



---

# Test-time Adaptation on Graphs via Adaptive Subgraph-based Selection and Regularized Prototypes

---

Yusheng Zhao<sup>1</sup> Qixin Zhang<sup>2</sup> Xiao Luo<sup>#3</sup> Junyu Luo<sup>1</sup> Wei Ju<sup>1</sup> Zhiping Xiao<sup>#4</sup> Ming Zhang<sup>#1</sup>

ICML 2025

Test time adaptation (TTA) in Euclidean data:

- **Self-training** often involve selecting reliable data for efficient self-training and aim to construct self-supervising signals via contrastive learning or pseudo-labeling.
- **Data-centric methods** aim to construct virtual samples related to the test distribution as data augmentation.

Domain adaptation in non-Euclidean data:

Some prior works investigate the problem of attempt to transfer the knowledge from the labeled training graphs to unlabeled test graphs via domain discrepancy minimization.

However, these methods typically require labeled source graphs to perform domain alignment, which is often impractical when **facing data privacy issues** (Tan et al., 2023).

# Motivation

- (1) How to overcome the label scarcity of test graphs more reliably in the face of structural shifts?
- (2) How to utilize the knowledge from unknown training graphs and the information from unlabeled test graphs more effectively? (avoid easily overfit and lose the knowledge in training graph)

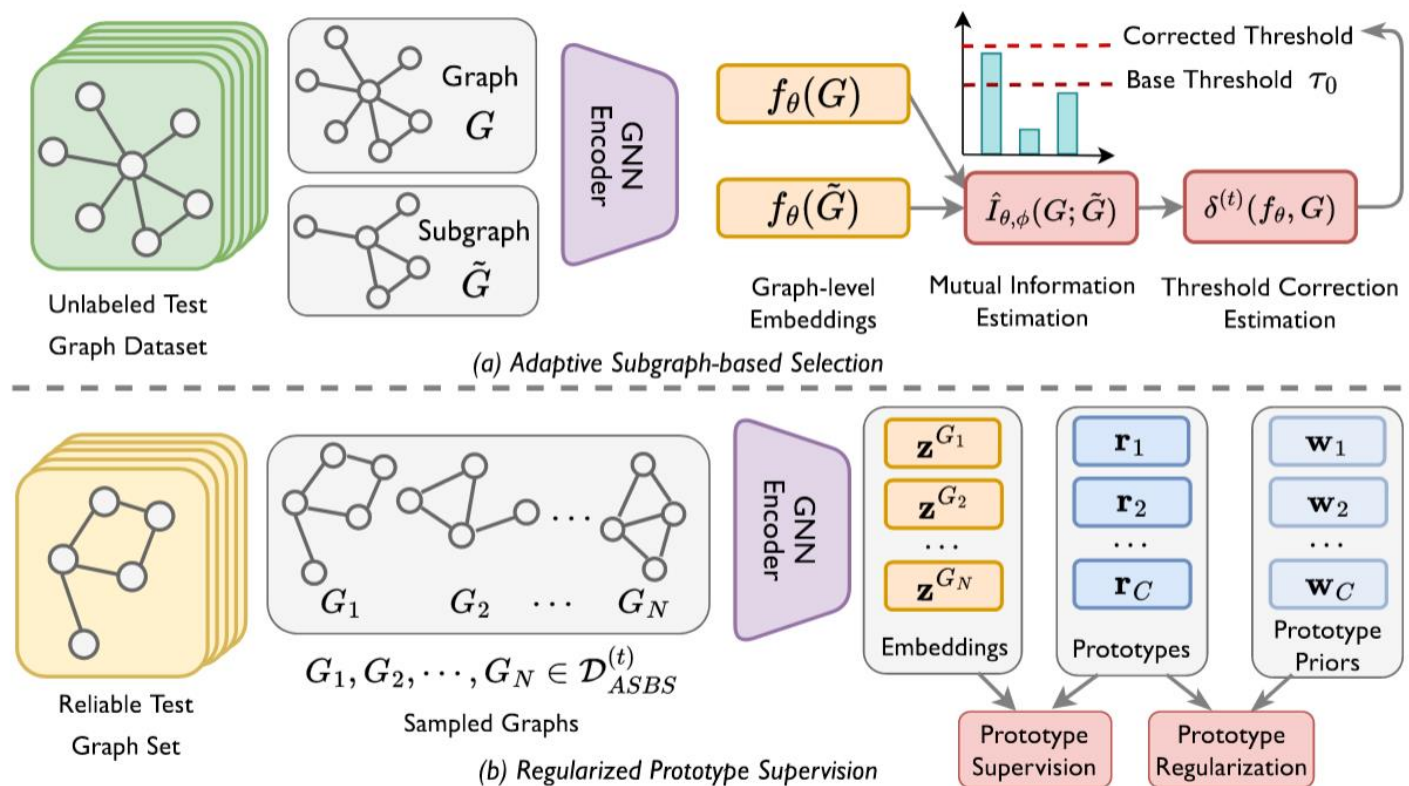


Figure 1. The overall framework of the proposed method. We first select reliable test graphs in the unlabeled test graph dataset using adaptive subgraph-based selection (ASBS), where we utilize mutual information to generate structure-aware, individual-level thresholds. Subsequently, we utilize these graphs for self-training with regularized prototype supervision (RPS), where the prototypes are regularized by prior knowledge and used for supervising the learned embedding of test graphs.

