Disentangling and Re-evaluating The Effectiveness of **Graph Structure Learning For GNNs**

Anonymous Author(s)

Affiliation Address email

Abstract

Graph Structure Learning (GSL) has been widely adopted in the design of Graph Neural Networks (GNNs), with similarity-based graph learning emerging as the most popular approach for node classification. However, which component of GSL really enhances GNN performance remains underexplored. In this paper, we disentangle its effects and present a comprehensive analysis. Specifically, we propose a novel framework that can decompose GSL into three steps: (1) GSL bases (i.e., processed node embeddings for construction) generation, (2) new graph construction, and (3) view fusion. Through empirical studies and theoretical analysis, we demonstrate that applying graph convolution to the newly constructed graphs does not increase the Mutual Information (MI) between node embeddings and labels. Our findings reveal that model performance is primarily driven by the quality of GSL bases rather than the graph construction methods. To validate them, we conduct extensive experiments with 450 GSL variants and benchmark them against GNN baselines within the same search space for GSL bases. Results show that similarity-based graph construction has negligible or even adverse impacts on GNN performance, while pre-trained GSL bases provide significant performance gains. These findings verify and confirm our analysis, underscoring the critical role of GSL bases and highlighting the need to simplify the other two GSL steps.

Introduction 19

2

3

5

6

7

8

10

11

12

13 14

15

16

17

18

34

Graph Neural Networks (GNNs) [17] are effective in capturing structural information from non-20 Euclidean data, which can be used in many applications such as recommendation [50, 49], telecom-21 munication [29], bio-informatics [54, 12, 13], and social networks [34, 24]. However, conventional 22 GNNs suffer from issues including heterophily [30, 31], over-squashing [5], adversarial attacks 23 [15, 22], and missing or noisy structures [21, 28]. To address these issues, Graph Structure Learning 24 (GSL), especially the similarity-based method that reconstructs or refines the original graph structures, 25 has been widely used in enhancing GNN performance and robustness [60]. Even though GSL is 26 believed to improve GNN performance, it introduces more hyperparameters and adds plenty of 27 computational cost in both the construction process and the learning process. In addition, recent 28 studies [39, 36] have shown that GSL methods cannot consistently outperform baseline GNNs with 29 the same hyperparameter tuning strategy. Therefore, an in-depth analysis of the effectiveness and 30 necessity of GSL is highly needed. 31

To have a detailed understanding of each component in GSL, we propose a new framework that can break down GSL into 3 steps: (1) GSL Bases Generation. GSL bases are the processed node embeddings that serve as inputs for the structure construction of new graphs. They are built by either graph-aware or graph-agnostic models with fixed or learnable parameters. (2) Graph Structure 35 **Construction.** Based on the GSL bases, new structures are constructed with similarity-based [14, 38], structural-based [55, 27], or optimization-based approaches [15] ¹, followed by graph refinements.

(3) **View Fusion.** To incorporate the original graph or combine multiple GSL-generated graphs, various view fusion strategies are applied, *e.g.*, late fusion [45], early fusion [22], or separation [28]. Compared with existing categorizations of GSL [40, 60, 61, 18] that mainly focus on **step (2)**, our proposed framework is more comprehensive, and is able to disentangle the effect of each component in GSL.

More specifically, for **step** (1), we argue that a fair comparison between GSL-enhanced GNNs and traditional GNNs should be made using the same GSL bases. Previous GSL studies often enhance node inputs with additional information before graph construction, such as pre-trained node embeddings [6, 51] or structural embeddings [38, 57]. However, these enhancements are typically absent for the inputs of GNN baselines, leading to potentially biased and invalid evaluations. For **step** (2), we examine the effectiveness of graph convolution operations with the similarity-based graphs. Our empirical and theoretical findings indicate that the Mutual Information (MI) between convolved node representations and labels does not increase after graph convolution. This suggests that the performance improvements observed in previous similarity-based GSL methods result from the processed GSL bases (*i.e.*,, enhanced node inputs) in step (1) rather than new graph construction in step (2)².

We conducted extensive experiments to validate our hypothesis. To thoroughly evaluate the performance of GSL-enhanced GNNs, we implemented these methods using six GNN backbones, five GSL bases, three GSL graph construction approaches, three view fusion methods, and two types of fusion strategies, resulting in 450 different GSL variants. The results demonstrate that, within the same search space of GSL bases, there are no significant performance differences between GSL-enhanced GNNs and the corresponding baseline GNNs on node classification tasks. In addition, the results show that the pre-trained GSL bases is the component which significantly enhance GNN performance on certain datasets. This aligns with our analysis and validate our claim. In summary, our main contributions are as follows:

- Comprehensive GSL Framework: In Section 3, we propose a novel framework that decomposes the GSL process into three steps. This decomposition provides a more comprehensive perspective than existing categorizations, offering valuable insights into the workings of GSL.
- Empirical and Theoretical Analysis: In Section 4, we present both empirical evidence and theoretical analysis demonstrating that the Mutual Information (MI) between node representations and labels does not increase after applying graph convolution on similarity-based GSL graphs. This finding suggests that similarity-based GSL methods may be unnecessary.
- Fair Re-Evaluation of GSL: In Section 5, we conduct a fair reassessment of GSL's impact on GNN performance. Our results highlight that GSL bases play a crucial role in improving GNN performance, while similarity-based graph construction has a negligible effect. Besides, we identify the key components for effective GSL, including pretrained GSL bases, parameter separation, and early fusion strategies.

2 Preliminary

Graphs. Suppose we have an undirected graph $\mathcal{G} = \{\mathcal{V}, \mathcal{E}\}$ with node set \mathcal{V} and edge set \mathcal{E} . Let $\mathbf{Y} \in \mathbb{R}^{N \times C}$ denote the node labels and $\mathbf{X} \in \mathbb{R}^{N \times M}$ represent the node features, where N is the number of nodes, C is the number of classes, and M is the number of features. The graph structure is represented by an adjacency matrix \mathbf{A} , where $\mathbf{A}_{u,v} = \mathbf{A}_{v,u} = 1$ indicates the existence of an edge $e_{uv}, e_{vu} \in \mathcal{E}$ between nodes u and v. The normalized adjacency matrix is given by $\hat{\mathbf{A}} = \tilde{\mathbf{D}}^{-\frac{1}{2}}\tilde{\mathbf{A}}\tilde{\mathbf{D}}^{-\frac{1}{2}}$, where $\tilde{\mathbf{D}} = \mathbf{D} + \mathbf{I_n}$ and $\tilde{\mathbf{A}} = \mathbf{A} + \mathbf{I_n}$ represent the degree matrix and adjacency matrix with added self-loops. The neighbors of node u is denoted as $\mathcal{N}_u = \{v | e_{uv} \in \mathcal{E}\}$. Graph Structure Learning (GSL) generates a new graph topology \mathbf{A}' , where the new neighbors of node u are denoted as \mathcal{N}'_u . Graph-aware models $\mathcal{M}^{\mathcal{G}}$, such as Graph Convolutional Networks (GCN) [17], are powerful in extracting structural information in graphs by message aggregation or graph filters [35]. In contrast, graph-agnostic models $\mathcal{M}^{\neg\mathcal{G}}$, such as Multilayer Perceptrons (MLP), only use \mathbf{X} without considering \mathcal{G} . For example, the updating process of node embeddings in GCN and MLP

¹Note that most of the construction methods are similarity-based, which is the main focus of our paper.

²As step (3) fuses the results from step (2), if step (2) is ineffective, then step (3) will also be ineffective.

can be represented as $H^l = \sigma(\hat{A}H^{l-1}W^{l-1})$ and $H^l = \sigma(H^{l-1}W^{l-1})$, respectively. Here, H^l and W^l are the node embeddings and weight matrix at the l-th layer, respectively, and $\sigma(\cdot)$ is an 91 activation function. 92

Graph Homophily. The concept of homophily originates from social network analysis and is 93 defined as the tendency of individuals to connect with others who have similar characteristics [16]. A 94 higher level of graph homophily makes the topological information of each node more informative, thereby improving the performance of graph-aware models $\mathcal{M}^{\mathcal{G}}$ [32, 33, 56]. Commonly used homophily metrics include edge homophily [2, 59] and node homophily [38]:

$$h_{\text{edge}}(\mathcal{G}, \mathbf{Y}) = \frac{\left| \left\{ e_{uv} \mid e_{uv} \in \mathcal{E}, Y_u = Y_v \right\} \right|}{|\mathcal{E}|} \tag{1}$$

$$h_{\text{edge}}(\mathcal{G}, \mathbf{Y}) = \frac{\left| \{ e_{uv} \mid e_{uv} \in \mathcal{E}, Y_u = Y_v \} \right|}{|\mathcal{E}|}$$

$$h_{\text{node}}(\mathcal{G}, \mathbf{Y}) = \frac{1}{|\mathcal{V}|} \sum_{v \in \mathcal{V}} \frac{\left| \{ u \mid u \in \mathcal{N}_v, Y_u = Y_v \} \right|}{|\mathcal{N}_v|}$$

$$(2)$$

Mutual Information. Mutual Information quantifies the amount of information obtained about one 100 random variable given another variable [1]. The mutual information between variable X and Y can 101 be expressed as: 102

$$I(\mathbf{X}; \mathbf{Y}) = \sum_{y \in \mathcal{Y}} \sum_{x \in \mathcal{X}} p(x, y) \log \frac{p(x, y)}{p(x)p(y)}$$
(3)

where p(x, y) is joint probability, and p(x) and p(y) are marginal probability. 103

Mutual information could be used to analyze the quality of input features by measuring how much 104 information the inputs X retain about the outputs Y. However, in graphs under the task of node 105 classification, the mutual information between a discrete variable Y and a continuous variable X106 cannot be directly measured by Eq. (3). Therefore, in this paper, we measure the mutual information 107 108 I(X; Y) based on entropy estimation from k-nearest neighbors distances following [19, 20, 41].

3 **Graph Structure Learning**

Existing studies and evaluations of GSL mainly focus on the structure construction method. However, 110 through extensive literature review, we find that it only constitutes one step of GSL [40, 60, 61]. To 111 comprehensively understand and disentangle GSL for GNN learning, we propose a new framework. 112 As shown in Figure 1, our framework includes three steps: GSL bases generation, new structure 113 construction, and view fusion. Then, the whole pipeline of GSL is: First, GSL bases B is constructed 114 based on node features X (and input graphs \mathcal{G}); Then, new graph structures \mathcal{G}' are constructed with 115 the GSL bases; At last, the information from \mathcal{G}' (sometimes with multiple views) and original graph 116 G are combined with different view fusion strategies for GNN training. We will introduce each 117 component in the following subsections.³ 118

3.1 GSL Bases

98

99

109

119

120

121

122

123

124

125

126

127

128

129

The GSL bases B is defined as the pre-processed node embeddings used for new structure construction. The quality of the GSL bases plays a crucial role for the graph construction step. For node classification tasks, an effective GSL bases B should exhibit consistency among intra-class nodes, as shown in Figure 2 (left). The construction of B can be categorized into non-parametric approaches [8, 38, 63], which generate fixed B, and parametric approaches [15, 6, 53], where Bis learnable during training. The construction of B can also be categorized into graph-agnostic [8, 15, 63] and graph-aware approaches [38, 53, 45], based on whether the original graph information will be contained in B. Combining these two perspectives, in Figure 1, we show the diagrams of four types of bases: B = X, $B = (A)^k X$, B = MLP(X), and B = GNN(X, A).

3.2 New Structure Construction

The construction of the new structure \mathcal{G}' , based on B, is a key element of GSL. Based on relation 130 extraction methods, the construction of \mathcal{G}' can be categorized into similarity-based [14, 38, 23], 131 structure-based [55, 27, 63], and parametric optimization-based [15, 28, 25] approaches. Similarity-132 based method is the most prevalent one, and the choice of similarity measurement, such as k-Nearest 133

³Please refer to Appendix A for a more detailed discussion of the representative GSL methods within our proposed GSL framework.



Figure 1: Our proposed GSL framework consists of three steps: GSL base generation, new structure construction, and view fusion.

Neighbors [8], cosine similarity [6], or Minkowski distance [28], plays a critical role in the quality of the reconstructed graphs. However, the initial \mathcal{G}' produced by these methods often results in a coarse graph structure, which may not be optimal for GNN training. Thus, further refinements are often necessary, such as sampling [55, 22, 28], symmetrization [53, 7, 28], normalization [14, 55, 28], or applying graph regularization [14, 15, 25].

3.3 View Fusion

For GSL methods which have already implicitly fused the information from the original graph structure \mathcal{G} into the reconstructed structure \mathcal{G}' [14, 7, 63], further view fusion is unnecessary. However, for other approaches, the fusion of information from \mathcal{G} and \mathcal{G}' is crucial. Based on the fusion stage, methods can be classified as early fusion [22, 21, 27], late fusion [45, 28, 57], and separation [28]. Early fusion, often seen as "graph editing", modifies \mathcal{G} by adding or removing edges with \mathcal{G}' before training. Late fusion keeps both views as input, fusing node embeddings either at each layer or in the final layer. Separation methods, typically paired with contrastive learning, maintain multiple views without embedding fusion during GNN training. Additionally, view fusion methods can be further distinguished by whether they involve parameter sharing across layers during training.

3.4 Training Mode

In addition to the previous three steps, the training mode of \mathcal{G}' plays a crucial role in GSL and can be categorized into static, joint, and 2-stage approaches. Most methods [15, 25, 51] use joint training where \mathcal{G}' and model parameters are optimized simultaneously. In contrast, some methods [8, 45, 27] follow a 2-stage mode, iteratively updating \mathcal{G}' and model parameters. While dynamic updates offer better flexibility for learning complex structures through parameter optimization, they also significantly increase computational complexity, especially during the bases and graph construction steps. To address this, other methods [43, 23, 57] opt for a static \mathcal{G}' during training. Although this fixed structure may limit performance, it avoids the time-consuming process of frequent graph updates.

4 Effectiveness of Graph Structure Learning

Based on the framework built in the previous section, in this section, we question the necessity of similarity-based graph construction methods with theoretical analysis and extensive experiments. We introduce the motivation with an example in Section 4.1. We then explore the impact of GSL on

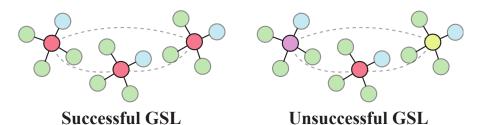


Figure 2: Examples of GSL that use the similarity of neighbor distribution as GSL bases for graph construction. The color of nodes indicates their labels. **Left**: new edges successfully connect intraclass (red) nodes, which share similar GSL bases (neighborhood pattern with 3 green nodes and 1 blue node). **Right**: new edges connect inter-class nodes, resulting in an unsuccessful GSL.

GNN performance through empirical observations in Section 4.2 and theoretical analysis in Section 4.3. Finally, the time complexity of GSL is discussed in Section 4.4.

4.1 Motivation

163

164

165

180

181

183

We revisit the effectiveness of GSL by the examples shown in Figure 2, where we use neighborhood 166 distributions as GSL bases. Suppose a successful new edge is the one that connect intra-class 167 nodes (Figure 2 (Left)), i.e., homophilic connection [31]; and an unsuccessful edge connects inter-168 class nodes (Figure 2 (Right)), i.e., heterophilic edge. We can see that successful and unsuccessful 169 connections both follow node similarity principle to build edges. In other words, the same construction 170 method can lead to totally different outcomes. Instead, the main difference comes from the GSL bases: 171 the left example has high-quality bases, where intra-class nodes share consistent representations; 172 however, the right example has low-quality bases, where inter-class nodes have similar embeddings. 173 On the other hand, this example also points out an awkward situation for GSL: when we have low-quality bases, GSL cannot work well; when we have high-quality bases, it means that the bases 175 themselves can already provide sufficiently informative and distinguishable node embeddings for 176 classification, and therefore, new graph construction may still be unnecessary. To further explore 177 the effectiveness of the new graph construction step in GSL, we conduct empirical and theoretical 178 analyses in the following subsections. 179

4.2 Empirical Observations on Synthetic Graphs

In this section, we investigate how GSL bases and the graph reconstruction methods influence GNN performance through experiments on synthetic graphs. The graph generation process is as follows.

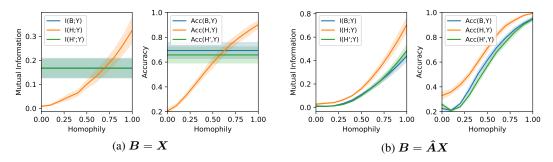


Figure 3: Mutual information and accuracy of node classification for GSL bases B, convoluted bases $H = \hat{A}B$, new graph convoluted bases $H' = \hat{A}'B$, across various homophily degrees. B is set to node features X in (a) and aggregated features $\hat{A}X$ in (b). Note that in (a), \hat{A}' only depends on B and does not change with homophily.

Settings Based on CSBM-H [33] (see details in Appendix B), we generate synthetic graphs with 10 random seeds for each homophily degree $h \in \{0, 0.1, ..., 1.0\}$ to mitigate randomness effects. Each graph \mathcal{G} contains 1000 nodes, with each node characterized by 10 features, 5 balanced classes,

and a degree sampled from the range [2, 10]. Then, we apply k-Nearest-Neighbors (kNN) on GSL bases \boldsymbol{B} with k=5 to generate new graphs, i.e., $\mathcal{G}'=\text{kNN}(\boldsymbol{B})$.

