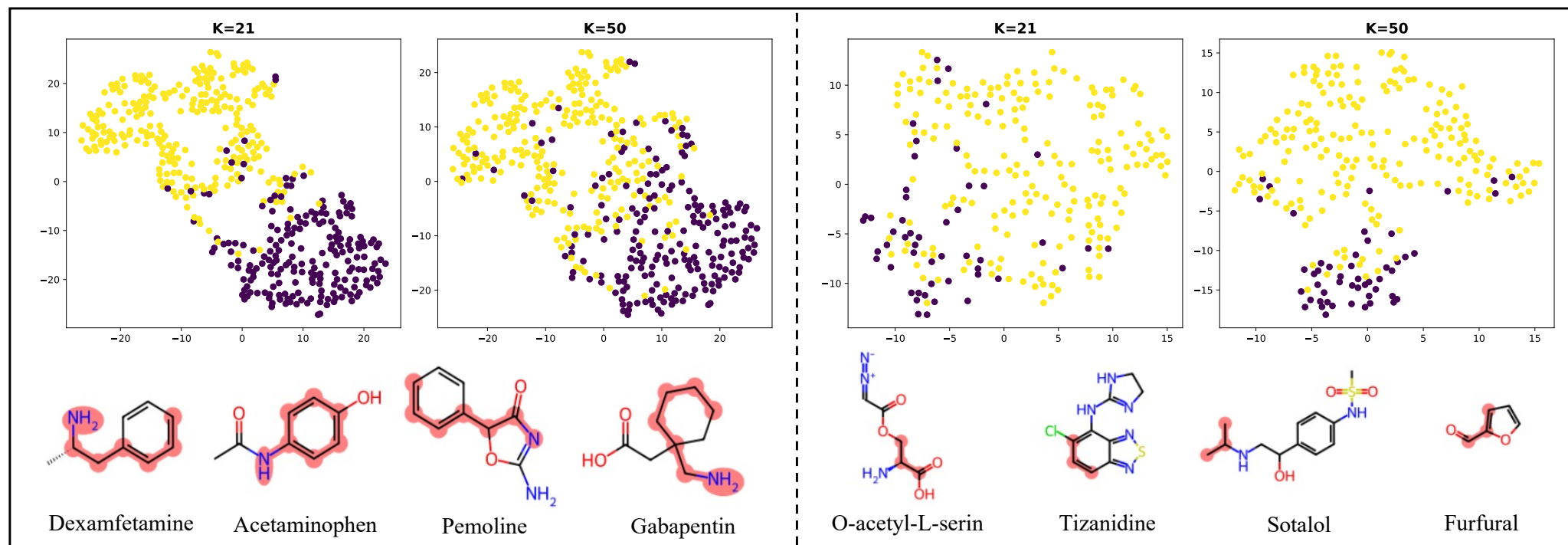## Experiment – Model Ensemble Strategy



The distribution shift between training and testing set may have an influence on the best ensemble strategy.
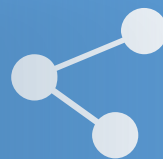
## Experiment – Prototype Visualization

BBBP            Tox21



The key substructure size of positive samples for discrimination of the BBBP dataset is larger than that from the Tox21 dataset. This implies that the classification of BBBP relies on high-level abstract features, whereas fine-grained substructure information is more helpful to the classification of the Tox21 dataset.

# Thanks!
# Q&A

Paper

Code