$$G = (V, A, X)$$

$V$  is the set of nodes

$$A \in \mathbb{R}^{|V| \times |V|}$$

$$X \in \mathbb{R}^{|V| \times d^f}$$

$$\mathcal{D}^{tr} = \{(G_i^{tr}, y_i^{tr})\}_{i=1}^{N_{tr}}$$

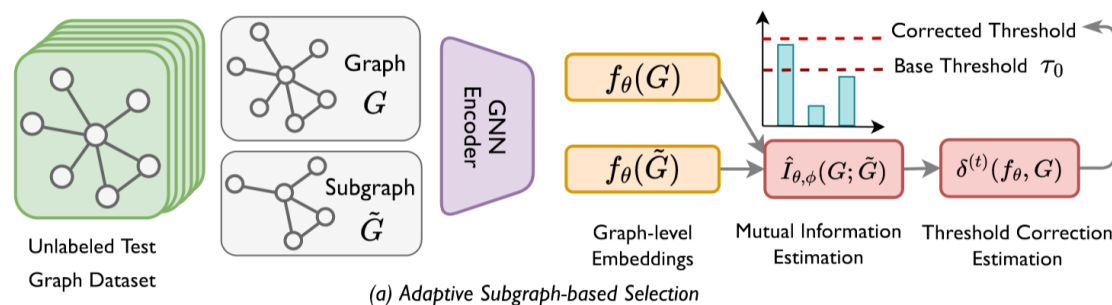
$$\mathcal{D}^{te} = \{G_j^{te}\}_{j=1}^{N_{te}}$$

$$\mathbf{z}^G = f_{\theta}(G)$$

$$\mathbf{p}^G = \text{softmax}(\mathbf{W}^T \mathbf{z}^G)$$

$$\mathbf{p}^G = \mathcal{M}(G)$$

# Method/Adaptive Subgraph-based Selection



$$s_{conf}^G = \max_{i \in \{1, \dots, C\}} p_i^G, \quad \mathcal{D}^{conf} = \{G \in \mathcal{D}^{te} | s_{conf}^G > \tau\},$$

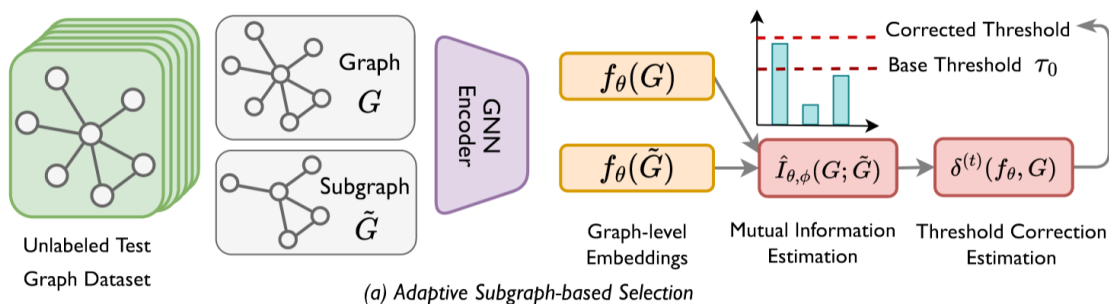
The threshold can be a shared threshold or class-wise threshold.

As the model can be inaccurate in exploring the unlabeled test data, especially when it comes to graphs with complex inherent structures.

Therefore, it is reasonable to use different selection thresholds for graph structures with different levels of complexity.

Graphs with complex structures are hard to learn for the model, and the model tends to be inaccurate. Thus, **a higher threshold should be used, and vice versa.**

# Method/ Adaptive Subgraph-based Selection



Instead of a shared threshold, we adopt a graph-specific, individual-level threshold

$$\tau^G = \tau_0 + \delta(G, f_\theta),$$

$\delta$  measures the model's ability to yield an accurate prediction from  $G$ .