The experiments are designed to answer two questions: **Q1**: Is the reconstructed graph necessary for GNN? **Q2**: How much does the reconstructed graph enhance GNN performance compared to the original graph structure?

Let us denote the original node representations (*i.e.*, original GSL bases) as $\bf B$, the original graph convoluted representations as $\bf H$, and the reconstructed graph convoluted embeddings as $\bf H'$. These bases will be separately fed into MLP, GCN, and GCN+GSL to compare model performance (prediction accuracy). We measure the quality of the bases using both the non-parametric metric mutual information $\bf I(\cdot)$ and the parametric metric $\bf Acc(\cdot)$. We test two different settings of GSL bases (1) graph-unaware GSL bases $\bf B = \bf X$ in Figure 3a, where $\bf H'$ does not rely on graph homophily; (2) and graph-aware GSL bases $\bf B = \hat{\bf A} \bf X$ in Figure 3b, where $\bf H'$ depends on graph homophily.

To answer **Q1**, we can compare I(B; Y) vs. I(H'; Y) and Acc(B; Y) vs. Acc(H'; Y). They compare the models which have and do not have graph reconstruction step, under the same GSL bases. To answer **Q2**, we can compare I(H; Y) vs. I(H'; Y) and Acc(H; Y) vs. Acc(H'; Y), which compare the performance of GCN and GSL-enhanced GCN. Through extensive experiments, we have the following observations.

Observation 1. Mutual information is an effective non-parametric measure of model performance. As shown in Figure 3a and 3b, the shape of mutual information $I(\cdot)$ curves (left) are highly similar to the curves for model accuracy $ACC(\cdot)$ (right). This shows that mutual information can effectively measure the quality of the embeddings. We will use it for theoretical analysis in the next section.

Observation 2. Graph construction does not make significant difference. In Figure 3, the mutual information I(B; Y) and classification accuracy ACC(B, Y) are close to I(H'; Y) and ACC(H', Y), respectively, across both graph-agnostic and graph-aware GSL bases. This suggests that the model performance does not improve significantly after applying graph convolution on the reconstructed graph \mathcal{G}' , which aligns with our analysis in the previous section.

Observation 3. GSL-enhanced GCN only outperforms GCN in heterophilous graphs under graph-agnostic bases. With graph-agnostic bases in Figure 3a, I(H; Y) and ACC(H, Y) increase as homophily increases, while I(H'; Y) and ACC(H', Y) remain constants across homophily degrees. As the reconstructed graph does not depend on homophily, the harmful connections in graphs with low homophily will not cause negative impact on H'. Thus, GSL-enhanced GCN can outperform GCN. However, this effect is observed only when B = X.

Note that when B = AX in Figure 3b, even when GCN+GSL outperforms GCN, its performance still remains close to MLP under the same GSL bases. This means that GSL-enhanced GNN cannot outperform the simple baselines significantly. Recent studies [36, 39] also indicate that under consistent hyperparameter tuning, GSL does not always consistently outperform classic GNN baselines. This leads us to reconsider the necessity of GSL. In addition to the above empirical observations, we proceed with a theoretical analysis on the effectiveness of GSL in the following section.

4.3 Theoretical Analysis

213

214

217

218

226

To explain the above empirical observations, in this section, we first prove that the mutual information I(Y;H) between label Y and aggregated features H can serve as a non-parametric measurement of the effect of graph convolution. Following this, we compare the mutual information between the node labels Y and either the original GSL bases B or the aggregated GSL bases H' (on G'), to reveal the impact of GSL on model performance.

Theorem 4.1. Given a graph $\mathcal{G} = \{\mathcal{V}, \mathcal{E}\}$ with node labels \mathbf{Y} and node features \mathbf{X} , the accuracy of graph convolution on node classification P_A is upper bounded by the mutual information of node label Y and aggregated node features H = AX:

$$P_A \le \frac{I(Y;H) + \log 2}{\log(C)} \tag{4}$$

Proposition 4.2. Consider a graph $\mathcal{G} = \{\mathcal{V}, \mathcal{E}\}$ characterized by node labels Y and n-dimensional node bases $\mathbf{B} = (B_1, B_2, \dots, B_n)$ with C classes. Each base B_i is independent and follows a class-dependent Gaussian distribution, i.e., $B_i \sim \mathcal{N}(\mu_Y, \sigma_Y)$. A new graph $\mathcal{G}' = \{\mathcal{V}, \mathcal{E}'\}$ is generated

using a non-parametric method based on the bases \mathbf{B} . For the aggregated bases $\mathbf{B'}$ on $\mathcal{G'}$, we have inf $I(Y; \mathbf{B'}) \leq \inf I(Y; \mathbf{B})$.

240 where the proofs are shown in Appendix C.

Theorem 4.1 shows that the mutual information I(Y;H) provides an upper bound on the accuracy of graph convolution for node classification, which justifies why mutual information serves as an effective measure of model performance, as demonstrated in Observation 1.

Based on the conclusion of mutual information in Theorem 4.1, we analyze the effectiveness of GSL. Proposition 4.2 shows that the graph convolution on new graphs generated by GSL does not increase the lower bound of mutual information. This explains why MLP performs similarly to, or slightly better than, GCN+GSL in Observation 2 and the dilemma of GSL in Figure 2

To further explain Observation 3 in Section 4.2, we refer again to Proposition 4.2. In conjunction with previous studies on graph homophily [38, 32, 56], we know that the performance of GCN could

with previous studies on graph homophily [38, 32, 56], we know that the performance of GCN could be inferior to MLP on heterophilous graphs. Since GCN+GSL is upper bounded by the MLP on the same GSL bases, when MLP outperforms GCN, GCN+GSL may also outperform GCN, as seen in Figure 3a. However, even when GCN+GSL surpasses GCN in some cases, it still lags behind MLP, a much simpler model, on the same GSL bases. Therefore, we hypothesize that previous GSL improvements stem from the construction of the GSL bases or the introduction of additional model parameters. A fair comparison of GSL with other GNNs or MLP baselines should be conducted using the same GSL bases, as demonstrated in our experiments.

4.4 Complexity Analysis

257

272

286

After investigating the difference in the performance of GCN+GSL and GCN, we then analyze the 258 time complexity of some representative methods of GSL, such as IDGL [6], GRCN [53], GAug [55], 259 and HOG-GCN [46], as shown in Table 2. Assume the dimension of node representation is F for all 260 the layers, the additional time complexity introduced by GSL generally includes: 1. Construction 261 of GSL bases: $O(|\mathcal{E}|F + |\mathcal{V}|F^2)$ for graph-aware bases or $O(|\mathcal{V}|F^2)$ for graph-agnostic bases, 2. 262 Graph construction: $O(|\mathcal{V}|^2 F)$, 3. Graph refinement: $O(|\mathcal{V}|^2)$, and 4: View Fusion $O(|\mathcal{V}|^2)$. Apart from the complexity of the new graph construction in GSL, during the graph convolution, compared 264 with GNNs without using GSL, the additional complexity is further introduced by single view GSL 265 $O(|\mathcal{E}'|F)$ or multiple view GSL $O((N_{\mathcal{G}}-1)(|\mathcal{E}|F+|\mathcal{V}|F^2))$, where $|\mathcal{E}'|$ is the additional edges 266 introduced in GSL and N_G is the number of views in GSL. Consider the fact that $|\mathcal{V}|^2 \gg |\mathcal{E}|$, we 267 have the total additional complexity of GSL by summing up all these terms: $O(|\mathcal{V}|^2F + |\mathcal{V}|F^2)$. 268 Compared with the complexity in normal GCN $O(|\mathcal{E}|F + |\mathcal{V}|F^2)$ [4], this additional complexity 269 $O((|\mathcal{V}|^2 - |\mathcal{E}|)F)$ adds tremendous training time and grows exponentially with the number of nodes 270 in graphs, which is shown in our experiments. 271

5 Experiments

In this section, we examine the effectiveness of Graph Structure Learning (GSL) through extensive experiments. To explore GSL's impact on Graph Neural Networks (GNNs), we compare the performance of 450 GSL variants integrated with various GNN backbones in Section 5.1. Furthermore, we analyze the influence of different components on GSL through an ablation study of each GSL component in Section 5.2.

Settings. Our experiments include six popular GNNs as backbones: GCN [17], SGC [47], GraphSage 278 [10], and GAT [44], Mixhop [2], and ACMGNN [32]. The datasets used in our experiments 279 include heterophilous graphs: Squirrel, Chameleon, Actor, Texas, Cornell, Wisconsin, Roman-empire, 280 and Amazon-ratings [38, 42, 39], and homophilous graphs: Cora, PubMed, and Citeseer [52], 281 Minesweeper, Tolokers, and Questions [39]. We show more dataset details in Appendix D. The model 282 performance is measured by accuracy for multi-class datasets or AUC-ROC for binary-class datasets 283 on node classification tasks. We use 50%/25%/25% random splits for training/validation/test sets. 284 285 For each experiment, we report the mean and standard deviation across 10 splits.

5.1 Performance Comparison

We investigate the impact of GSL on GNNs by the comparison of GNNs and the corresponding GSL-enhanced GNNs (GNN+GSL). As GSL introduces significant variations in three key aspects, we aim to comprehensively evaluate all possible GSL configurations through a combination of

Table 1: Performance Comparison of MLP, GNNs and the corresponding GSL-enhanced GNNs. For each GNN backbone, the best-performing method is highlighted in red, while the second-best method is highlighted in blue.

Model	Construct	Fusion	Param Sharing	Mines.	D	A	Tolokers	0	Squirrel	Chameleon	A	Texas	Cornell	Wisconsin	Cora	CiteSeer	PubMed	Rank
		rusion	Param Snaring	79 55+1 23	Roman. 65 45+0 99	Amazon. 46 65±0 83	75 94+1 38	Questions 74 92+1 39	39 29+2 22	43 57+4 18	Actor 35 40+1 38	80 46±6 44		85 88+7 78		76.68±2.10	87 39+2 18	3 93
MLP GCN	None None	-	-	79.55±1.23	65.45±0.99 81.46±1.25	46.65±0.83 50.89±1.16	75.94±1.38 84.61±0.99	77.68±1.10		43.5/±4.18 43.24±3.86	35.40±1.38 34.34±1.17	73.08±8.68	73.78±7.34 67.03±10.54	85.88±7.78 78.24±8.32	87.97±1.80 87.97±1.51	76.68±2.10 76.75±2.30	87.39±2.18 89.47±0.64	1.36
GCN		{G'}		77 91±5.79	67.40±1.23	46.72±1.51		72.56±1.14	38.15±2.45	39.87±4.87	33.47±1.61	63.06±8.68	65.68±7.76	72.75±5.70	85.21±1.39	75.52±1.14	89.47±0.64 89.03±0.42	
GCN	cos-graph cos-graph	{g,g'}	$\theta_1 = \theta_2$	52.53±6.45	62.57±0.81	41.29±1.61	74.22±1.79		37.62±1.74		32.74±0.92	57.88±8.75	66.49±9.12	73.14±5.92	64.68±1.61	67.32±1.14	86.43±0.76	
GCN	cos-graph	[g,g']	$\theta_1 \neq \theta_2$	88.70±0.86	69.90±2.38	47.35±0.83	82.85±0.95	75.29±1.38		40.30±4.31		65.47±8.48	62.97±10.89	75.29±6.54	85.51±1.87	75.23±1.14	88.74±0.59	
GCN	cos-graph cos-node	{G'}	01 ≠ 02	85.57±6.63	68 24+2 49	47.56±1.32	77 26+1 44	74.16±1.80					61 08+8 16	71 18+6 98	86.06±1.95	75.76±1.19	88 92±0.59	
GCN	cos-node	{q,q'}	$\theta_1 = \theta_2$	52.53±6.45	62.57±0.81					39.78±4.00			66.49±9.12	73.14±5.92	64.68±1.61	67.32±1.89	86.43±0.76	
GCN	cos-node	{g,g'}	$\theta_1 \neq \theta_2$	89 17±0 68	72.63±1.45		82.91±0.97			39.94±4.49			63.24±9.47	73.92±7.51	85.69±1.73	75.49±1.42	88 72±0.70	4.29
GCN	kNN	{G'}	01702	82.89±6.66	68 44+0 83	47 13±1 00	78.92±1.79	73.90±1.73		40.22±3.82		63.03±8.53	61.35±9.28	72.16±7.41	86.08±1.62	75.56±1.42	88 59±0 58	5.93
GCN	kNN	$\{g, g'\}$	$\theta_1 = \theta_2$	52.53±6.45	62.57±0.81	41.29±1.61	74.22±1.79			39.78±4.00			66.49±9.12	73.14±5.92	64.68±1.61	67.32±1.89	86,43±0,76	9.39
GCN	kNN	{G,G'}	$\theta_1 \neq \theta_2$	88.96±0.73	72.44±1.61	47.06±0.83	83.10±0.80	75.61±1.19	37.63±1.93	40.18±4.76	33.84±1.94	63.87±9.68	62.16±9.77	75.49±7.29	85.82±1.55	75.50±1.30	88.54±0.55	5.00
MLP	None			79 55+1 23	65.45+0.99	46.65±0.83	75 94+1 38	74.92±1.39	39 29+2 22	43 57+4 18	35.40±1.38	80.46±6.44	73.78±7.34	85 88+7 78	87.97±1.80	76.68±2.10	87 39+2 18	3.71
SGC	None			83.45±4.47	78.04±0.69	51.38±0.68		77.39±1.23		42.35±4.10			70.27±9.91	80.59±5.13	88.10±1.89	77.52±2.20	89.39±0.62	1.57
SGC	cos-graph	{G'}	-	73.76±4.46	67.17±0.81	47.15±0.88	76.28±1.63	73.93±2.66	38.66±2.53	40.07±4.39	33.87±1.45	71.19±7.38	67.57±9.19	77.65±6.08	86.95±2.01	76.12±1.29	89.10±0.43	5.79
SGC	cos-graph	$\{g, g'\}$	$\theta_1 = \theta_2$	52.53±4.89	62.97±0.78	42.42±1.57	74.29±1.79	70.56±1.27	37.56±2.25	39.33±3.60			66.49±10.37	71.57±4.46	64.82±2.11	67.55±1.80	86.58±0.72	9.64
SGC	cos-graph	$\{g,g'\}$	$\theta_1 \neq \theta_2$	79.70±1.21	62.02±2.06	47.24±0.93	83.22±1.52	77.19±0.99	38.32±1.80	40.85±4.61	33.51±1.50	70.34±7.31	64.86±9.01	75.29±6.82	87.47±1.70	75.70 ± 1.28	88.65±0.49	6.14
SGC	cos-node	$\{\mathcal{G}'\}$		79.03±3.76	67.84±1.87	47.93±0.94	78.09±1.84	75.46±1.43	38.61±2.20	40.50±4.10	34.03±1.27	70.08 ± 6.84	68.11±9.23	77.45±4.63	87.47±1.86	76.36±1.27	89.37±0.41	4.54
SGC	cos-node	$\{\mathcal{G},\mathcal{G}'\}$	$\theta_1 = \theta_2$	52.53±4.89	62.97±0.78	42.42±1.57	74.29±1.79	70.56±1.27	37.56±2.25	39.33±3.60	32.85±0.90	57.60±7.53	66.49±10.37	71.57±4.46	64.82±2.11	67.55±1.80	86.58±0.72	9.57
SGC	cos-node	$\{G, G'\}$	$\theta_1 \neq \theta_2$	80.12±1.36	66.90±1.66	48.04±0.97		77.11±1.09		40.20±4.66	34.20±1.79	68.47±8.11	64.59±9.74	75.29±6.05	87.54±1.63	75.88 ± 1.26	88.68±0.43	5.11
SGC	kNN	$\{G'\}$	-	75.53±4.98	67.94±0.70	47.68±0.84	79.45±2.06	74.22±2.47		39.92±3.91		72.81 ± 6.15	70.00±7.98	77.84±6.02	87.82±1.77	76.54±1.44	89.19±0.42	
SGC	kNN	$\{G, G'\}$	$\theta_1 = \theta_2$	52.53±4.89	62.97±0.78	42.42±1.57	74.29±1.79		37.56±2.25	39.33±3.60			66.49±10.37	71.57±4.46	64.82±2.11	67.55±1.80	86.58±0.72	9.50
SGC	kNN	$\{G, G'\}$	$\theta_1 \neq \theta_2$	80.78±1.08	64.59±1.93	47.48±0.99	83.17±1.43	76.80±1.09	36.53±2.06	40.17±4.24	34.23±1.72	69.26±6.77	65.95±8.87	76.08±5.92	87.38±1.49	76.02±1.22	88.77±0.45	5.79
MLP	None	-	-	79.55±1.23	65.45±0.99	46.65±0.83	75.94±1.38	74.92±1.39	39.29±2.22	43.57±4.18	35.40±1.38	80.46±6.44	73.78±7.34	85.88±7.78	87.97±1.80	76.68±2.10	87.39±2.18	4.14
SAGE	None	-	-	90.66±0.88	85.02±0.97	52.93±0.83	83.31±1.12	75.95±1.41		42.95±5.37	34.83±1.20	80.17±6.90	75.68±7.52	86.27±6.67	88.13±1.77	76.65±2.00	89.18±0.65	1.71
SAGE	cos-graph	$\{G'\}$	-	80.39±4.66	70.13±1.05	47.55±1.17	76.77±1.28		39.03±2.69		34.75±1.39	70.91 ± 8.58	70.00±7.56	78.24±6.87	83.64±2.03	75.53 ± 1.36	89.18±0.35	6.07
SAGE	cos-graph	$\{G, G'\}$	$\theta_1 = \theta_2$	53.02±6.49	59.98±1.73	39.99±2.29	71.57±2.28	66.01±3.58		38.49±3.68		60.30±7.05	67.57±4.59	76.47±5.92			85.53±0.51	9.93
SAGE	cos-graph	$\{G, G'\}$	$\theta_1 \neq \theta_2$	90.67±0.66	79.02±1.21	52.10±0.84	82.17±0.89			40.64±6.06		76.08±6.30	70.27±6.62	79.41±5.71	83.60±1.78	74.39 ± 1.35	88.88±0.50	3.86
SAGE	cos-node	$\{G'\}$	-	85.26±4.64	71.25±1.76	48.96±0.87		73.01±1.11					68.11±7.87	75.49±6.32	84.88±1.90		89.17±0.35	
SAGE	cos-node	$\{G, G'\}$	$\theta_1 = \theta_2$	53.02±6.49	59.98±1.73	39.99±2.29	71.59±2.28		35.05±2.41	38.49±3.68		60.30±7.05	67.57±4.59	76.47±5.92	64.58±1.74	67.77±1.31	85.53±0.51	9.79
SAGE	cos-node	$\{G, G'\}$	$\theta_1 \neq \theta_2$	90.64±0.65	78.60±0.98	52.08±0.90	82.02±0.88		39.18±2.54	40.86±6.17			69.73±7.43	80.00±5.68	83.96±1.65	74.63±1.26	88.93±0.64	3.93
SAGE	kNN	$\{\mathcal{G}'\}$		82.86±3.14	70.74±0.80	48.40±1.01	78.12±2.17		38.93±2.84	39.68±5.40		70.91±9.05	68.92±6.88	75.69±6.73	84.40±1.75	75.68±1.43	88.86±0.44	6.50
SAGE SAGE	kNN kNN	$\{g,g'\}$ $\{g,g'\}$	$\theta_1 = \theta_2$	53.02±6.49	59.98±1.73 79.16±1.15	39.99±2.29 51.56±1.07	71.59±2.28 81.66±0.87	66.01±3.58 75.22±0.97	35.05±2.41 39.20±2.39	38.49±3.68 40.44±5.82	31.32±1.04 35.13±1.38	60.30±7.05 74.17±6.31	67.57±4.59 70.54±7.32	76.47±5.92 79.61±6.61	64.58±1.74 84.05±1.63	67.77±1.31 74.59±1.25	85.53±0.51 88.67±0.55	9.86 4.57
		{9,9}	$\theta_1 \neq \theta_2$	90.61±0.63														
MLP	None	-	-	79.55±1.23	65.45±0.99	46.65±0.83	75.94±1.38		39.29±2.22	43.57±4.18		80.46±6.44	73.78±7.34	85.88±7.78	87.97±1.80	76.68±2.10	87.39±2.18	3.86
GAT	None	-	-	90.41±1.34	84.51±0.84	52.00±2.84	84.37±0.96	77.78±1.27		43.83±3.66		75.28±8.12	65.41±12.14	77.84±7.41	88.02±1.92	76.77±2.02	89.21±0.67	2.04
GAT	cos-graph	$\{G'\}$		80.78±8.24		45.79±1.10	74.84±1.84			40.21±3.53		62.73±9.06	67.57±7.03	77.06±7.29	86.03±1.85	75.46±1.49	88.63±0.59	
GAT	cos-graph	$\{g, g'\}$	$\theta_1 = \theta_2$	53.16±7.93	63.67±1.08	44.83±2.04	73.46±1.07		37.14±2.13	39.85±2.87	32.06±1.12	57.03±8.70	67.30±4.67	75.10±5.85	64.84±1.45	67.82±0.62	86.47±0.66	
GAT	cos-graph	$\{G, G'\}$	$\theta_1 \neq \theta_2$	89.97±0.80 87.64±8.40	76.08±1.70 68.80±2.39	49.61±0.73 46.37±1.06	82.75±0.90 77.77±1.86	77.13±1.20 73.65±1.47		40.40±3.30			66.76±7.23 65.41±8.48	78.82±6.76	86.60±1.75	75.05±1.36 75.59±1.49	87.85±0.72 88.59±0.49	
	cos-node	{g'}				46.37±1.06 44.83±2.04	73.46±1.07		38.65±2.46	40.33±3.25	33.43±0.94	64.64±9.09		75.10±6.13	87.08±1.66			9.46
GAT	cos-node cos-node	$\{g, g'\}$	$\theta_1 = \theta_2$	53.16±7.93 90.03±0.78	63.67±1.08	44.83±2.04 50.36±0.70	73.46±1.07 82.72±1.16	68.92±1.53 76.83±1.16	37.14±2.13	39.85±2.87 40.56±3.77	32.06±1.12 33.49±1.35	57.03±8.70 70.39±7.34	67.30±4.67 65.95±6.77	75.10±5.85 78.63±6.59	64.84±1.45 86.64±1.78	67.82±0.62 75.32±1.04	86.47±0.66 87.87±0.61	9.46 4.21
GAT	kNN	$\{G, G'\}$ $\{G'\}$	$\theta_1 \neq \theta_2$	90.03±0.78 84.27±5.25	68.73±1.47	46.05±0.70		71.58±1.62			33.49±1.35 33.84±1.07	61.68±8.71	65.95±6.77 62.97±7.43	74.90±5.86	86.64±1.78 86.77±1.90	75.52±1.04 75.64±1.45	87.87±0.61 88.29±0.48	
GAT	knn			84.27±5.25 53.16±7.93	68.73±1.47 63.67±1.08	46.05±0.90 44.83±2.04	73.46±1.75		38.82±2.33 37.14±2.13			57.03±8.70	62.97±7.43 67.30±4.67	75.10±5.85	64.84±1.45	67.82±0.62	86.47±0.66	
GAT	knn	$\{g,g'\}$ $\{g,g'\}$	$\theta_1 = \theta_2$ $\theta_1 \neq \theta_2$		77.23±1.63								67.30±4.67 65.95±6.52			75.20±1.55		
JAI	n. tit	(n, n)	v ₁ ≠ v ₂	G7.7G±0.79	//.EJ.11.03	45.75±0.72	02.7010.93	70.07±1.13	57.0522.70	**.**±3.92	JJ.J-4±1.30	70.30±1.22	05.75±0.32	77.0-1.7.23	O0.71±1.73	73.20 ± 1.33	07.77.10.31	7.70