we propose to use mutual information (MI) between graph  $G$  and its subgraph

$$\tilde{G} = (V_{idx}, \mathbf{A}_{idx,idx}, \mathbf{X}_{idx,:})$$

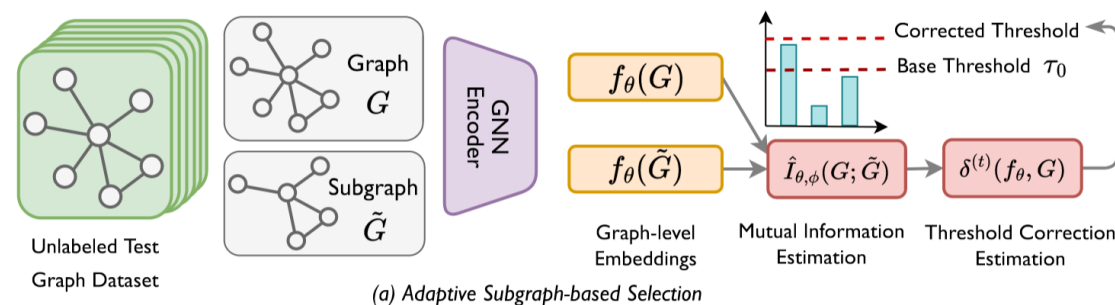
$$I_{\theta,\phi}(G; \tilde{G}) = -\mathbb{E}_{G \sim \mathcal{D}^{te}} [\text{sp}(-g_\phi(f_\theta(G), f_\theta(\tilde{G})))]$$

```

172 ... class DiscriminatorMI(nn.Module):
173     def __init__(self, embed_dim):
174         super(DiscriminatorMI, self).__init__()
175         self.fc_1 = nn.Sequential(
176             nn.Linear(embed_dim, embed_dim),
177             nn.ReLU(inplace=True),
178             nn.Linear(embed_dim, embed_dim)
179         )
180         self.fc_2 = nn.Sequential(
181             nn.Linear(embed_dim, embed_dim),
182             nn.ReLU(inplace=True),
183             nn.Linear(embed_dim, embed_dim)
184         )
185
186     def forward(self, embeddings_1, embeddings_2):
187         emb_1 = self.fc_1(embeddings_1) + embeddings_1
188         emb_2 = self.fc_2(embeddings_2) + embeddings_2
189         score = torch.sum(emb_1 * emb_2, dim=1)
190         return score

```

# Method/Adaptive Subgraph-based Selection



$$\tau^G = \tau_0 + \delta(G, f_\theta),$$

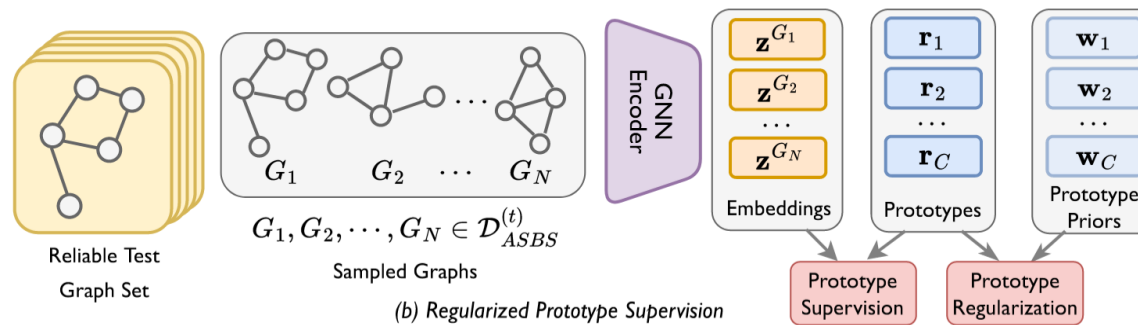
$$\delta(G, f_\theta) = -\omega \hat{I}_{\theta, \phi}(G; \tilde{G}),$$

$$\delta^{(t)}(G, f_\theta) = \omega \sum_{i=0}^{t-1} \beta(1 - \beta)^i \hat{I}_{\theta, \phi}^{(t-i)}(G; \tilde{G}),$$

leads to a stable estimation of mutual information.  
β controlling the speed of updates

$$\mathcal{D}_{ASBS}^{(t)} = \{G \in \mathcal{D}^{te} | s_{conf}^G > \tau^G = \tau_0 + \delta^{(t)}(G, f_\theta)\}.$$

# Method/Regularized Prototype Supervision



initialize them with the weight matrix of the last layer of the model

$$\mathbf{R} = [\mathbf{r}_1, \mathbf{r}_2, \dots, \mathbf{r}_C] \in \mathbb{R}^{d \times C}$$

$$\mathbf{R}^* = \underset{\mathbf{R}}{\operatorname{argmax}} \log p(\mathbf{R}|G), \quad G \in \mathcal{D}_{ASBS}^{(t)}.$$

$$\log p(\mathbf{R}|G) = \log p(\mathbf{R}) + \log p(G|\mathbf{R}) + C_0, \quad \mathbf{R} \sim \mathcal{N}(\mathbf{R}; \mathbf{W}, \sigma_p^2 \mathbf{I})$$

$$\begin{aligned} \log p(\mathbf{R}) &= \log \left( \frac{1}{\sqrt{(2\pi\sigma_p^2)^d}} \exp \left( -\frac{1}{2\sigma_p^2} \|\mathbf{R} - \mathbf{W}\|_2^2 \right) \right) \\ &= C_1 \|\mathbf{R} - \mathbf{W}\|_2^2 + C_2, \end{aligned}$$





Using the Bayes rule, this objective can be further decomposed into prior and likelihood, which is

$$\log p(\mathbf{R}|G) = \log p(\mathbf{R}) + \log p(G|\mathbf{R}) + C_0, \quad (8)$$

where  $C_0$  is a constant. For the first term, we assume that the prototypes  $\mathbf{R}$  follow isotropic Gaussian distribution, i.e.  $\mathbf{R} \sim \mathcal{N}(\mathbf{R}; \mathbf{W}, \sigma_p^2 \mathbf{I})$ , where the mean is set to the initial values of the prototypes, and  $\sigma_p^2$  is the shared variance of the prototypes. Under this assumption, the prior knowledge constraint of the prototypes can be written as:

$$\begin{aligned} \log p(\mathbf{R}) &= \log \left( \frac{1}{\sqrt{(2\pi\sigma_p^2)^d}} \exp \left( -\frac{1}{2\sigma_p^2} \|\mathbf{R} - \mathbf{W}\|_2^2 \right) \right) \\ &= C_1 \|\mathbf{R} - \mathbf{W}\|_2^2 + C_2, \end{aligned} \quad (9)$$

where  $C_1$  and  $C_2$  are constants.

For the second term, we transform it into a self-supervised objective. Inspired by the mixture models (McLachlan & Basford, 1988; Reynolds et al., 2009), we have:

$$\begin{aligned} \log p(G|\mathbf{R}) &= \log p(G|\{\mathbf{r}_i\}_{i=1}^C) = \log \sum_{i=1}^C q(G, i) p(G|\mathbf{r}_i) \\ &\geq \sum_{i=1}^C q(G, i) \log p(G|\mathbf{r}_i) \end{aligned} \quad (10)$$

where  $q(G, i)$  denotes the assignment of graph  $G$  belonging to prototype  $\mathbf{r}_i$ , and the inequality comes from **Jensen's inequality**. Thus, we have a lower bound for the likelihood. For the estimation of  $q(G, i)$ , a common approach is to compute the similarity between the representation of the graph  $\mathbf{z}^G$  and the prototypes  $\mathbf{r}_i$ . Denoting  $\mathbf{q}^G = [q(G, 1), q(G, 2), \dots, q(G, C)]^T$ , this simple objective can be written as

$$\underset{\mathbf{q}^G}{\operatorname{argmin}} \langle \mathbf{q}^G, \mathbf{R}^T \mathbf{z}^G \rangle, \text{ s.t. } \langle \mathbf{q}^G, \mathbf{1}_C \rangle = 1, \quad (11)$$

where  $\langle \cdot, \cdot \rangle$  denotes inner product measuring the similarity, and  $\mathbf{1}_C \in \mathbb{R}^C$  is an all-one vector. However, this simple objective can result in unbalanced estimation, where the number of assignments to some classes is substantially larger than others. To solve this problem, another constraint is needed. Concretely, given a set of graphs  $G_1, G_2, \dots, G_N$ , where  $N$  is the number of graphs, we hope to achieve a balanced assignment, which can be written as

$$\underset{\mathbf{Q}}{\operatorname{argmin}} \operatorname{tr}(\mathbf{Q}^T \mathbf{R}^T \mathbf{Z}), \text{ s.t. } \mathbf{Q} \mathbf{1}_N = \frac{N}{C} \mathbf{1}_C, \mathbf{Q}^T \mathbf{1}_C = \mathbf{1}_N, \quad (12)$$

where  $\mathbf{Q} = [\mathbf{q}^{G_1}, \mathbf{q}^{G_2}, \dots, \mathbf{q}^{G_N}] \in \mathbb{R}^{C \times N}$  is the stacked assignments, and  $\mathbf{Z} = [\mathbf{z}^{G_1}, \mathbf{z}^{G_2}, \dots, \mathbf{z}^{G_N}] \in \mathbb{R}^{d \times N}$  is the stacked graph representations. Solving this problem,

however, requires solving a large linear program (Genevay et al., 2019), and a common workaround is to add entropy regularization that encourages the diversity of assignment. A Lagrangian objective can be written as follows