various GSL components, which include (1) five GSL bases: original features X, aggregated features $\hat{A}X$, MLP-pretrained features MLP(X), GCN-pretrained features GCN(X, A), GCL (Graph Contrastive Learning)-pretrained features [62] GCL(X, A); (2) three similarity-based graph construction methods: graphs are constructed via cosine similarity of GSL bases with threshold from the graph level (cos-graph) and node level (cos-node), and k-nearest neighbors (kNN); and (3) three view fusion methods: early fusion $\{\mathcal{G}'\}$, late fusion $\{\mathcal{G},\mathcal{G}'\}$ with parameter sharing $\theta_1 = \theta_2$ or not $\theta_1 \neq \theta_2$. To ensure a fair comparison of the performance between GNN+GSL, GNN, and MLP, we consider all five GSL bases as input choices and train all models on each GSL bases. The details of all these modules can be found in Appendix E.

Table 1 reports the performance of MLP, GNN baselines, and GNN+GSL across eight datasets under the best of five GSL bases. Notably, under fair comparison conditions, all six baseline GNNs outperform their GNN+GSL counterparts⁴. This suggests that **the incorporation of GSL does not consistently yield performance improvements of GNNs, and in some situations, it even lead to worse results.** Besides, under the same search space of GSL bases, MLP outperforms most GNN+GSL in average rank. This result verifies the dilemma in Section 4.1 that **high-quality GSL bases already provide informative node representations without newly constructed graphs.** Since GSL-based methods may require specific training procedures or more complex model designs, we further examine the performance of state-of-the-art (SOTA) GSL approaches to evaluate the potential of GSL in Appendix F.3, where the results also indicate GSL makes no significant improvement.

As previously mentioned, besides boosting model performance, GSL is often used to enhance the robustness of GNNs [15]. Therefore, under our proposed framework, we conduct fair experiments to study the robuseness of GSL-enhanced GNNs with the same GSL search space. Figure 4 demonstrates the performance of GNNs alongside their GSL-enhanced counterparts on perturbed graphs, incorporating feature noise, edge addition, and edge removal, as suggested by (author?) [26]. The curves for GSL-enhanced GNNs (dotted lines) is close to those of the original GNNs (solid lines) across three types of perturbations and four GNN backbones, indicating that the baseline GNNs perform comparably to their GSL-enhanced versions. Therefore, the similarity-based graph construction may not be indispensable for enhancing model robustness. See Appendix F.8 for more details.

5.2 Ablation Study on Each GSL Component

Since the performance of GNN and GNN+GSL models is comparable under the same bases, we further investigate how different components of GSL influence GNNs in Figure 5, where each result is the averaged performance of four GNN backbones, including GCN, GAT, SGC, and GraphSAGE.

⁴Due to page limitation, results of the other two heterophily-oriented GNNs are shown in Appendix F.5, where we can derive the same conclusion as in Table 1.

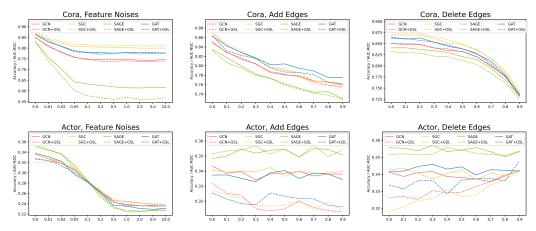


Figure 4: Response to feature noise, edge additions, and edge removals in GNN baselines and their GSL-enhanced counterparts.

The results indicate that: (1) Pretrained node representations, such as $MLP(\mathbf{X})$ and $GCN(\mathbf{X}, \mathbf{A})$, significantly enhance GNN performance 5 , (2) GSL graph generation has minimal impact on model performance, (3) two view fusion with parameter separation improves GNN performance, and (4) early fusion generally outperforms late fusion. Especially, GSL bases influence model performance most among all the GSL components, verifying our analysis in Section 4 that the quality of GSL bases greatly influences GNN performance, while graph construction has little impact.

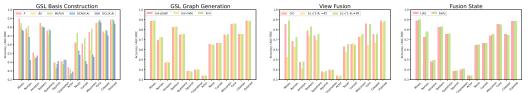


Figure 5: The influences of different GSL components on GNN+GSL.

6 Conclusion

In this paper, we disentangle the impact of GSL in GNN performance through our proposed GSL framework. Motivated by the dilemma associated with GSL, we show that it is the pretrained node features that really improve GNN performance instead of the similarity-based graph construction methods. Our research contributes to a deeper understanding of GSL and provides insights for re-evaluating essential components in future GNN designs. Although this paper primarily focuses on the impact of GSL on model performance in node classification tasks, future research could expand this analysis to other graph-related tasks and different types of graphs, as well as theoretically examine the effects of GSL under broader assumptions.

⁵See more discussion of GSL bases in Appendix F.2.1.

References

- 338 [1] Aug 2024.
- [2] S. Abu-El-Haija, B. Perozzi, A. Kapoor, N. Alipourfard, K. Lerman, H. Harutyunyan,
 G. Ver Steeg, and A. Galstyan. Mixhop: Higher-order graph convolutional architectures
 via sparsified neighborhood mixing. In *international conference on machine learning*, pages
 21–29. PMLR, 2019.
- [3] N. J. Beaudry and R. Renner. An intuitive proof of the data processing inequality. *Quantum Info. Comput.*, 12(5–6):432–441, May 2012.
- ³⁴⁵ [4] D. Blakely, J. Lanchantin, and Y. Qi. Time and space complexity of graph convolutional networks. *Accessed on: Dec*, 31:2021, 2021.
- [5] S. Brody, U. Alon, and E. Yahav. How attentive are graph attention networks? In *International Conference on Learning Representations*, 2021.
- [6] Y. Chen, L. Wu, and M. Zaki. Iterative Deep Graph Learning for Graph Neural Networks:
 Better and Robust Node Embeddings. In *Advances in Neural Information Processing Systems*,
 volume 33, pages 19314–19326. Curran Associates, Inc., 2020.
- [7] B. Fatemi, L. El Asri, and S. M. Kazemi. SLAPS: Self-Supervision Improves Structure Learning
 for Graph Neural Networks. In *Advances in Neural Information Processing Systems*, volume 34,
 pages 22667–22681. Curran Associates, Inc., 2021.
- [8] L. Franceschi, M. Niepert, M. Pontil, and X. He. Learning Discrete Structures for Graph Neural
 Networks, June 2020. arXiv:1903.11960 [cs, stat].
- [9] S. Gerchinovitz, P. Ménard, and G. Stoltz. Fano's inequality for random variables. 2020.
- 158 [10] W. Hamilton, Z. Ying, and J. Leskovec. Inductive representation learning on large graphs.

 Advances in neural information processing systems, 30, 2017.
- [11] D. He, C. Liang, H. Liu, M. Wen, P. Jiao, and Z. Feng. Block Modeling-Guided Graph
 Convolutional Neural Networks, Dec. 2021. arXiv:2112.13507 [cs].
- [12] C. Hua, S. Luan, M. Xu, Z. Ying, J. Fu, S. Ermon, and D. Precup. Mudiff: Unified diffusion for complete molecule generation. In *Learning on Graphs Conference*, pages 33–1. PMLR, 2024.
- [13] C. Hua, B. Zhong, S. Luan, L. Hong, G. Wolf, D. Precup, and S. Zheng. Reactzyme: A benchmark for enzyme-reaction prediction. In *The Thirty-eight Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2024.
- [14] B. Jiang, Z. Zhang, D. Lin, J. Tang, and B. Luo. Semi-Supervised Learning With Graph
 Learning-Convolutional Networks. In 2019 IEEE/CVF Conference on Computer Vision and
 Pattern Recognition (CVPR), pages 11305–11312, June 2019. ISSN: 2575-7075.
- [15] W. Jin, Y. Ma, X. Liu, X. Tang, S. Wang, and J. Tang. Graph structure learning for robust
 graph neural networks. In *Proceedings of the 26th ACM SIGKDD International Conference* on Knowledge Discovery & Data Mining, KDD '20, page 66–74, New York, NY, USA, 2020.
 Association for Computing Machinery.
- 174 [16] K. Z. Khanam, G. Srivastava, and V. Mago. The homophily principle in social network analysis: A survey. *Multimedia Tools and Applications*, 82(6):8811–8854, 2023.
- [17] T. N. Kipf and M. Welling. Semi-supervised classification with graph convolutional networks. arXiv preprint arXiv:1609.02907, 2016.
- 18] L. Kolbeck, S. Vilgertshofer, J. Abualdenien, and A. Borrmann. Graph rewriting techniques in engineering design. *Frontiers in built environment*, 7:815153, 2022.
- [19] L. F. Kozachenko and N. N. Leonenko. Sample estimate of the entropy of a random vector.
 Problemy Peredachi Informatsii, 23(2):9–16, 1987.

- [20] A. Kraskov, H. Stögbauer, and P. Grassberger. Estimating mutual information. *Physical Review E—Statistical, Nonlinear, and Soft Matter Physics*, 69(6):066138, 2004.
- D. Lao, X. Yang, Q. Wu, and J. Yan. Variational Inference for Training Graph Neural Networks in Low-Data Regime through Joint Structure-Label Estimation. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, KDD '22, pages 824–834, New York, NY, USA, Aug. 2022. Association for Computing Machinery.
- K. Li, Y. Liu, X. Ao, J. Chi, J. Feng, H. Yang, and Q. He. Reliable Representations Make A
 Stronger Defender: Unsupervised Structure Refinement for Robust GNN. In *Proceedings of the* 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, pages 925–935,
 Washington DC USA, Aug. 2022. ACM.
- [23] S. Li, D. Kim, and Q. Wang. Restructuring graph for higher homophily via adaptive spectral
 clustering. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pages
 8622–8630, 2023.
- 395 [24] X. Li, L. Sun, M. Ling, and Y. Peng. A survey of graph neural network based recommendation in social networks. *Neurocomputing*, 549:126441, 2023.
- [25] X. Li, R. Zhu, Y. Cheng, C. Shan, S. Luo, D. Li, and W. Qian. Finding global homophily in graph neural networks when meeting heterophily. In *International Conference on Machine Learning*, pages 13242–13256. PMLR, 2022.
- Z. Li, X. Sun, Y. Luo, Y. Zhu, D. Chen, Y. Luo, X. Zhou, Q. Liu, S. Wu, L. Wang, and J. X.
 Yu. GSLB: the graph structure learning benchmark. In A. Oh, T. Naumann, A. Globerson,
 K. Saenko, M. Hardt, and S. Levine, editors, Advances in Neural Information Processing
 Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS
 2023, New Orleans, LA, USA, December 10 16, 2023, 2023.
- [27] N. Liu, X. Wang, L. Wu, Y. Chen, X. Guo, and C. Shi. Compact Graph Structure Learning via
 Mutual Information Compression, Jan. 2022. arXiv:2201.05540 [cs].
- [28] Y. Liu, Y. Zheng, D. Zhang, H. Chen, H. Peng, and S. Pan. Towards Unsupervised Deep Graph
 Structure Learning, Jan. 2022. arXiv:2201.06367 [cs].
- [29] Q. Lu, S. Luan, and X.-W. Chang. Gcepnet: Graph convolution-enhanced expectation propagation for massive mimo detection. *In IEEE GLOBECOM 2024 Conference Proceedings*, 2024.
- 412 [30] Q. Lu, J. Zhu, S. Luan, and X.-W. Chang. Flexible diffusion scopes with parameterized laplacian for heterophilic graph learning. In *The Third Learning on Graphs Conference*, 2024.
- 414 [31] S. Luan, C. Hua, Q. Lu, L. Ma, L. Wu, X. Wang, M. Xu, X.-W. Chang, D. Precup, R. Ying, et al. The heterophilic graph learning handbook: Benchmarks, models, theoretical analysis, applications and challenges. *arXiv preprint arXiv:2407.09618*, 2024.
- [32] S. Luan, C. Hua, Q. Lu, J. Zhu, M. Zhao, S. Zhang, X.-W. Chang, and D. Precup. Revisiting
 heterophily for graph neural networks. *Advances in neural information processing systems*,
 35:1362–1375, 2022.
- 420 [33] S. Luan, C. Hua, M. Xu, Q. Lu, J. Zhu, X.-W. Chang, J. Fu, J. Leskovec, and D. Precup. When 421 do graph neural networks help with node classification? investigating the homophily principle 422 on node distinguishability. *Advances in Neural Information Processing Systems*, 36, 2024.
- 423 [34] S. Luan, M. Zhao, X.-W. Chang, and D. Precup. Break the ceiling: Stronger multi-scale deep graph convolutional networks. *Advances in neural information processing systems*, 32, 2019.
- [35] S. Luan, M. Zhao, C. Hua, X.-W. Chang, and D. Precup. Complete the missing half: Augmenting aggregation filtering with diversification for graph convolutional networks. In *NeurIPS 2022 Workshop: New Frontiers in Graph Learning*, 2022.
- [36] Y. Luo, L. Shi, and X.-M. Wu. Classic gnns are strong baselines: Reassessing gnns for node classification. *arXiv preprint arXiv:2406.08993*, 2024.