$$\begin{aligned} \mathcal{E} &= \operatorname{tr}(\mathbf{Q}^T \mathbf{R}^T \mathbf{Z}) + \epsilon \mathcal{H}(\mathbf{Q}) + \langle \mathbf{a}, \mathbf{Q} \mathbf{1}_N - \frac{N}{C} \mathbf{1}_C \rangle \\ &\quad + \langle \mathbf{b}, \mathbf{Q}^T \mathbf{1}_C - \mathbf{1}_N \rangle, \end{aligned} \quad (13)$$

where  $\mathcal{E}$  is the objective to be minimized,  $\epsilon$  is the temperature parameters controlling the degree of regularization,  $\mathbf{a}$  and  $\mathbf{b}$  are Lagrangian multipliers,  $\mathcal{H}(\mathbf{Q}) = -\sum_{i,j} \mathbf{Q}_{i,j} (\log \mathbf{Q}_{i,j} - 1)$  is the entropy regularization. We take the derivative with respect to  $\mathbf{Q}_{i,j}$ , which gives the following result:

$$\frac{\partial \mathcal{E}}{\partial \mathbf{Q}_{i,j}} = [\mathbf{R}^T \mathbf{Z}]_{i,j} - \epsilon \log \mathbf{Q}_{i,j} + \mathbf{a}_i + \mathbf{b}_j. \quad (14)$$

Therefore, the closed-form solution can be written as:

$$\mathbf{Q}^* = \operatorname{diag} \left( \exp \left( \frac{\mathbf{a}}{\epsilon} \right) \right) \exp \left( \frac{\mathbf{R}^T \mathbf{Z}}{\epsilon} \right) \operatorname{diag} \left( \exp \left( \frac{\mathbf{b}}{\epsilon} \right) \right). \quad (15)$$

In practice, we utilize the iterative Sinkhorn-Knopp algorithm (Asano et al., 2019; Caron et al., 2020; Zheng et al., 2021) to solve this problem (more details in Appendix B). Specifically, repeated row normalization and column normalization are conducted on  $\exp \left( \frac{\mathbf{R}^T \mathbf{Z}}{\epsilon} \right)$ , and the algorithm converges quickly. Previous studies (Genevay et al., 2019; Luo et al., 2023) indicate satisfactory performance within 3 iterations, and the soft assignment of the entropy is beneficial.

Then, we provide an estimation of  $p(G|\mathbf{r}_i)$ . Although this probability takes the form of generative modeling, actually training a generative model can be costly, and when facing insufficient data, it can be inaccurate. Therefore, we estimate this term in the latent space using the feature extractor  $f_\theta$  with isotropic Gaussian distribution, which is  $p(G|\mathbf{r}_i) \propto \exp \left( -\frac{1}{2\sigma_z^2} \|f_\theta(G) - \mathbf{r}_i\|_2^2 \right)$ , where  $\sigma_z^2$  is the shared variance of the latent features. Combining the two terms, the loss function can be written as

$$\mathcal{L}_{RPS} = \sum_{i=1}^C q(G, i) \|f_\theta(G) - \mathbf{r}_i\|_2^2 + \alpha_1 \|\mathbf{R} - \mathbf{W}\|_2^2, \quad (16)$$

where  $q(G, i)$  is obtained via the Sinkhorn-Knopp iterations within a batch of graphs, and  $\alpha_1$  is the hyperparameter balancing the two terms. The first term serves as a self-training objective utilizing the posterior information from the unlabeled data, whereas the second term serves as regularization from the prior knowledge of the unknown training graphs. As the prototypes  $\mathbf{r}_i$  (with their matrix form  $\mathbf{R}$ ) are used to provide supervision signals in the self-training and regularized by the prior knowledge, we call this method regularized prototype supervision.

```
if args.update_sinkhorn:
```

```
    sim_mat = F.normalize(x_emb, dim=1) @ F.normalize(prototypes, dim=1).T # B x C
    sh_prob = sinkhorn(sim_mat.detach(), eps=args.sinkhorn_eps, n_iters=5) # B x C
```

$$\mathcal{L} = \mathcal{L}_{RPS} + \alpha_2 \mathcal{L}_{MI},$$

$$\mathcal{L}_{MI} = -I_{\theta, \phi}(\bar{G}; \tilde{G})$$

---

**Algorithm 1** Optimization Algorithm of ASSESS

---

**Input:** The well-trained GNN-based model on the training graphs  $\mathcal{M}^{(0)}(\cdot)$ , test graphs  $\mathcal{D}^{te}$ ,