- 430 [37] Y. Ma, X. Liu, N. Shah, and J. Tang. Is homophily a necessity for graph neural networks? *arXiv* preprint arXiv:2106.06134, 2021.
- [38] H. Pei, B. Wei, K. C.-C. Chang, Y. Lei, and B. Yang. Geom-gcn: Geometric graph convolutional networks. *arXiv preprint arXiv:2002.05287*, 2020.
- 434 [39] O. Platonov, D. Kuznedelev, M. Diskin, A. Babenko, and L. Prokhorenkova. A critical look at the evaluation of gnns under heterophily: Are we really making progress? *arXiv preprint arXiv:2302.11640*, 2023.
- [40] L. Qiao, L. Zhang, S. Chen, and D. Shen. Data-driven graph construction and graph learning: A review. *Neurocomputing*, 312:336–351, 2018.
- 439 [41] B. C. Ross. Mutual information between discrete and continuous data sets. *PloS one*, 9(2):e87357, 2014.
- [42] B. Rozemberczki, C. Allen, and R. Sarkar. Multi-scale attributed node embedding. *Journal of Complex Networks*, 9(2):cnab014, 2021.
- [43] S. Suresh, V. Budde, J. Neville, P. Li, and J. Ma. Breaking the limit of graph neural networks by improving the assortativity of graphs with local mixing patterns. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*, pages 1541–1551, 2021.
- [44] P. Velicković, G. Cucurull, A. Casanova, A. Romero, P. Lio, and Y. Bengio. Graph attention networks. *arXiv preprint arXiv:1710.10903*, 2017.
- 448 [45] R. Wang, S. Mou, X. Wang, W. Xiao, Q. Ju, C. Shi, and X. Xie. Graph Structure Estimation
 449 Neural Networks. In *Proceedings of the Web Conference 2021*, pages 342–353, Ljubljana
 450 Slovenia, Apr. 2021. ACM.
- 451 [46] T. Wang, D. Jin, R. Wang, D. He, and Y. Huang. Powerful graph convolutional networks with adaptive propagation mechanism for homophily and heterophily. In *Proceedings of the AAAI*452 conference on artificial intelligence, volume 36, pages 4210–4218, 2022.
- [47] F. Wu, A. Souza, T. Zhang, C. Fifty, T. Yu, and K. Weinberger. Simplifying graph convolutional networks. In *International conference on machine learning*, pages 6861–6871. PMLR, 2019.
- [48] Q. Wu, W. Zhao, Z. Li, D. Wipf, and J. Yan. NodeFormer: A Scalable Graph Structure Learning
 Transformer for Node Classification, June 2023. arXiv:2306.08385 [cs].
- 458 [49] S. Wu, F. Sun, W. Zhang, X. Xie, and B. Cui. Graph neural networks in recommender systems: a survey. *ACM Computing Surveys*, 55(5):1–37, 2022.
- [50] S. Wu, Y. Tang, Y. Zhu, L. Wang, X. Xie, and T. Tan. Session-based recommendation with graph
 neural networks. In *Proceedings of the AAAI conference on artificial intelligence*, volume 33,
 pages 346–353, 2019.
- 463 [51] Y. Yan, M. Hashemi, K. Swersky, Y. Yang, and D. Koutra. Two sides of the same coin: Het-464 erophily and oversmoothing in graph convolutional neural networks. In 2022 IEEE International 465 Conference on Data Mining (ICDM), pages 1287–1292. IEEE, 2022.
- 466 [52] Z. Yang, W. Cohen, and R. Salakhudinov. Revisiting semi-supervised learning with graph 467 embeddings. In *International conference on machine learning*, pages 40–48. PMLR, 2016.
- Long Lang MarketterLong Mark
- 470 [54] X.-M. Zhang, L. Liang, L. Liu, and M.-J. Tang. Graph neural networks and their current applications in bioinformatics. *Frontiers in genetics*, 12:690049, 2021.
- T. Zhao, Y. Liu, L. Neves, O. Woodford, M. Jiang, and N. Shah. Data Augmentation for Graph Neural Networks, Dec. 2020. arXiv:2006.06830 [cs, stat].
- 474 [56] Y. Zheng, S. Luan, and L. Chen. What is missing in homophily? disentangling graph homophily for graph neural networks. *arXiv preprint arXiv:2406.18854*, 2024.

- 476 [57] Y. Zheng, J. Xu, and L. Chen. Learn from heterophily: Heterophilous information-enhanced graph neural network. *arXiv preprint arXiv:2403.17351*, 2024.
- Z. Zhiyao, S. Zhou, B. Mao, X. Zhou, J. Chen, Q. Tan, D. Zha, Y. Feng, C. Chen, and C. Wang.
 Opengsl: A comprehensive benchmark for graph structure learning. Advances in Neural
 Information Processing Systems, 36, 2024.
- [59] J. Zhu, Y. Yan, L. Zhao, M. Heimann, L. Akoglu, and D. Koutra. Beyond homophily in graph
 neural networks: Current limitations and effective designs. *Advances in Neural Information Processing Systems*, 33, 2020.
- 484 [60] Y. Zhu, W. Xu, J. Zhang, Y. Du, J. Zhang, Q. Liu, C. Yang, and S. Wu. A survey on graph structure learning: Progress and opportunities. *arXiv preprint arXiv:2103.03036*, 2021.
- ⁴⁸⁶ [61] Y. Zhu, W. Xu, J. Zhang, Q. Liu, S. Wu, and L. Wang. Deep graph structure learning for robust representations: A survey. *arXiv preprint arXiv:2103.03036*, 14:1–1, 2021.
- 488 [62] Y. Zhu, Y. Xu, F. Yu, Q. Liu, S. Wu, and L. Wang. Deep graph contrastive representation learning. *arXiv preprint arXiv:2006.04131*, 2020.
- [63] D. Zou, H. Peng, X. Huang, R. Yang, J. Li, J. Wu, C. Liu, and P. S. Yu. SE-GSL: A General and Effective Graph Structure Learning Framework through Structural Entropy Optimization. In *Proceedings of the ACM Web Conference 2023*, pages 499–510, Apr. 2023. arXiv:2303.09778 [cs].

NeurIPS Paper Checklist

502

503

504

505

506

507

508

509

520

521

522

523

524

525

526

527

528

530

531

532

533

534

535

536

537

538

539

541

The checklist is designed to encourage best practices for responsible machine learning research, 495 496 addressing issues of reproducibility, transparency, research ethics, and societal impact. Do not remove the checklist: The papers not including the checklist will be desk rejected. The checklist should 497 follow the references and follow the (optional) supplemental material. The checklist does NOT count 498 towards the page limit. 499

Please read the checklist guidelines carefully for information on how to answer these questions. For 500 each question in the checklist: 501

- You should answer [Yes], [No], or [NA].
- [NA] means either that the question is Not Applicable for that particular paper or the relevant information is Not Available.
- Please provide a short (1–2 sentence) justification right after your answer (even for NA).

The checklist answers are an integral part of your paper submission. They are visible to the reviewers, area chairs, senior area chairs, and ethics reviewers. You will be asked to also include it (after eventual revisions) with the final version of your paper, and its final version will be published with the paper.

The reviewers of your paper will be asked to use the checklist as one of the factors in their evaluation. 510 While "[Yes]" is generally preferable to "[No]", it is perfectly acceptable to answer "[No]" provided a 511 proper justification is given (e.g., "error bars are not reported because it would be too computationally 512 expensive" or "we were unable to find the license for the dataset we used"). In general, answering 513 "[No] " or "[NA] " is not grounds for rejection. While the questions are phrased in a binary way, we acknowledge that the true answer is often more nuanced, so please just use your best judgment and write a justification to elaborate. All supporting evidence can appear either in the main paper or the supplemental material, provided in appendix. If you answer [Yes] to a question, in the justification 517 please point to the section(s) where related material for the question can be found. 518

IMPORTANT, please: 519

- Delete this instruction block, but keep the section heading "NeurIPS Paper Checklist",
- Keep the checklist subsection headings, questions/answers and guidelines below.
- Do not modify the questions and only use the provided macros for your answers.

1. Claims

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

Answer: [Yes] Justification:

Guidelines:

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

2. Limitations

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

Answer: Yes 540 Justification:

Guidelines:

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

3. Theory assumptions and proofs

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

Answer: [Yes]
Justification:

Guidelines:

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

4. Experimental result reproducibility

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

Answer: [Yes]
Justification:
Guidelines:

The answer NA means that the paper does not include experiments.

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

5. Open access to data and code

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

Answer: [Yes]
Justification:
Guidelines:

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

• Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental setting/details

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

Answer: [Yes]
Justification:
Guidelines:

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

7. Experiment statistical significance

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

Answer: [Yes]
Justification:

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error
 of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments compute resources

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

Answer: [Yes]

Justification: We put it in the main paper

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.

- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code of ethics

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

Answer: [Yes]
Justification:

Guidelines:

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

10. Broader impacts

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

Answer: [NA]

Justification: We have carefully reevaluated our project and we do not have such negative impact.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

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

Answer: [NA]

Justification: We have carefully reevaluated our project and we do not have such risk.

Guidelines:

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

12. Licenses for existing assets

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

Answer: [Yes]
Justification:

Guidelines:

749

750

751

752

753

754

755

756

759

760

761

762

763

764

765

766

767

768

769

770

771

772

773

776 777

778

779

780

781

782

783

784

785

786

787

789

790

791

792

793

794

795

798

799

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

13. New assets

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

Answer: [NA] Justification:

Guidelines:

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

14. Crowdsourcing and research with human subjects

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

Answer: [NA]

Justification: We do not have such problem

Guidelines:

The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.

Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be

included in the main paper.
According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data

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

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

Answer: [Yes]

Justification: We have carefully reevaluated our project and we do not have such risk.

Guidelines:

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

16. Declaration of LLM usage

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

Answer: [NA]

Justification: LLM only for grammar check.

Guidelines:

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

A Taxonomy of Graph Structure Learning Methods

We present several representative GSL-based GNNs within our proposed GSL framework in Table 2.
Below, we provide a detailed description of each method.

Table 2: Representative GSL methods under our proposed GSL framework		α					CCI	
	Table 2. Redieschianve	\mathbf{u}	memous	unuci	Our	DIODOSCU	OOL.	manicwork

Method	Bases	Construct	Refinement	View Fusion	Training Mode
LDS [8]	X	$\{\mathcal{E}' = kNN(\mathbf{B})\}$ +Opt.	Bernoulli(\mathcal{E}')	Late Fusion, $\{G'_1, G'_2, \dots, G'_m\}$, $\theta_1 = \theta_2$	2-stage
Geom-GCN [38]	Isomap/Poincare/ Struc2vec(X, A)	$\{\mathcal{E}' e'_{ij}= oldsymbol{B}_i-oldsymbol{B}_j \}$	$threshold(\mathcal{E}')$	Late Fusion, $\{\mathcal{G},\mathcal{G}'\}$, $\theta_1 \neq \theta_2$	Static
ProGNN [15]	ϵ	$\{\mathcal{E}'=\mathrm{Opt}(\epsilon)\}$	Low Rank+Sparsity +Original	No Fusion, $\{\mathcal{G}'\}$	Joint
IDGL [6] GRCN [53]	MLP(X) GCN(X,A)	$\{\mathcal{E}' e'_{ij} = \cos(\mathbf{B}_i, \mathbf{B}_j)\}\$ $\{\mathcal{E}' e'_{ij} = \sigma(\mathbf{B}_i \mathbf{B}_j^T)\}$	$topk(\mathcal{E}')$ $topk(\mathcal{E}')$, $sym(\mathcal{E}')$	Early Fusion, $\{G + G'\}$ Early Fusion, $\{G + G'\}$	Joint Joint
GAug-M [55]	$\mathrm{GCN}^{(2)}(\boldsymbol{X}, \boldsymbol{A})$	$\{\mathcal{E}' e'_{ij} = \sigma(\mathbf{B}_i\mathbf{B}_j^T)\}$	$G'_{+} = topk(E'),$ $G'_{-} = bottom(E')$	Early Fusion, $\{G + G'_+ - G'\}$	Joint
GAug-O [55] SLAPS [7]	$oldsymbol{X} ext{MLP}(oldsymbol{X})$	$ \begin{aligned} \{\mathcal{E}' e'_{ij} &= p(e_{ij} \mathrm{GAE}(\boldsymbol{B},\boldsymbol{A}))\} \\ \{\mathcal{E}' &= kNN(\boldsymbol{B})\} \end{aligned} $	Gumbel(\mathcal{E}') norm(\mathcal{E}'),sym(\mathcal{E}')	Early Fusion, $\{G + G'\}$ No Fusion, $\{G'\}$	Joint Joint
CoGSL [27]	$GCN(X, \{A, kNN(X), PPR(X), Subgraph(X)\})$	$\{\mathcal{E}' e'_{ij} = p(e_{ij} \text{MLP}(\boldsymbol{B}, \boldsymbol{A}))\}$	-	Early Fusion, $\{\mathcal{G}^* \min \mathcal{L}_{CL}(\mathcal{G}, \mathcal{G}')\}, \theta_1 \neq \theta_2$	2-stage
GEN [45]	GCN(X, A)	$\{\mathcal{E}' = kNN(\mathbf{B})\}$	-	Late Fusion, $\{G_1', G_2', \dots, G_m'\}$, $\theta_1 \neq \theta_2$	2-stage
STABLE [22]	$\mathrm{GCL}(\boldsymbol{X},\boldsymbol{A})$	$\{\mathcal{E}' e'_{ij} = \cos(\mathbf{B}_i, \mathbf{B}_j)$ or $\cos(\mathbf{B}_i, \mathbf{B}_j)\}$	$G'_{+} = \text{topk}(\mathcal{E}'),$ $G'_{-} = \text{threshod}(\mathcal{E}')$	Early Fusion, $\{\mathcal{G} + \mathcal{G}'_+ - \mathcal{G}'\}$	Joint
SEGSL [63]	X	$\{\mathcal{E}' \min \mathcal{H}_S, e'_{ii} \in \text{EncTree}(kNN(\mathbf{B}))\}$	-	No Fusion, $\{\mathcal{G}'\}$	Joint
SUBLIME [28]	$\mathrm{GCN}(\boldsymbol{X},\boldsymbol{A})$	$\{\mathcal{E}' = \text{Opt}(\epsilon)\}\ \text{or}$ $\{\mathcal{E}' e'_{ij} = \cos/\text{Minkowski}(\boldsymbol{B}_i, \boldsymbol{B}_j)\}$	$topk(\mathcal{E}'),\!sym(\mathcal{E}'),\!norm(\mathcal{E}')$	Separation, $\{\mathcal{G}, \mathcal{G}'\}$, $\theta_1 = \theta_2$	Joint
BM-GCN [11]	$\hat{Y} = MLP(X),$ $\min \mathcal{L}_{CE}(\hat{Y}, Y)$	$\{\mathcal{E}' = oldsymbol{B}oldsymbol{Q}oldsymbol{B}^T\}$	$\text{norm}(\mathcal{E}')$	Early Fusion, $\{\mathcal{G}\odot\mathcal{G}'\}$	Joint
WSGNN [21]	MLP(X)	$\{\mathcal{E}' e'_{ij} = cos(\mathbf{B}_i, \mathbf{B}_j)\}$	-	Early Fusion, $\{G + G'\}$	Joint
GLCN [14]	\boldsymbol{X}	$\{E' e'_{ij} = \phi(B_i - B_j)\}$	$norm(\mathcal{E}')$, Original +Sparsity+Smoothness	No Fusion, $\{\mathcal{G}'\}$	Joint
ASC [23] WRGAT [43] HOG-GCN [46]	$ \begin{array}{c} \operatorname{SpectralCluster}(X) \\ \operatorname{GCN}(X,A) \\ \operatorname{GCN}(X,A) \end{array} $	$ \begin{aligned} \{\mathcal{E}' e'_{ij} &= \ \boldsymbol{B}_i - \boldsymbol{B}_j\ \} \\ \{\mathcal{E}' e'_{ij} \cdot Opt(\boldsymbol{B}) \} \\ \{\mathcal{E}' e^{\prime}_{ij} &= \sigma(\boldsymbol{B}_i\boldsymbol{B}_j^T) \} \end{aligned} $	$topk(\mathcal{E}')$ Sparsity + MultiHop Sparsity + Smoothness	No Fusion, $\{\mathcal{G}'\}$ Early Fusion $\{\mathcal{G} + \mathcal{G}'\}$ No Fusion $\{\mathcal{G}'\}$	Static Static Joint
GGCN [51] GloGNN [25]	MLP(X) $MLP(X)$	$\{\mathcal{E}' e'_{ij} = \cos(\mathbf{B_i}, \mathbf{B_j})\}\$ $\{\mathcal{E}' = \mathrm{Opt}(\mathbf{B})\}$	Low Rank + Sparsity Sparsity+MultiHop	Early Fusion, $\{G + G'\}$ No Fusion, $\{G'\}$	Joint Joint
HiGNN [57]	$\hat{Y} = GCN(\hat{X}, A),$ $\min \mathcal{L}_{CE}(\hat{Y}, Y)$	$\{\mathcal{E}'=e'_{ij}=\cos(\boldsymbol{B_i},\boldsymbol{B_j}))\}$	$topk(\mathcal{E}')$, $sym(\mathcal{E}')$	Late Fusion, $\{\mathcal{G},\mathcal{G}'\}$, $\theta_1 \neq \theta_2$	Static

LDS [8]. The GSL bases in LDS is constructed as node features X and the GSL graph \mathcal{G}' is initialized using a k-Nearest-Neighbors algorithm based on B. Then, \mathcal{G}' is updated with a loss function of node classification. Then multiple graphs are sampled based on \mathcal{G}' with a Bernoulli function and used to update the model parameters. The \mathcal{G}' construction and model parameters are updated as a 2-stage mode.

Geom-GCN [38]. Geom-GCN constructs the GSL bases from several graph-aware node embedding strategies using both of the X and A: Isomap [], Poincare [], and struc2vec []. Then, new graphs are constructed by filtering node pairs with a higher similarity measured by Euclidean distance $\{\mathcal{E}'|e'_{ij}=|B_i-B_j|<\delta\}$ where δ is a threshold. Finally, both of the aggregated message from \mathcal{G} and \mathcal{G}' are fused after applying graph convolution layers with no parameter sharing. The \mathcal{G}' is not updated through the training process.

ProGNN [15]. The \mathcal{G}' in ProGNN is purely learned by optimization without GSL bases. It optimizes the \mathcal{G}' using low rank, sparsity, and similarity with the original graphs \mathcal{G} . It outputs a single graph \mathcal{G}' without fusion and updates the \mathcal{G}' together with model parameters.

IDGL [6]. The GSL bases in LDS is constructed by linear transformation of node features MLP(X). Then, a GSL graph \mathcal{G}' is constructed using cosine similarity with topk threshold refinement. The early fusion is applied by fusing GSL graph \mathcal{G}' with original graph \mathcal{G} before training. The GSL graph \mathcal{G}' is trained with model parameters jointly.

GRCN [53]. GRCN constructs GSL bases by node embeddings of graph convolution GCN(X, A). Then, the GSL graph \mathcal{G}' is constructed by a kernel function with topk and symmetrization refinement $\{\mathcal{E}'|e'_{ij}=\sigma(B_iB_j)>\delta\}$. The final graph is obtained by early fusion and the GSL graph \mathcal{G}' is updated together with model parameters.

GAug-M and **GAug-O** [55]. GAug-M constructs GSL bases using a 2-layer graph convolution $GCN^{(2)}(X, A)$. Then, the GSL graph \mathcal{G}' is constructed by a kernel function. The final graph is obtained by adding some edges with highest probabilities and removing some edges with lowest probabilities on \mathcal{G} . GAug-O selects node features as GSL bases X, then trains a Graph Auto-Encoder to predict edges as \mathcal{G}' . Then, after gumbel sampling, the GSL graph \mathcal{G}' is fused with original graph \mathcal{G} before training. The \mathcal{G}' in both of the GAug-M and GAug-O is updated together with model parameters.

SLAPS [7]. SLAPS constructs the GSL bases by applying MLP(X) followed by a k-nearest neighbors (kNN) algorithm based on node feature similarities. The GSL graph \mathcal{G}' is then processed by an adjacency processor that symmetrizes and normalizes the adjacency matrix to ensure non-negativity and symmetry. The final graph is obtained of the generated graph \mathcal{G}' with the node features without fusion.

Additionally, a self-supervised denoising autoencoder $L_{DAE} = L(X_i, GNN_{DAE}(\hat{X}_i; \theta_{GNN_{DAE}}))$ is introduced to address the supervision starvation problem, updating \mathcal{G}' together with the model parameters.

 CoGSL [27]. CoGSL constructs GSL bases using two views, one of them is the Origin graph. Another is selected from the Adjacency matrix A, Diffusion matrix PPR(X), the KNN graph KNN(X) and the Subgraph of the Origin. GCNs are applied to these views to obtain node embeddings. The GSL graph is constructed by applying a linear transformation to the node embeddings of each node pair to estimate the connection probability between them. This connection probability is then added to the original view to finalize the graph. The refinement $\mathcal{E}'|e'_{ij} = p(e_{ij}|\text{MLP}(\mathbf{B},\mathbf{A}))$ step involves maximizing the mutual information between the two selected views and the newly constructed graph. InfoNCE loss is used to optimize the connection probability, where the same node serves as a positive sample, and different nodes serve as negative samples. The final graph \mathcal{G}' is obtained via early fusion of the selected views, and the GSL graph is updated with model parameters.

GEN [45]. GEN constructs the GSL bases by generating kNN graphs though several GCN layer, utilizing node representations from different layers. These kNN graphs are then combined using a Stochastic Block Model (SBM) to create a new graph \mathcal{G}' . The GSL graph \mathcal{G}' is refined iteratively through Bayesian inference to maximize posterior probabilities $P(G, \alpha, \beta|O, Z, Y_l) = \frac{P(O|G,\alpha,\beta)P(G,\alpha,\beta)P(O,Z,Y_l)}{P(O,Z,Y_l)}$, considering both the original graph and node embeddings. The final graph is obtained by feeding the graph Q back into the GCN for further optimization. The iterative process updates both the GSL graph and GCN parameters as a 2-stage mode, providing mutual reinforcement between the graph estimation and model learning.

STABLE [22]. STABLE constructs the GSL bases by generating augmentations based on node similarity through kNN graph and perturbing edges to simulate adversarial attacks. The GSL graph \mathcal{G}' is constructed by refining the structure using contrastive learning between positive samples (slightly perturbed graphs) and negative samples (undesirable views generated by feature shuffling). The refinement step applies a top-k filtering strategy on the node similarity matrix to retain helpful edges while removing adversarial ones. The final graph is obtained through early fusion, and the GSL graph \mathcal{G}' is updated together with model parameters during joint training

SE-GSL [63]. SE-GSL constructs the GSL bases using a kNN graph fused with the original graph. The GSL graph \mathcal{G}' is constructed through a structural entropy minimization process that extracts hierarchical community structures in the form of an encoding tree. The final graph is optimized by sampling node pairs from the encoding tree and generating new edges based on the minimized entropy structure. The refined graph is then used for downstream tasks, and the GSL graph \mathcal{G}' is updated jointly with model parameters during training.

SUBLIME [28]. SUBLIME constructs the GSL bases using both an anchor view (original graph) and a learner view (new graph). The new graph is initialized through kNN and further optimized either by parameter-based methods (using models like MLP, GCN, or GAT) or by non-parameter-based approaches (using cosine similarity or Minkowski distance). After obtaining the new graph, post-processing operations such as top-k filtering, symmetrization, and degree-based regularization are applied to ensure the graph's sparsity and structure. The GSL graph \mathcal{G}' is refined by applying contrastive learning between the anchor and learner views, incorporating edge drop and feature masking to generate node embeddings. The final graph is used in downstream tasks, and both views are updated together with model parameters in a joint training process.

BM-GCN [11]. BM-GCN constructs the GSL bases by introducing soft labels for nodes enbedding $\mathbf{B} = softmax(\sigma(MLP(X)))$ via a multilayer perceptron $\mathcal{L}_{MLP} = \sum_{v_i \in \mathcal{V}} f(B_i, Y_i)$. These soft labels are then used to compute a block matrix (H), which models the connection probabilities between different node classes. The GSL graph \mathcal{G}' is constructed by creating a block similarity matrix $Q = HH^T$ from the block matrix $Y_s = Y_i, B_i | \forall v_i \in \mathcal{T}_y, \forall v_j \notin \mathcal{T}_y, H = (Y_s^T A Y_s) \circ (Y_s^T A E)$, reflecting similarities between classes. The new graph is optimized using BQB^T and further fused with the original graph $A + \beta I$ for downstream tasks. The final graph is obtained by optimizing \mathcal{G}'

through degree-based regularization and top-k filtering. The GSL graph \mathcal{G}' is updated together with model parameters during joint training.

WSGNN [21]. WSGNN introduces a two-branch graph structure learning method, where each branch operates on different aspects of the graph: Branch AZ learns node labels from the new graph structure, while Branch ZA learns the new graph structure from the labels. The GSL bases is constructed using the observed graph A_{obs} and node features X. The new graph A' is inferred via cosine similarity between node embeddings. After constructing two separate views from each branch, the final graph is obtained by averaging the graphs from both branches. The refinement process ensures sparsity through cosine-based edge calculation $\mathcal{E}'|e'_{ij} = cos(\mathbf{B}_i, \mathbf{B}_j)$. Finally, both views undergo early fusion, with graph structure and node labels optimized jointly using a composite loss function that includes ELBO for structure prediction and cross-entropy loss for label prediction. The final GSL graph \mathcal{G}' is updated during joint training.

GLCN [14]. GLCN constructs the GSL bases by computing pairwise distances between node features and passing them through an MLP to obtain a block similarity score. This score is then processed with a softmax function to generate an $n \times n$ probability matrix that serves as the learned graph structure. The graph is refined using regularization techniques to ensure sparsity and feature smoothness $L_{GL} = \sum_{i,j=1}^{n} ||x_i - x_j||_2^2 S_{ij} + \gamma ||S||_F^2 + \beta ||S - A||_F^2$. The learned graph is then used for downstream graph tasks, where the task loss and the graph regularization loss are jointly optimized during joint training

ASC [23]. ASC constructs the GSL bases is formed by using pseudo-eigenvectors from spectral clustering. They divide the Laplacian spectrum into slices, with each slice corresponding to an embedding matrix. The GSL graph \mathcal{G}' is constructed by adaptive spectral clustering, where pseudo-eigenvectors are weighted based on alignment with node labels Where $f_i^{\mathcal{Z}}$. For refinement, they apply top-K edge selection by minimizing node embedding distance and maximizing homophily argmin $\sum_{i,j\in V_Y} (d(f_i^Z, f_j^Z), 1(y_i, y_j))$. This final restructured graph is training without fusion. Finally, the GSL graph is updated together with the model parameters.

WRGAT [43]. WRGAT constructs the GSL bases using the node features and a weighted relational GNN (WRGNN) framework that fuses structural and proximity information. A multi-relational graph is built by assigning different types of edges based on the structural equivalence of nodes at various neighborhood levels. This framework adapts to both assortative and disassortative mixing patterns, which helps improve node classification tasks. The GSL graph \mathcal{G}' is refined through attention-based message passing across these relational edges, and early fusion of proximity and structural features is used. The GSL graph \mathcal{G}' is trained jointly with the model parameters to optimize the node classification task.

HOG-GCN [46]. HOG-GCN constructs the GSL bases by incorporating both topological information and node attributes to estimate a homophily degree matrix $S = BB^T$, $B = softmax(Z_m)$, $Z_m^{(l)} = \sigma(Z_m^{(l-1)}W_m^{(l)})$. The GSL graph \mathcal{G}' is constructed using a homophily-guided propagation mechanism, which adapts the feature propagation weights between neighborhoods based on the homophily degree matrix $Z^{(l)} = \sigma(\mu Z^{(l-1)}W_e^{(l)} + \xi \hat{D}^{(-1)}A_k \odot HZ^{(l-1)}W_n^{(l)})$. For refinement, the graph incorporates both k-order structures and class-aware information to model the homophily and heterophily relationships between nodes. The final graph is obtained through joint fusion of topological and attribute-based homophily degrees, and both graph structure and model parameters are updated during joint training.

GGCN [51]. GGCN constructs the GSL bases using node features and structural properties such as node-level homophily h_i and relative degree \bar{r}_i . It incorporates structure-based edge correction by learning new edge weights derived from structural properties like node degree, and feature-based edge correction by learning signed edge weights from node features, allowing for positive and negative influences between neighbors. The GSL graph \mathcal{G}' is constructed by combining signed and unsigned edge information, aiming to capture both homophily and heterophily. The refinement process uses edge correction and decaying aggregation to mitigate oversmoothing and heterophily problems. The final graph is updated with early fusion, and the GSL graph \mathcal{G}' is optimized during joint training

GloGNN [25]. GloGNN constructs its GSL bases using node embeddings derived from MLP, combining both low-pass and high-pass convolutional filters. A coefficient matrix $Z^{(l)}$ is used to characterize the relationship between nodes and is optimized to capture both feature and structural

similarities $H_X^{(0)}=(1-\alpha)H_X^{(0)}+\alpha H_A^{(0)}$. Refinement is achieved via top-k selection based on the multi-hop adjacency matrix, and the matrix is symmetrized. The final graph is obtained through 981 982 global aggregation of nodes, capturing both local and distant homophilous nodes. This graph is then 983 used in downstream tasks, where the GSL graph \mathcal{G}' is jointly optimized with the model parameters. 984 **HiGNN** [57]. HiGNN constructs its GSL bases by utilizing heterophilous information as node 985 neighbor distributions, which represent the likelihood of neighboring nodes belonging to different 986 classes $\mathcal{H}_u = [p_1, p_2, ..., p_c]$, where $p_i = \frac{|v|v \in \mathcal{N}_u, y_c = i|}{|\mathcal{N}_u|}$. A new graph structure \mathcal{G}' is constructed by linking nodes with similar heterophilous distributions using cosine similarity. The refinement 987 988 involves selecting top-k edges based on the similarity score and applying symmetrization. The final 989 graph is fused with the original adjacency matrix A and the newly constructed adjacency matrix 990 A' via late fusion during message passing, where the node embeddings from both A and A' are 991 combined with a balance parameter λ . The graph \mathcal{G}' and node embeddings are updated during static 992 training. 993

B Contextual Stochastic Block Models with Homophily

To study the behavior of GNNs, CSBM-H [33, 37] have been proposed to create synthetic graphs with a controlled homophily degree. Specifically, in CSBM-H, for a node u with label y, its features $X_u \in \mathbb{R}^M$ are sampled from a class-wised Gaussian distribution $X_u \sim N_{Y_u}(\mu_{Y_u}, \Sigma_{Y_u})$ with $\mu_{Y_u} \in \mathbb{R}^F$ and $\Sigma_{Y_u} \in \mathbb{R}^{F \times F}$, where each dimension of X_u is independent from each other, e.e., $\Sigma_{Y_u} = \text{diag}(\mathbb{R}^n_{\geq 0})$. Then, to generate graph structure $\mathcal G$ with given homophily degree h with the range of [0,1], the node u has the probability h to connect intra-class nodes and the probability h to connect inter-class nodes. After applying neighbor sampling, both of the node homophily h_{node} and edge homophily h_{edge} in $\mathcal G$ are approximately equal to h.

C Proof of Theorem

994

1003

Theorem 4.1 Given a graph $\mathcal{G} = \{\mathcal{V}, \mathcal{E}\}$ with node labels \mathbf{Y} and node features \mathbf{X} , the accuracy of graph convolution in node classification is upper bounded by the mutual information between the node label Y and the aggregated node features H:

$$P_A \le \frac{I(Y;H) + \log 2}{\log(C)} \tag{5}$$

1007 Proof. For an arbitrary node u, the aggregated node features can be derived as $H_u = \frac{1}{|\mathcal{N}_u|} \sum_{v \in \mathcal{N}_u} X_v$ 1008 following the graph convolution operation. For a classifier predicting labels based on H_u , we have 1009 $\hat{Y}_u = \operatorname{cls}(H_u)$. Consequently, the Markov chain $Y \to H \to \hat{Y}$ holds. By applying Fano's inequality 1010 [9], we obtain

$$H(Y|H) \le H_b(P_E) + P_E \log(C - 1) \tag{6}$$

where P_E represents the error rate and $H_b(\cdot)$ is the binary entropy function. Rearranging this inequality gives us a lower bound on P_E :

$$P_E \ge \frac{H(Y|H) - H_b(P_E)}{\log(C - 1)} \tag{7}$$

Since $H(Y|H) = H(Y) - I(Y;H) = \log(C) - I(Y;H)$ and $H_b(P_E) \le \log 2$, we can substitute these terms into the equation:

$$P_E \ge 1 - \frac{I(Y;H) + \log 2}{\log(C)} \tag{8}$$

Finally, by expressing the accuracy rate P_A , we find:

$$P_A = 1 - P_E \le \frac{I(Y; H) + \log 2}{\log(C)}$$
 (9)

1016 This concludes the proof of Theorem 4.1.

Proposition 4.2 Consider a graph $\mathcal{G} = \{\mathcal{V}, \mathcal{E}\}$ characterized by node labels Y and n-dimensional node bases $\mathbf{B} = \{B_1, B_2, \dots, B_n\}$ with C classes. Each base B_i is independent and follows a class-dependent Gaussian distribution, i.e., $B_i \sim \mathcal{N}(\mu_Y, \sigma_Y)$. A new graph $\mathcal{G}' = \{\mathcal{V}, \mathcal{E}'\}$ is generated using a non-parametric method based on the bases \mathbf{B} . For the aggregated bases \mathbf{B}' on \mathcal{G}' , we have inf $I(Y; \mathbf{B}') \leq \inf I(Y; \mathbf{B})$.