**Output:** GNN-based model after  $T$  epochs  $\mathcal{M}^{(T)}(\cdot)$

- 1: Initialize the prototypes  $R \leftarrow W$ ;
  - 2: Initialize the reliable test graph set  $\mathcal{D}_{ASBS}^{(0)} = \mathcal{D}^{conf}$  using Eq. 1;
  - 3: **for**  $i = 1, 2, \dots, T$  **do**
  - 4:   **for** each batch **do**
  - 5:     Sample a mini-batch of target graphs;
  - 6:     Calculate  $\mathcal{L}_{RPS}$  using Eq. 16 over  $\mathcal{D}_{ASBS}^{(i)}$ ;
  - 7:     Calculate  $\mathcal{L}_{MI}$  using Eq. 3 over  $\mathcal{D}^{te}$ ;
  - 8:     Calculate the loss objective using Eq. 17;
  - 9:     Update parameters of  $\mathcal{M}^{(i)}$  through back-propagation;
  - 10:   **end for**
  - 11:   Update the threshold correction function using Eq. 5;
  - 12:   Obtain reliable graphs  $\mathcal{D}_{ASBS}^{(i+1)}$  using Eq. 6;
  - 13: **end for**
-





# Experiment

Table 2. The classification accuracy (in %, training  $\rightarrow$  test) on PROTEINS (P0, P1, P2) and NCII (N0, N1, N2).

Methods	PROTEINS							NCII						
	P0 $\rightarrow$ P1	P0 $\rightarrow$ P2	P1 $\rightarrow$ P0	P1 $\rightarrow$ P2	P2 $\rightarrow$ P0	P2 $\rightarrow$ P1	AVG	N0 $\rightarrow$ N1	N0 $\rightarrow$ N2	N1 $\rightarrow$ N0	N1 $\rightarrow$ N2	N2 $\rightarrow$ N0	N2 $\rightarrow$ N1	AVG
GCN	57.0 $\pm$ 4.8	41.9 $\pm$ 0.7	64.8 $\pm$ 2.0	54.7 $\pm$ 3.3	62.5 $\pm$ 1.3	61.7 $\pm$ 3.5	57.1 $\pm$ 2.6	70.1 $\pm$ 3.1	54.8 $\pm$ 1.3	52.6 $\pm$ 2.0	56.3 $\pm$ 1.9	50.1 $\pm$ 3.3	63.0 $\pm$ 1.5	57.8 $\pm$ 2.2
GraphSAGE	53.9 $\pm$ 5.3	41.5 $\pm$ 0.5	70.3 $\pm$ 3.2	57.3 $\pm$ 1.9	69.0 $\pm$ 0.9	62.1 $\pm$ 2.1	59.0 $\pm$ 2.3	69.2 $\pm$ 1.0	55.8 $\pm$ 0.8	51.2 $\pm$ 1.2	59.5 $\pm$ 2.7	48.8 $\pm$ 0.8	61.6 $\pm$ 3.7	57.7 $\pm$ 1.7
GIN	58.9 $\pm$ 1.1	41.7 $\pm$ 0.6	64.2 $\pm$ 1.9	53.1 $\pm$ 2.2	59.2 $\pm$ 0.7	68.2 $\pm$ 2.4	57.6 $\pm$ 1.5	73.0 $\pm$ 2.5	56.9 $\pm$ 1.5	55.6 $\pm$ 4.7	59.1 $\pm$ 2.2	50.3 $\pm$ 0.9	64.5 $\pm$ 1.9	59.9 $\pm$ 2.3
GAT	63.2 $\pm$ 4.3	41.5 $\pm$ 0.5	77.1 $\pm$ 2.9	63.8 $\pm$ 4.1	64.4 $\pm$ 2.9	61.0 $\pm$ 1.8	61.8 $\pm$ 2.8	70.5 $\pm$ 6.9	55.6 $\pm$ 0.4	51.9 $\pm$ 1.4	57.4 $\pm$ 1.9	50.7 $\pm$ 1.2	63.9 $\pm$ 2.9	58.3 $\pm$ 2.5
MeanTeacher	72.0 $\pm$ 7.2	65.3 $\pm$ 5.2	70.5 $\pm$ 2.8	66.5 $\pm$ 7.3	73.3 $\pm$ 5.1	67.2 $\pm$ 2.6	69.1 $\pm$ 5.0	71.5 $\pm$ 8.6	61.9 $\pm$ 3.7	64.7 $\pm$ 10.5	75.4 $\pm$ 2.5	59.4 $\pm$ 12.3	66.1 $\pm$ 1.5	66.5 $\pm$ 6.5
GraphCL	79.2 $\pm$ 6.0	73.7 $\pm$ 6.4	75.2 $\pm$ 5.2	73.9 $\pm$ 4.1	73.5 $\pm$ 5.1	77.0 $\pm$ 7.4	75.4 $\pm$ 5.7	76.2 $\pm$ 3.3	67.6 $\pm$ 1.3	69.0 $\pm$ 1.1	77.3 $\pm$ 1.5	66.2 $\pm$ 3.2	73.1 $\pm$ 3.2	71.6 $\pm$ 2.3
SHOT	68.2 $\pm$ 5.4	62.9 $\pm$ 4.2	71.6 $\pm$ 4.2	65.9 $\pm$ 7.0	67.8 $\pm$ 6.0	68.8 $\pm$ 7.0	67.5 $\pm$ 5.6	64.4 $\pm$ 4.4	56.6 $\pm$ 2.6	73.9 $\pm$ 1.1	65.5 $\pm$ 3.4	69.6 $\pm$ 2.7	64.0 $\pm$ 2.8	65.7 $\pm$ 2.8
TAST	68.9 $\pm$ 7.5	41.9 $\pm$ 1.0	75.1 $\pm$ 2.2	62.1 $\pm$ 2.6	63.8 $\pm$ 2.5	61.9 $\pm$ 1.6	62.3 $\pm$ 2.9	72.2 $\pm$ 5.1	56.5 $\pm$ 0.8	55.2 $\pm$ 2.0	57.7 $\pm$ 2.3	48.9 $\pm$ 0.7	66.0 $\pm$ 3.7	59.4 $\pm$ 2.4
RNA	61.8 $\pm$ 9.1	65.8 $\pm$ 4.1	79.1 $\pm$ 2.3	71.1 $\pm$ 9.1	74.2 $\pm$ 3.1	69.3 $\pm$ 1.9	70.2 $\pm$ 4.9	72.0 $\pm$ 1.8	61.3 $\pm$ 4.7	76.4 $\pm$ 3.7	72.4 $\pm$ 1.1	65.7 $\pm$ 4.4	74.9 $\pm$ 0.5	70.5 $\pm$ 2.7
Ours	81.9 $\pm$ 3.7	74.7 $\pm$ 4.1	86.7 $\pm$ 1.7	77.3 $\pm$ 2.8	82.4 $\pm$ 4.0	70.4 $\pm$ 6.6	78.9 $\pm$ 3.8	79.6 $\pm$ 2.1	69.1 $\pm$ 2.5	81.9 $\pm$ 2.0	78.5 $\pm$ 1.1	73.4 $\pm$ 2.0	70.3 $\pm$ 1.1	75.5 $\pm$ 1.8