1022 *Proof.* Let's first consider the mutual information for i-th node base B_i . For a non-parametric GSL method, we have the probability that class k connects with class j as:

$$p_{k,j} = \frac{g(B_i^k, B_i^j)}{\sum_{q=1}^C g(B_i^k, B_i^q)}$$
(10)

where $g(\cdot)$ is a non-parametric measurement of the probability of new connections, such as cosine similarity or Minkowski Distance. Then, we can get aggregated bases from the new graph by the operation of graph convolution [37, 33]:

$$B_i^{\prime k} = \sum_{q=1}^{C} p_{k,q} B_i^q \tag{11}$$

Therefore, the Markow chain $Y \to B_i \to B_i'$ holds. From data processing inequality [3], we have

$$I(Y; B_i') \le I(Y, B_i) \tag{12}$$

To extend this conclusion to multi-dimensional variables, we apply the chain rule of mutual information

$$I(Y; \mathbf{B}) = I(Y; \{B_1, \dots, B_n\}) = \sum_{i=1}^{n} I(Y; B_i \mid \{B_1, \dots, B_{i-1}\})$$

$$I(Y; \mathbf{B}') = I(Y; \{B'_1, \dots, B'_n\}) = \sum_{i=1}^{n} I(Y; B'_i \mid \{B'_1, \dots, B'_{i-1}\})$$
(13)

Due to the property that conditioning reduces entropy, we have

$$I(Y; B_i \mid \{B_1, \dots, B_{i-1}\}) \ge I(Y; B_i)$$

$$I(Y; B'_i \mid \{B'_1, \dots, B'_{i-1}\}) \ge I(Y; B'_i)$$
(14)

1031 Thus, we have

$$\inf I(Y; \mathbf{B}) = \sum_{i=1}^{n} I(Y; B_i) \text{ and } \inf I(Y; \mathbf{B}') = \sum_{i=1}^{n} I(Y; B'_i)$$
 (15)

where inf represents infimum. Since $I(Y; B'_i) \leq I(Y, B_i)$ holds for each i, we have

$$\inf I(Y; \mathbf{B'}) < \inf I(Y; \mathbf{B}) \tag{16}$$

This concludes the proof of Proposition 4.2.

1034 D Dataset Details

The datasets used in our experiments include heterophilous graphs: Squirrel, Chameleon, Actor, Texas, Cornell, and Wisconsin [38, 42], homophilous graphs: Cora, PubMed, and Citeseer [52], and Minesweeper, Roman-empire, Amazon-ratings, Tolokers, and Questions [39]. The dataset statistics are shown in 3. The descriptions of all the datasets are given below:

Cora, Citeseer, and Pubmed datasets are widely used citation networks in graph structure learning research. In each dataset, nodes represent academic papers, while edges capture citation relationships

Table 3: Dataset Statistics

Dataset	#Nodes	#Edges	#Classes	#Features	Edge Homophily
Cora	2,708	5,278	7	1,433	0.81
Pubmed	19,717	44,324	3	500	0.80
Citeseer	3,327	4,552	6	3,703	0.74
Roman-empire	22,662	32,927	18	300	0.05
Amazon-ratings	24,492	93,050	5	300	0.38
Minesweeper	10,000	39,402	2	7	0.68
Tolokers	11,758	529,000	2	10	0.59
Questions	48,921	153,540	2	301	0.84
Cornell	183	295	5	1,703	0.30
Chameleon	2,277	36,101	5	2,325	0.23
Wisconsin	251	466	5	1,703	0.21
Texas	183	309	5	1,703	0.11
Squirrel	5,201	216,933	5	2,089	0.22
Actor	7,600	33,544	5	931	0.22

between them. The node features are bag-of-words vectors derived from the paper's content, and each node is assigned a label based on its research topic. These datasets offer a structured framework to evaluate GNN models on classification tasks within citation networks.

Roman-Empire is constructed from the Roman Empire Wikipedia article, with nodes representing words and edges formed by either word adjacency or dependency relations. It contains 22.7K nodes and 32.9K edges. The task is to classify words by their syntactic roles, and node features are fastText embeddings. The graph is chain-like, with an average degree of 2.9 and a large diameter of 6824. Adjusted homophily is low ($h_{adj} = -0.05$), making it useful for GNN evaluation under low homophily and sparse connectivity.

Amazon-Ratings is based on Amazon's product co-purchasing network, this dataset includes nodes as products (books, CDs, DVDs, etc.) and edges linking frequently co-purchased items. It consists of the largest connected component of the graph's 5-core. The goal is to predict product ratings grouped into five classes.

Minesweeper is a synthetic dataset resembling the Minesweeper game, nodes in a 100x100 grid represent cells, with edges connecting adjacent cells. The task is to identify mines (20% of nodes). Node features indicate neighboring mine counts, with 50% of features missing. The average degree is 7.88, and the graph has near-zero homophily due to random mine placement.

Tolokers is derived from the Toloka crowdsourcing platform, where nodes represent workers connected by shared tasks. The graph has 11.8K nodes and an average degree of 88.28. The task is to predict which workers have been banned, using profile and task performance features. The graph is much denser than others in the benchmark.

Questions is based on user interactions from Yandex Q, this dataset focuses on users interested in medicine. Nodes are users, and edges represent questions answered between users. It contains 48.9K nodes with an average degree of 6.28. The task is to predict user activity at the end of a one-year period, with fastText embeddings from user descriptions as features. The graph is highly imbalanced (97% active users).

Texas, **Wisconsin**, **Cornell** are part of the WebKB project, representing web pages from university computer science departments. Nodes correspond to web pages, and edges represent hyperlinks between them. The node features are bag-of-words vectors from the web page content, and the labels classify each page into one of five categories: student, project, course, staff, and faculty.

Chameleon, **Squirrel** are page-page networks based on specific topics from Wikipedia. Nodes represent web pages, and edges correspond to mutual links between them. Node features are derived from the page content, and the classification task is based on average monthly traffic. These datasets are characterized by high heterophily, making them challenging for traditional GNN models.

Actor is an induced subgraph from a film-director-actor-writer network. Nodes represent actors, and edges are created when two actors co-occur on the same Wikipedia page. The task is to classify actors into five categories based on the keywords associated with their Wikipedia pages.

E Implementation Details

We implement GSL on 6 baseline GNNs with a variety of GSL approaches from the perspective of GSL bases, GSL graph construction, and view fusion. The baseline GNNs include:

- GCN [17] performs layer-wise propagation of node features and aggregates information
 from neighboring nodes to capture local graph structures. Each layer applies a convolution
 operation to update node embeddings, combining the node's features with its neighbors.
- GAT [44] employs self-attention to learn dynamic attention coefficients between nodes and their neighbors. These coefficients are normalized using softmax, and the final node representation is computed as a weighted sum of the neighbor features. Multi-head attention is used to enhance stability and expressiveness, with the number of attention heads set to 8 by default in our experiments.
- **SAGE** [10] uses an inductive framework to aggregate features from a node's local neighborhood, allowing it to generalize to unseen nodes. The aggregation function, set to mean in our experiments, efficiently combines neighbor information at each layer.
- SGC [47] simplifies the GCN model by removing non-linear activations and collapsing multiple layers into a single linear transformation. This reduction in complexity accelerates training. Node features are propagated using precomputed matrices, making the model faster and more efficient. In our experiments, the number of k-hops in SGC is set to 2 by default.
- **MixHop** [2] extends traditional GNNs by allowing nodes to aggregate information from neighbors at multiple distances within a single layer. Instead of only considering immediate neighbors, MixHop raises the adjacency matrix to different powers, capturing diverse topological signals. In our experiments, we follow the original paper's setup by using three propagation levels.
- ACMGCN [32] introduces an adaptive channel mixing mechanism to dynamically learn
 and combine information from different channels of node features. By leveraging attentionbased feature transformation, ACMGCN enhances representation learning for graphs with
 diverse structural properties. In our experiments, we use the default channel mixing setup as
 described in the original paper.

The GSL bases B includes the following options:

- B = X: The original node features are used as the GSL bases.
- $B = \hat{A}X$: Aggregated node features from 1-hop neighbors, normalized by node degree, are used as the GSL bases.
- B = MLP(X): Pretrained MLP embeddings are used as the GSL bases. A 2-layer MLP is trained using node features and labels on the training set for 1000 epochs per run. The hidden layer size is set to 128, the learning rate to $1e^{-2}$, the dropout rate to 0.5, and the weight decay to $5e^{-4}$. All parameters are optimized with Adam. After training, node embeddings are extracted from the last hidden layer, with a dimension of 128, prior to classifier input.
- B = GCN(X, A): Pretrained node embeddings are obtained from a 2-layer GCN model, following the same training procedure as for the MLP embeddings.
- B = GCL(X, A): Pretrained node embeddings are derived from a Graph Contrastive Learning (GCL) model without supervision, following the same training process as the MLP embeddings. GRACE [62] is used as the GCL model, with 2 views and 2 layers. The edge and feature dropout rates in each view are set to 0.2.

The approaches for the construction of GSL graph \mathcal{G}' includes:

• Cos-graph: $\mathcal{G}' = \{e_{ij} | \cos(\boldsymbol{B_i}, \boldsymbol{B_j}) > \delta, i \in \mathcal{V}, j \in \mathcal{V}\}$. This method calculates the cosine similarity between all node pairs in the original graph \mathcal{G} . Node pairs with a similarity higher than the threshold δ are selected as the edge set for the GSL graph \mathcal{G}' .

- Cos-node: $\mathcal{G}' = \bigcup_{i \in \mathcal{V}} \{\{e_i j\} | \cos(\boldsymbol{B_i}, \boldsymbol{B_j}) > \delta_i, j \in \mathcal{N}_i\}$. Unlike Cos-graph, which operates at the graph level, Cos-node constructs \mathcal{G}' at the node level. To prevent nodes from being left without neighbors (which may occur in Cos-graph), Cos-node selects neighbors based on node-level cosine similarity, ensuring each node has sufficient connections.
 - kNN: $\mathcal{G}' = \text{kNN}(B)$. This method constructs a kNN graph using the k-Nearest Neighbors algorithm based on the GSL bases B.

1131 The view fusion in GSL includes:

1129

1130

1132

1133

1134

1135

1136

1137

1139

1140

1141

1142

1155

- $\{\mathcal{G}'\}$: This approach uses only the GSL graph \mathcal{G}' for subsequent GNN training, completely ignoring the original graph \mathcal{G} .
 - $\{\mathcal{G}, \mathcal{G}'\}, \theta_1 = \theta_2$. Both the GSL graph \mathcal{G}' and the original graph \mathcal{G} are used for GNN training, with parameter sharing across each layer of the GNN.
 - $\{\mathcal{G}, \mathcal{G}'\}, \theta_1 \neq \theta_2$. Both the GSL graph \mathcal{G}' and the original graph \mathcal{G} are used for GNN training, but with separate model parameters for each graph.

1138 Especially, for graphs with two views, the fusion stage in GSL includes:

- Early Fusion: G + G'. Combine the two graphs, G and G', into a single new graph prior to GNN training.
- Late Fusion: H + H'. After training the GNN on the original graph \mathcal{G} and the GSL graph \mathcal{G}' , merge the node embeddings, \mathbf{H} and \mathbf{H}' , before passing them to the classifiers.

In addition to the original models based on 4 baseline GNNs, we implement GNN+GSL (GSL-augmented GNNs) by combining the aforementioned GSL modules, resulting in multiple variants for each type of GNN. For all models, we explore hyperparameters including hidden dimensions from the set $\{64, 128, 256\}$, learning rates from $\{1e-2, 1e-3, 1e-4\}$, weight decay values from $\{0, 1e-5, 1e-3\}$, the number of layers from $\{2, 3\}$, and dropout rates from $\{0.2, 0.4, 0.6, 0.8\}$. All the experiments are conducted on a Linux server(Operation system: Ubuntu 16.04.7 LTS) with one NVIDIA Tesla V100 card.

For GSL graph generation, we also search for additional hyperparameters to ensure the performance quality of the GSL-augmented GNN. Specifically, for Cos-graph and Cos-node, we control the parameter δ to vary the ratio of the number of edges in \mathcal{G}' to the number of edges in \mathcal{G} across the set $\{0.1, 0.5, 1, 5\}$. For kNN, we investigate the number of neighbors from the set $\{2, 3, 5, 10\}$.

1154 F Additional Experiment Results

F.1 Impact of GSL Bases on GNN baselines

In Figure 6, we illustrate the influence of 5 GSL bases on the performance of 4 GNNs across both 1156 homophilous and heterophilous graphs. The results indicate that MLP-pretrained features, denoted as MLP(X), significantly enhance GNN performance compared to raw features X across 6 out of 9 1158 datasets. These improvements stem from the self-training process applied to node inputs, suggesting 1159 that various self-training strategies could be employed with different graph datasets to further enhance 1160 GNN performance. Many GSL-enhanced GNNs leverage trained GSL bases to improve model 1161 performance, whereas GNN baselines utilize raw node features as GSL bases for comparison. This 1162 raises concerns about the fairness of previous comparisons between GNNs using original node 1163 features and those employing GNN+GSL, underscoring the importance of high-quality GSL bases. 1164 Additionally, we observe that GCN and GCL-pretrained features tend to degrade GNN performance 1165 on heterophilous datasets. This degradation is attributed to the increased noise within heterophilous 1166 datasets, leading to lower-quality GSL bases that can negatively impact GNN performance. 1167

1168 F.2 Impact of each GSI component on GNN+GSL

1169 F.2.1 GSL Bases

In addition to the analysis of the impact of GSL bases shown in Figure 5, Figure 7 presents further results on the performance of various GSL bases (\mathbf{X} , $\mathbf{\hat{A}X}$, $MLP(\mathbf{X})$, $GCN(\mathbf{X}, \mathbf{A})$, $GCL(\mathbf{X}, \mathbf{A})$)

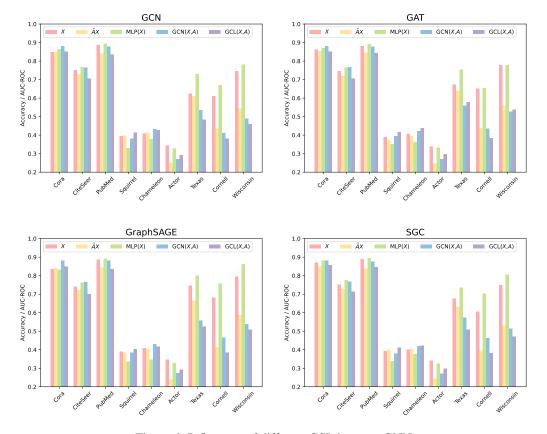


Figure 6: Influences of different GSL bases to GNNs.

across GAT, SGC, and GraphSAGE. The results are consistent with those observed in GCN and MLP, where the original node features do not always yield the best input. Some pretrained features, such as $MLP(\mathbf{X})$ on the Texas, Cornell, and Wisconsin datasets, demonstrate significant improvement compared to the original features \mathbf{X} , highlighting the necessity of self-training. Since many GSL methods [57, 43] utilize self-training during the training process, a fair comparison of these GSL methods and baseline GNNs should be conducted in the context of self-training, such as by using pretrained node features as input, as shown in Table 1.

F.2.2 GSL Graph Generation

Figure 8 compares the Cos-graph, Cos-node, and kNN methods for GSL graph generation. Across most datasets, the performance differences among these methods are minimal. In certain datasets, such as Roman-empire and Pubmed, the models exhibit comparable performance regardless of the graph generation technique employed. This suggests that variations in graph generation have a limited effect on overall performance.

F.2.3 View Fusion

Figure 9 illustrates the impact of different view fusion approaches, comparing the use of only the GSL graph \mathcal{G}' , the combination of the original graph \mathcal{G} with \mathcal{G}' using shared parameters $\theta_1 = \theta_2$, and the use of separate parameters $\theta_1 \neq \theta_2$. Notably, using only the GSL graph \mathcal{G}' underperforms compared to employing both graph views with separate model parameters. This indicates that incorporating information from the original graph \mathcal{G} is beneficial for maximizing GNN+GSL performance. Furthermore, for the two graph views, parameter sharing significantly underperforms parameter separation. We speculate that the messages aggregated under \mathcal{G} and \mathcal{G}' differ considerably, suggesting that different graphs should be treated with distinct model parameters.