Table 3. The classification results (in %, training  $\rightarrow$  test) on IMDB-BINARY (I0, I1, I2).

Methods	I0 $\rightarrow$ I1	I0 $\rightarrow$ I2	I1 $\rightarrow$ I0	I1 $\rightarrow$ I2	I2 $\rightarrow$ I0	I2 $\rightarrow$ I1	AVG
GCN	78.5 $\pm$ 2.0	66.0 $\pm$ 2.9	60.6 $\pm$ 1.8	62.4 $\pm$ 1.7	55.2 $\pm$ 1.3	57.0 $\pm$ 3.3	63.3 $\pm$ 2.2
GraphSAGE	78.2 $\pm$ 1.8	60.6 $\pm$ 4.1	58.5 $\pm$ 3.2	63.9 $\pm$ 1.5	56.7 $\pm$ 0.9	56.7 $\pm$ 2.7	62.4 $\pm$ 2.4
GIN	75.8 $\pm$ 1.7	66.6 $\pm$ 2.2	58.8 $\pm$ 5.6	64.5 $\pm$ 1.7	54.3 $\pm$ 0.7	56.7 $\pm$ 2.1	62.8 $\pm$ 2.4
GAT	77.3 $\pm$ 1.1	64.5 $\pm$ 1.7	63.6 $\pm$ 0.7	62.1 $\pm$ 2.0	54.9 $\pm$ 2.2	58.8 $\pm$ 0.7	63.5 $\pm$ 1.4
MeanTeacher	70.1 $\pm$ 6.8	65.2 $\pm$ 2.5	63.2 $\pm$ 4.3	70.1 $\pm$ 2.4	59.7 $\pm$ 2.4	61.2 $\pm$ 4.2	64.9 $\pm$ 3.8
GraphCL	66.2 $\pm$ 3.1	70.6 $\pm$ 3.1	61.7 $\pm$ 1.9	70.1 $\pm$ 2.4	61.7 $\pm$ 4.9	46.3 $\pm$ 4.4	62.8 $\pm$ 3.3
SHOT	68.9 $\pm$ 5.1	66.6 $\pm$ 4.5	65.7 $\pm$ 3.4	67.8 $\pm$ 6.8	61.5 $\pm$ 7.3	60.6 $\pm$ 9.4	65.2 $\pm$ 6.1
TAST	80.3 $\pm$ 2.2	65.1 $\pm$ 3.8	58.2 $\pm$ 2.1	64.2 $\pm$ 1.9	57.6 $\pm$ 0.7	58.2 $\pm$ 2.5	63.9 $\pm$ 2.2
RNA	64.2 $\pm$ 3.2	67.2 $\pm$ 4.9	59.7 $\pm$ 3.2	69.2 $\pm$ 2.5	61.7 $\pm$ 3.5	68.7 $\pm$ 3.2	65.1 $\pm$ 3.4
Ours	82.1 $\pm$ 3.7	66.7 $\pm$ 3.7	63.7 $\pm$ 1.9	63.7 $\pm$ 1.9	69.7 $\pm$ 3.1	61.2 $\pm$ 3.2	67.8 $\pm$ 2.9

Three types:

1. GNN.
2. Un/semi-supervised.
3. TTA

The latest baseline is RNA (Luo et al., 2024)

The unsupervised / semisupervised learning methods are not the optimal solution. Lack of consideration of test-time distribution shifts.

TTA methods designed for Euclidean data are also not very satisfactory on graphs.

# Experiment

Table 4. Ablation studies on three datasets, *i.e.* FRANKENSTEIN (FRAN), Mutagenicity (MUTA), and NCI1. Average accuracy over six settings (in %) is reported.

Experiments	FRAN	MUTA	NCI1	AVG
w/o ASBS	59.8 $\pm$ 2.0	72.4 $\pm$ 1.3	73.9 $\pm$ 1.4	68.7
w/o RPS-a	59.4 $\pm$ 1.5	71.3 $\pm$ 1.1	74.0 $\pm$ 1.7	68.2
w/o RPS-b	59.4 $\pm$ 1.2	70.6 $\pm$ 1.5	74.1 $\pm$ 1.8	68.0
Full Model	60.3 $\pm$ 1.9	73.5 $\pm$ 1.5	75.5 $\pm$ 1.8	69.8

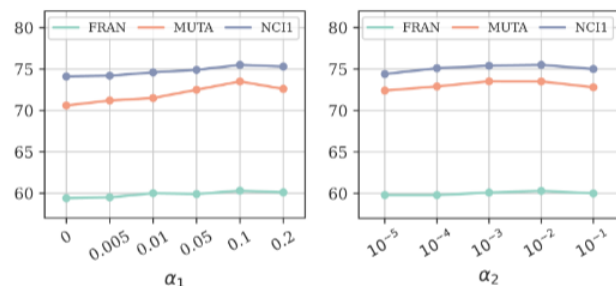


Figure 2. The model's sensitivity to hyperparameters (*i.e.*  $\alpha_1$  and  $\alpha_2$ ). The experiments are performed on FRANKENSTEIN (FRAN), Mutagenicity (MUTA), and NCI1. Average accuracy (in %) are reported over six settings.

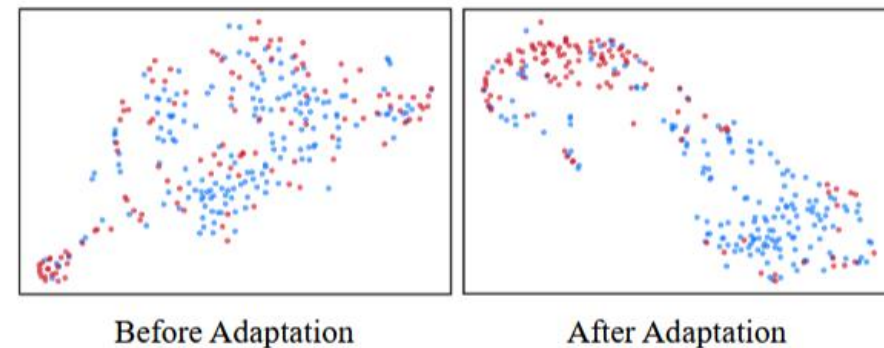


Figure 3. Visualization of learned representations on the PROTEINS dataset, before and after adaptation in the P1  $\rightarrow$  P0 setting.

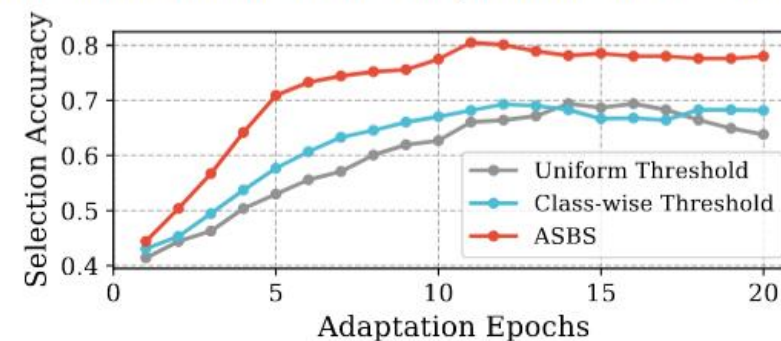


Figure 4. The selection accuracy using uniform threshold, class-wise threshold, and the proposed ASBS on the PROTEINS dataset.



南京航空航天大学  
NANJING UNIVERSITY OF AERONAUTICS AND ASTRONAUTICS

**Thank you**