F.2.4 Fusion Stage

Figure 10 compares early fusion and late fusion for GNN+GSL with multiple graph views. The performance difference between the two fusion states is often minimal. While early fusion tends to perform slightly better on complex datasets like Actor and Pubmed, the overall impact of switching between early and late fusion is limited across most datasets. For simpler datasets like Minesweeper and Amazon, both fusion methods yield nearly identical performance, indicating that the choice of fusion state does not drastically alter the model's outcome in most cases.

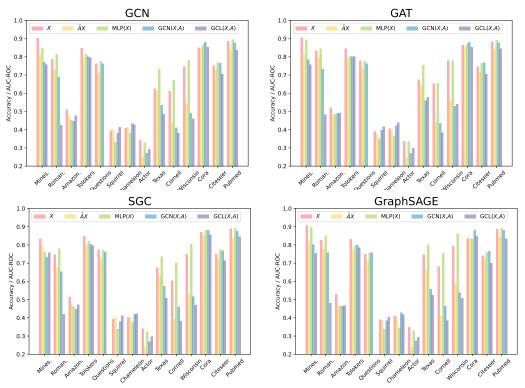


Figure 7: Influences of different GSL bases to more GNNs.

F.3 Removing GSL in SOTA GNNs

Settings To fairly reassess the impact of GSL in state-of-the-art (SOTA) methods, we compare the performance of SOTA models with their SOTA-GSL counterparts within the same hyperparameter search space. These GSL-based SOTA models include GAug [55], GEN [45], GRCN [53], IDGL [6], NodeFormer [48], GloGNN [25], WRGAT [43], and WRGCN [43]. Corresponding to the analysis of GCN and MLP in Section 4.2, the SOTA-GSL methods include two variants: (1) SOTA, $\mathcal{G}' = \mathcal{G}$, which replaces the GSL graph \mathcal{G}' with the original graph \mathcal{G} ; and (2) SOTA, $\mathcal{G}' = \text{MLP}$, which substitutes the graph convolution layers of GSL \mathcal{G}' with MLP layers. We train each model for 1000 epochs and search the hidden dimensions from the set {16, 32, 64, 128, 256, 512}, learning rate from {1e-1, 1e-2, 1e-3, 1e-4, 1e-5}, weight decay values from {5e-4, 5e-5, 5e-6, 5e-7, 0}, the number of layers from {1, 2, 3}, and dropout rates from {0.2, 0.4, 0.6, 0.8}. The hyperparameters of the above methods are shown in Table 5. The model-specific hyperparameters are shown as follows:

In **GRCN**, the hyperparameter K determines the number of nearest neighbors used to create a sparse graph from a dense similarity graph which helps balance efficiency and accuracy. We set the k as 5.

In **GAug**, the alpha is a hyperparameter that regulates the influence of the edge predictor on the original graph. We set the alpha as 0.1.

In **IDGL**, The parameter graph_learn_num_pers defines the number of perspectives for evaluating node similarities in the graph learning process. The parameter num_anchors specifies the number of

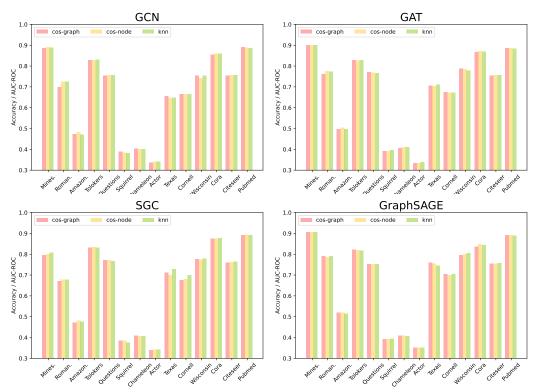


Figure 8: Influences of the approaches of GSL generation to GNN+GSL.

anchor points used to reduce computational complexity and improve scalability in graph structure learning. The graph_skip_conn parameter controls the proportion of skip connections, preserving information from the original graph during new graph structure learning. The update_adj_ratio parameter determines the proportion of the adjacency matrix updated at each iteration, influencing the dynamic adjustment of the graph structure. We set the graph_learn_num_pers as 6, num_anchors as 500, graph_skip_conn as 0.7, and update_adj_ratio as 0.3.

In **NodeFormer**, The parameter k determines the number of neighbors considered for each node in constructing the local graph structure, influencing the strength of node connections and the propagation of features. The parameter tolerance controls the degree of error tolerance during optimization. A larger tolerance allows more flexibility in the search space near local optima, while a smaller one results in stricter convergence. The number of attention heads in a graph attention network (GAT). Multi-head attention enables the model to focus on different subspace representations simultaneously, enhancing the diversity and stability of the representations. We set the k as 10, lambda as 0.01, and n_heads as 4.

In **GEN**,the parameter K in KNN refers to the number of nearest neighbors used to construct the graph structure, determining how many adjacent nodes are selected. The parameter tolerance defines the acceptable range of error during optimization, controlling the convergence criteria of the model. The parameter threshold determines the edge weight threshold in the graph, deciding which edges to retain in the graph structure. We set the k as 10, tolerance as 0.01, and threshold as 0.5.

In **GloGNN**, we set the Delta as 0.9, Gamma as 0.8, alpha as 0.5, beta as 2000, and orders as 5. Delta adjusts the balance between local and global node embeddings. Gamma controls the significance of global aggregation versus local information. Alpha balances the contributions of node features and graph structure. Beta regularizes the model, preventing overfitting. Order defines how many layers of neighbors are considered.

In **WRGAT**, we set the number of attention heads as 2 and the negative slope as 0.2. The number of attention heads determines how many attention mechanisms are used. The negative slope modifies the LeakyReLU activation.

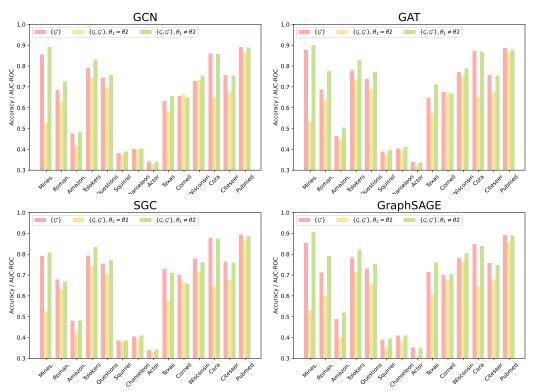


Figure 9: Influences of the approaches of view fusion in GSL to GNN+GSL.

Results. The results are presented in Table 4, where "OOM" denotes out-of-memory. It is evident that removing GSL does not diminish model performance; in fact, it is often comparable to or even exceeds the original results. Furthermore, GSL-based SOTA methods require significantly more GPU memory and longer running times compared to their non-GSL counterparts. Based on these findings, we conclude that GSL not only fails to enhance performance across most datasets but also increases model complexity. In conjunction with the results in Table 1, we assert that GSL may be unnecessary for effective GNN design in most cases.

Table 4: Model Performance and training time per epoch of SOTA methods and SOTA-GSL. The results for methods marked with "*" are reported in (author?) [58].

	Questio	ns	Mineswe	eper	Roman-en	pire	Amazon-ra	atings	Toloke	rs	Cora		Pubme	d	Citeseer	
Model	AUC	Time	AUC	Time	Acc	Time	Acc	Time	AUC	Time	Acc	Time	Acc	Time	Acc	Time
$GAug^*$ GAug, $G' = GGAug$, $G' = MLP$	OOM OOM OOM	-	77.93±0.64 80.56±0.36 64.31±1.40	- 11s 4.8s	OOM OOM OOM	-	48.42±0.39 48.45±0.37 48.05±0.66	12s 37s	OOM OOM OOM	-	82.48±0.66 81.73±0.38 78.90±0.00	7s 1s 3.2s	78.73±0.77 79.38±0.46 77.40±0.00	20s 6s 8.1s	72.79±0.86 72.34±0.18 72.91±0.32	10s 2s 9s
GEN* GEN, $G' = G$ GEN, $G' = MLP$	OOM OOM OOM	-	79.56±1.09 80.81±0.23 71.81±0.98	260s 75s 12s	OOM OOM OOM	-	49.17±0.68 50.08±0.30 49.29±0.65	130s 49s	OOM OOM OOM	-	81.66±0.91 82.16±0.37 80.20±0.00	214s 39s 140s	78.49±3.98 80.49±0.13 66.80±0.00	1384s 114s 1592s	73.21±0.62 71.52±0.34 73.50±0.00	470s 25s 310s
GRCN* $GRCN, G' = G$ $GRCN, G' = MLP$	74.50±0.84 75.69±0.52 63.59±2.35	8s 3.9s	72.57±0.49 71.15±0.05 72.18±1.09	60s 10s 2s	44.41±0.41 45.84±0.52 45.89±0.83	180s 8s 7.5s	50.06±0.38 46.07±1.02 48.77±0.60	220s 10s 8.1s	71.27±0.42 71.73±0.42 70.45±1.39	37s 10s 8s	84.61±0.34 81.66±1.10 79.40±0.00	13s 2s 1.3s	79.30±0.34 79.35±0.26 78.10±0.00	17s 3s 5s	72.34±0.34 69.55±1.28 71.40±0.00	20s 2s 4.2s
$IDGL^*$ IDGL, $G' = GIDGL$, $G' = MLP$	OOM OOM OOM	-	50.00±0.00 50.00±0.00 79.56±1.26	157s 51s 13.7s	47.10±0.65 41.24±0.86 50.35±0.36	186s 42s 35s	45.87±0.58 OOM 39.93±0.88	- 15s	50.00±0.00 50.00±0.00 71.55±1.08	279s 52s 11s	84.19±0.61 82.43±0.45 83.20±0.00	123s 13s 6.6s	82.78±0.44 73.50±1.85 79.20±0.00	146s 23s 13s	73.26±0.53 73.13±0.49 72.60±0.00	332s 36s 13.9s
NodeFormer* NodeFormer, $G' = G$ NodeFormer, $G' = MLP$	OOM OOM OOM	-	77.29±1.71 80.66±0.82 80.04±1.42	215s 21s	56.54±3.73 68.37±1.95 53.08±2.37	236s 7.2s	41.33±1.25 OOM 71.55±1.08	- 26s	OOM OOM OOM	-	78.81±1.21 77.01±1.99 78.82±0.00	213s 152s 8s	78.38±1.94 OOM 76.30±0.00	- 127s	70.39±2.04 70.82±0.13 72.80±0.00	219s 139s 15s
GloGNN GloGNN, $G' = G$ GloGNN, $G' = MLP$	68.67±1.07 68.32±1.23 69.69±0.22	66.6s 49.4s 25.7s	52.45±0.30 52.30±0.21 52.30±0.20	13.0s 3.6s 2.1s	66.21±0.17 66.03±0.14 66.49±0.16	26.1s 15.3s 12.4s	50.72±0.88 50.23±0.83 49.56±0.73	31.1s 21.7s 12.3s	79.81±0.20 80.02±0.16 74.85±0.12	47.4s 25.1s 2.8s	78.07±1.66 73.49±2.01 73.93±1.81	6.6s 5.1s 3.2s	87.88±0.26 87.62±0.20 87.64±0.27	18.2s 14.4s 10.2s	71.95±1.90 72.27±2.08 72.09±1.81	21.8s 21.2s 13.8s
WRGAT WRGAT, $G' = G$ WRGAT, $G' = MLP$	OOM 74.67±0.95 68.07±2.62	64.1s 75.8s	90.22±0.64 89.79±0.37 87.08±2.11	168.0s 18.6s 16.2s	OOM OOM OOM	-	OOM 50.41±0.53 41.38±1.46	49.9s 24.4s	78.69±1.21 78.81±0.89 76.41±1.25	153.0s 47.0s 37.7s	84.28±1.52 83.48±1.48 76.99±1.10	19.5s 3.4s 2.9s	88.82±0.50 88.92±0.43 80.27±6.23	421.6s 26.5s 23.9s	73.50±1.41 73.22±1.90 65.28±2.11	22.1s 4.7s 4.5s
WRGCN WRGCN, $G' = G$ WRGCN, $G' = MLP$	74.70±1.71 75.91±1.30 64.59±1.48	358.3s 43.3s 23.1s	90.63±0.64 90.65±0.49 70.66±1.37	40.9s 5.5s 7.7s	OOM OOM OOM	-	52.76±0.95 52.54±0.56 37.05±0.46	508.4s 50.1s 8.0s	82.68±0.82 82.65±0.86 69.10±0.91	52.3s 15.6s 12.2s	88.30±1.46 88.32±0.79 70.00±3.59	23.7s 3.9s 2.2s	OOM 89.26±0.45 67.29±2.49	- 19.4s 9.9s	73.74±1.60 74.45±1.51 70.84±1.36	54.2s 10.5s 4.1s

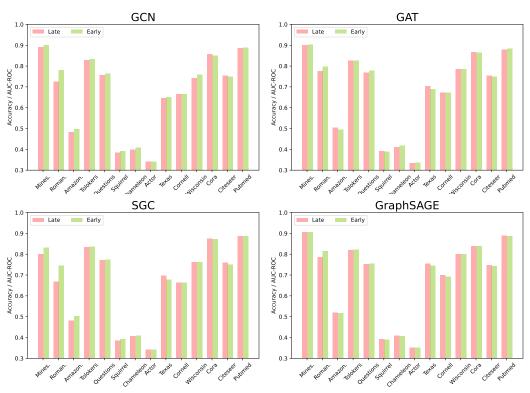


Figure 10: Influences of the states of view fusion in GSL to GNN+GSL.

F.4 Quality of GSL Graphs

Previous studies [25, 57] suggest that GSL constructs graphs with properties that improve intra-class node connectivity, which can be measured by homophily. This improvement can be visualized by inspecting graph structures with nodes sorted by their class labels. A graph that appears closer to a block diagonal matrix indicates stronger intra-class connectivity. However, this enhancement may not always be essential and can be achieved through non-GSL methods as well. In Figure 11, we visualize the original and reconstructed structures of a heterophilous graph from the Wisconsin dataset. The GSL graphs are constructed using various bases: \mathbf{X} , $\mathbf{\hat{A}}\mathbf{X}$, $\mathbf{MLP}(\mathbf{X})$, $\mathbf{GCN}(\mathbf{X},\mathbf{A})$, and $\mathbf{GCL}(\mathbf{X},\mathbf{A})$. We also include reconstructed graphs using a simple method that samples edges between nodes of the same class based on label predictions, *i.e.*, $\hat{\mathbf{Y}} = \mathbf{GCN}(\mathbf{X},\mathbf{A})$ or $\hat{\mathbf{Y}} = \mathbf{MLP}(\mathbf{X},\mathbf{A})$. Figure 11 demonstrates that, although GSL improves intra-class connectivity, the improvement is not as substantial as that achieved by non-GSL methods, as seen in the last two subfigures. Thus, the improvement in homophily within GSL graphs is unnecessary, as it can be easily achieved through simple methods.

F.5 Heterophily-oriented GNN with GSL

We also include heterophily-oriented GNNs, specifically ACMGNN [32] and MixHop [2], in our experiments that incorporate GSL into GNN baselines. These experiments follow the same setup as described in Table 1. The results, presented in Table 6, demonstrate that, under fair comparison conditions, both ACMGNN and MixHop outperform their GNN+GFS counterparts. This suggests that adding GSL to these heterophily-oriented GNNs may be unnecessary.

F.6 Trainable GSL

In Table 7, we present the results of applying trainable GSL to baseline GNNs. Specifically, we select the best-performing GSL variants, as shown in Tables 1 and 6, for each backbone GNN. The best-performing method is highlighted in bold, while the runner-up is indicated with an underline.

Table 5: Hyperparameters for GSL-enhanced SOTA methods and their counterparts by replacing or removing new graphs.

Color Colo	remo	ving new g	rapns.											
Color Colo			Learning Rate	Weight Decay	Dropout		Num of Layers	Dataset	Model	Learning Rate	Weight Decay	Dropout	Hidden Dim	Num of Layers
CRICKY FMLP 16-3 5-3 0.5 256 26 2 16 16 16 16 16 16		GAug, $G' = G$ GAug, $G' = MLP$ GEN GEN, $G' = G$	1e-4 1e-4 1e-2 1e-2	5e-7 5e-7 5e-4 5e-4	0.8 0.8 0.5 0.5	512 512 16 16	2		GAug, $G' = G$ GAug, $G' = MLP$ GEN GEN, $G' = G$	1e-2 1e-2 1e-3 1e-3	5e-4 5e-4 5e-4 5e-4	0.5 0.5 0.2 0.2	128 128 32 32	2 2 2 2 2 2
Charles Cha		GRCN GRCN, $G' = G$ GRCN, $G' = MLP$ IDGL	1e-3 1e-3 1e-3 1e-2	5e-3 5e-3 5e-3 5e-4	0.5 0.5 0.5 0.5	256 256 256 512	2 2 2 2 2		GRCN GRCN, $G' = G$ GRCN, $G' = MLP$ IDGL	1e-3 1e-3 1e-3 1e-2	5e-3 5e-3 5e-3 5e-4	0.5 0.5 0.5 0.5	32 32 32 16	2 2 2 2 2
Checker	Cora	IDGL, $G' = MLP$ NodeFormer NodeFormer, $G' = G$ NodeFormer, $G' = MLP$	1e-2 1e-2 1e-2 1e-2	5e-4 5e-4 5e-4 5e-4	0.5 0.2 0.2 0.2	512 64 64 64	2 2 2 2 2 1	PubMed	IDGL, $G' = MLP$ NodeFormer NodeFormer, $G' = G$ NodeFormer, $G' = MLP$ GloGNN	1e-2 1e-3 1e-3 1e-3	5e-4 5e-4 5e-4 5e-4	0.5 0.2 0.2 0.2	16 64 64 32	2 2 2 2 2 2 3
Distance		Glognn, $G' = MLP$ WRGAT WRGAT, $G' = G$ WRGAT, $G' = MLP$	1e-2 1e-2 1e-2 1e-2	5e-5 1e-5 5e-5 1e-5	0.5 0.5 0.5 0.5	64 128 128 128	1 1 2 2 2 2		GloGNN, $G' = G$ GloGNN, $G' = MLP$ WRGAT WRGAT, $G' = G$ WRGAT, $G' = MLP$	1e-3 1e-3 1e-2 1e-2 1e-2	5e-5 5e-5 5e-5 1e-5 5e-5	0.7 0.7 0.5 0.5 0.5	64 64 64 64	3 2 2 2 2
Glorg		WRGCN, $G' = G$	1e-2	5e-5	0.5	128	2		WRGCN, $G' = G$	1e-2	5e-5	0.5	64	2 2
$ \begin{array}{c c c c c c c c c c c c c c c c c c c $	Dataset							Dataset	Model	Learning Ra	te Weight Decay	Dropout	Hidden Dim	Num of Layers
GRCN, G = MLP		GAug, $G' = G$ GAug, $G' = MLP$ GEN GEN G' = G	1e-4 1e-4 1e-2 1e-2	5e-7 5e-7 5e-4 5e-4	0.8 0.8 0.5 0.5	512 512 16 16	2 2 2 2		GAug, $G' = GGAug$, $G' = MLPGEN$, $G' = G$	1e-3 1e-3 1e-4 1e-4	5e-6 5e-6 5e-6 5e-6 5e-6	0.8 0.8 0.8 0.8	256 256 256 256 256	3 3 3 3
Discrete		GRCN GRCN, $G' = G$ GRCN, $G' = MLP$ IDGL	1e-3 1e-3 1e-2 1e-2	5e-3 5e-3 5e-3 5e-4	0.8 0.8 0.5 0.5	512 512 256 32	3 3		GRCN GRCN, $G' = G$ GRCN, $G' = MLP$ IDGL	1e-3 1e-3 1e-1	5e-7 5e-6 5e-6 5e-6	0.2 0.2 0.2 0.2	128 128 128 128	3 2 2 2 3 3
GIGCHN,	Citeseer	IDGL, $G' = MLP$ NodeFormer NodeFormer, $G' = G$ NodeFormer, $G' = MLP$ GloGNN	1e-3 1e-2 1e-2 1e-2	5e-4 5e-4 5e-4 5e-4	0.2 0.2 0.2	16 64 64 64		Mineswee	per $IDGL, \mathcal{G}' = MLP$ NodeFormer NodeFormer, $\mathcal{G}' = \mathcal{G}$ NodeFormer, $\mathcal{G}' = ML$ GloGNN	1e-1 1e-2 1e-2 P 1e-2	5e-6 5e-4 5e-4 5e-4 5e-4	0.2 0.8 0.8 0.8	128 32 32 32 32	3 2 2 2 2 5
WRCCN, g' = g 1e-2 5e-5 0.5 128		Glognn, $G' = G$ Glognn, $G' = MLP$ WRGAT WRGAT, $G' = G$ WRGAT, $G' = MLP$	1e-2 1e-2 1e-2 1e-2 1e-2	1e-5 1e-5 5e-5 5e-5 5e-5	0.7 0.5 0.5 0.5	64 64 128 128 128	2 2 2 2 2 2		GloGNN, $G' = MLP$ WRGAT WRGAT, $G' = G$ WRGAT, $G' = MLP$	1e-2 1e-2 1e-2 1e-2	5e-4 5e-5 5e-5 5e-5	0.5 0.5 0.5	512 128 128 128	5 5 2 2 2 2
$ \begin{array}{c c c c c c c c c c c c c c c c c c c $		WRGCN, $G' = G$	1e-2	5e-5	0.5	128	2		WRGCN, $G' = G$	1e-2	5e-5	0.5	128	2 2
$ \begin{array}{c c c c c c c c c c c c c c c c c c c $	Dataset							Dataset						Num of Layers
$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$		GAug, $G' = G$ GAug, $G' = MLP$ GEN GEN, $G' = G$ GEN, $G' = MLP$ GRCN	le-1 le-1 le-2 le-2 le-2 le-3	5e-5 5e-5 5e-7 5e-7 5e-7 5e-5	0.5 0.5 0.2 0.2 0.2 0.2	32 32 128 128 128 128	2 2 2 2 2 2		GAug, $G' = G$ GAug, $G' = MLP$ GEN GEN, $G' = G$ GEN, $G' = MLP$ GPCN	le-2 le-2 le-2 le-2 le-2 le-3	5e-7 5e-7 5e-7 5e-7 5e-7 5e-7	0.2 0.2 0.2 0.2 0.2 0.2	128 128 128 128 128 128 128	2 2 2 2 2 2 2 2 2
NodeFormer, G' = MLP 1e-3 5e-5 0.8 128 3 NodeFormer, G' = MLP 1e-4 5e-5 0.5 64	Roman-en	GRCN, $\mathcal{G}' = MLP$ IDGL IDGL, $\mathcal{G}' = \mathcal{G}$ IDGL, $\mathcal{G}' = \mathcal{G}$ IDGL, $\mathcal{G}' = MLP$ NodeFormer	1e-2 1e-1 1e-1 1e-3	5e-5 5e-5 5e-5 5e-5 5e-6	0.5 0.5 0.5 0.5	128 128 128 128 128	2 2 2 2 2 3 3	Amazon-r	GRCN, $G' = MLP$ IDGL IDGL, $G' = G$ IDGL, $G' = MLP$ NodeFormer	1e-2 1e-2 1e-2 1e-2 1e-4	5e-7 5e-7 5e-7 5e-7 5e-5	0.2 0.2 0.2 0.2	128 128 128 128 128	2 2 2 2 2 3 3
WRGAT, G' = MLP 1e-2 5e-5 0.5 128 2 WRGCN, G' = G 1e-2 5e-5 0.3 128 WRGCN, G' = G 1e-2 5e-5 0.5 128 2 WRGCN, G' = G 1e-2 5e-5 0.7 128 2 WRGCN, G' = MLP 1e-2 5e-5 0.7 128 2 WRGCN, G' = MLP 1e-2 5e-5 0.7 128 2 WRGCN, G' = MLP 1e-2 5e-5 0.7 128 2 WRGCN, G' = MLP 1e-2 5e-5 0.7 128 2 WRGCN, G' = MLP 1e-2 5e-5 0.5 32 WRGCN, G' = G 1e-2 5e-6 0.5 64 WRGCN, G' = G 1e-2 5e-5 0.3 128 WRGCN, G' = G 1e-2 5e-5 0.3 128 WRGCN, G' = G 1e-2 5e		NodeFormer, $G' = M$ GloGNN GloGNN, $G' = G$ GloGNN, $G' = MLP$ WRGAT	LP 1e-3 1e-2 1e-2 1e-2 1e-2	5e-5 5e-5 5e-5 5e-5 5e-5	0.8 0.7 0.7 0.7 0.5	128 128 128 128 128	3 3 3 3 2		NodeFormer, $G' = M$ GloGNN GloGNN, $G' = G$ GloGNN, $G' = MLP$ WRGAT	LP 1e-4 1e-2 1e-2 1e-2 1e-2	5e-5 5e-5 5e-5 5e-5 5e-5	0.5 0.3 0.3 0.3	64 128 128 128 128	3 3 3 3 2
$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$		WRGAT, $G' = MLP$ WRGCN WRGCN, $G' = G$	1e-2 1e-2 1e-2 1e-2	5e-5 5e-5 5e-5 5e-5	0.5 0.5 0.5 0.5	128 128 128 128	2 2		WRGAT, $G' = MLP$ WRGCN WRGCN, $G' = G$ WRGCN, $G' = MLP$	1e-2 1e-2 1e-2 1e-2	1e-5 5e-5 5e-5 1e-5	0.3 0.7 0.7 0.7	128 128 128 128	2 3 3 3
$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$	Dataset						Num of Layers	Dataset						Num of Layers
$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$		GAug, $\mathcal{G}' = \mathcal{G}$ GAug, $\mathcal{G}' = MLP$ GEN $\mathcal{G}' = \mathcal{G}$ GEN, $\mathcal{G}' = \mathcal{G}$ GEN, $\mathcal{G}' = MLP$ GRCN GRCN, $\mathcal{G}' = \mathcal{G}$ GRCN, $\mathcal{G}' = MLP$	1e-2 1e-2 1e-2 1e-2 1e-2 1e-2 1e-2 1e-2	5e-4 5e-4 5e-7 5e-7 5e-7 5e-6 5e-6 5e-6	0.5 0.2 0.2 0.2 0.2 0.5 0.5	64 64 256 256 256 64 64	3 2 2 2 2 2		GAug, $\mathcal{G}' = \mathcal{G}$ GAug, $\mathcal{G}' = MLP$ GEN, $\mathcal{G}' = \mathcal{G}$ GEN, $\mathcal{G}' = \mathcal{G}$ GEN, $\mathcal{G}' = MLP$ GRCN, $\mathcal{G}' = \mathcal{G}$ GRCN, $\mathcal{G}' = \mathcal{G}$	1e-1 1e-1 1e-2 1e-2 1e-2 1e-2 1e-2 1e-1	5e-5 5e-5 5e-6 5e-6 5e-5 5e-6 5e-6	0.5 0.5 0.2 0.2 0.2 0.5 0.5	32 32 128 128 128 32 32 34 64	2 2 2 2 2 2 2 2 2 2
WRGAT 5e-3 5e-5 0.3 64 2 WRGAT 1e-2 5e-5 0.5 128 WRGAT 3f' = g 5e-3 1e-5 0.3 64 2 WRGAT 3f' = g 1e-2 1e-5 0.5 128 WRGAT 3f' = MILP 1e-2 1e-5 0.5 128 WRGAT 4f' = MILP 1e-2 1e-5 0.5 128 WRGAT 5f' = MILP 1e-2 1e-5 0.5 128 WRGCN 5e-3 5e-5 0.7 64 2 WRGCN 1e-2 5e-5 0.5 128 WRGCN 1e-2 5e-5 0.7 64 2 WRGCN 1e-2 5e-5 0.5 128 WRGCN 1e-2 5e-5 0.7 64 2 WRGCN 1e-2 5e-5 0.5 128 WRGCN 1e-2 5e-5 0.7 64 2 WRGCN 1e-2 5e-5 0.5 128	Questions	IDGL IDGL, $\mathcal{G}' = \mathcal{G}$ IDGL, $\mathcal{G}' = MLP$ NodeFormer NodeFormer, $\mathcal{G}' = \mathcal{G}$ NodeFormer, $\mathcal{G}' = MLP$ GloGNN	le-2 le-2 le-4 le-4 le-4 le-2	5e-7 5e-7 5e-3 5e-3 5e-3 5e-5	0.2 0.2 0.5 0.5 0.5 0.7	128 128 128 64 64 128	3 3 3	Tolokers	IDGL IDGL, $\mathcal{G}' = \mathcal{G}$ IDGL, $\mathcal{G}' = \text{MLP}$ NodeFormer NodeFormer, $\mathcal{G}' = \mathcal{G}$ NodeFormer, $\mathcal{G}' = \text{MLP}$ GloGNN	1e-2 1e-2 1e-2 1e-2 1e-2 1e-2	5e-4 5e-4 5e-4 5e-4 5e-4 5e-5	0.5 0.5 0.2 0.2 0.2 0.2	64 64 64 64 64 128	2 2 2 2 2 2 2 3
		Glognn, $\mathcal{G}' = MLP$ WRGAT WRGAT, $\mathcal{G}' = \mathcal{G}$ WRGAT, $\mathcal{G}' = MLP$ WRGCN WRGCN, $\mathcal{G}' = \mathcal{G}$	1e-2 5e-3 5e-3 5e-3 5e-3 5e-3	5e-5 5e-5 1e-5 5e-5 5e-5	0.7 0.3 0.3 0.3 0.7 0.7	128 64 64 64 64 64	3 2 2 2 2 2 2		WRGAT WRGAT, $G' = G$ WRGAT, $G' = MLP$ WRGCN	1e-2 1e-2 1e-2 1e-2 1e-2 1e-2	5e-5 5e-5 1e-5 5e-5 5e-5 5e-5	0.3 0.5 0.5 0.5 0.5	128 128 128 128 128 128	3 2 2 2 2 1 2 2

"OOM" refers to "out of memory." The results demonstrate the following: (1) The average rank indicates that trainable GSL improves GNN performance on 5 out of 6 GNN backbones; (2) Although trainable GSL outperforms non-trainable GSL, it remains inferior to GNN backbones without GSL, indicating that GSL could be unnecessary in improving GNN performance on node classification.

F.7 Performance on Graph classification

1281

1282

1283

1285

In addition to the node classification experiments, we further investigate whether GSL consistently improves GNN performance in graph classification. Specifically, we conduct ablation experiments by replacing the GSL graph with the original graph, following the methodology outlined in [26]. As shown in 8, removing GSL from 4 state-of-the-art GNNs, including ProGNN [15], GEN [45], GRCN [53], and IDGL [6], results in significantly reduced training time. At the same time, the GNN

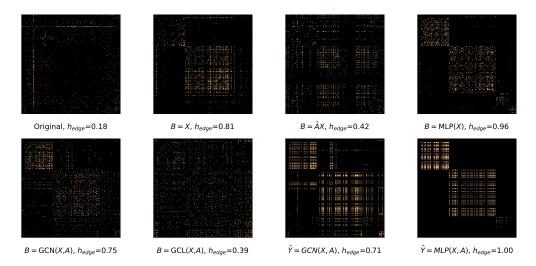


Figure 11: Visualization of original graph and reconstructed graphs on Wisconsin

Table 6: Performance of heterophily-oriented GNNs with GNN+GSL

Model	Construct	Fusion	Param Sharing	Mines.	Roman.	Amazon.	Tolokers	Questions	Squirrel	Chameleon	Actor	Texas	Cornell	Wisconsin	Cora	CiteSeer	PubMed	Rank
MLP	-	-	-	79.55±1.23	65.45±0.99	46.65±0.83	75.94±1.38	74.92±1.39	39.29±2.22	43.57±4.18	35.40±1.38	80.46±6.44	73.78±7.34	85.88±7.78	87.97±1.80	76.68±2.10	87.39±2.18	2.93
ACMGNN	-	-		90.56 ± 1.03	84.86±0.73	52.07±1.72	84.41±1.12	77.72±1.59	41.53±2.43	44.65±4.43	34.86±1.22	82.62±5.97	75.68±8.99	87.65±7.15	88.23±1.81	76.63±2.34	89.37±0.56	1.21
ACMGNN	cos-graph	$\{G'\}$		47.36±3.47	60.97±0.76	41.50±0.75	70.21±1.51	67.32±1.37	38.12±1.92	39.90±3.64	33.43±0.95	59.80±6.99	59.46±8.35	71.57±6.68	77.47±2.41	73.68±0.97	87.19±0.38	7.64
ACMGNN	cos-graph	$\{G, G'\}$	$\theta_1 = \theta_2$	52.74±5.22	51.18±2.12	33.11±1.38	69.06±4.65	62.30±3.23	31.58±4.39	38.79±4.73	29.06±2.60	54.10±7.59	59.19±8.87	70.39 ± 9.58	59.74±1.87	65.17±1.94	79.53±1.69	9.86
ACMGNN	cos-graph	$\{G, G'\}$	$\theta_1 \neq \theta_2$	87.46±1.02	74.63±0.76	49.35±0.58	81.63±0.87	73.84±1.41	38.54±1.89			67.67±5.97	70.00±5.90	80.78±5.21	80.83±1.84	73.43±1.47	88.98±0.47	3.64
ACMGNN	cos-node	$\{G'\}$	-	52.83±3.52	61.26±0.62	42.47±0.53	74.14±1.14	72.23±1.36	38.23±1.97	40.77±3.68		61.45±6.13	63.51±5.87	74.31±6.43	75.84±2.93	73.05±1.18	87.22±0.41	6.21
ACMGNN	cos-node	$\{\mathcal{G}, \mathcal{G}'\}$	$\theta_1 = \theta_2$	52.74±5.22	51.18±2.12	33.11±1.38	69.06±4.65	62.30±3.23	31.58±4.39	38.79±4.73	29.06±2.60	54.10±7.59	59.19±8.87	70.39 ± 9.58	59.74±1.87	65.17±1.94	79.53±1.69	9.86
ACMGNN	cos-node	$\{G, G'\}$	$\theta_1 \neq \theta_2$	87.80±0.97	73.55±0.51	49.04±0.57	80.74±0.92	74.11±1.40	39.19±2.12	40.28±4.30		69.86±5.56	69.46±7.21	80.39±5.23	80.33±1.90	73.31±1.26	88.94±0.36	4.07
ACMGNN	kNN	$\{G'\}$	-	51.68±3.38	60.86±0.87	41.68±0.95	71.31±0.64	69.56±1.41	38.58±1.96	40.56±2.34	34.88±0.77	62.51±6.16	62.70±5.95	76.47±4.43	75.99±2.85	70.20±1.51	87.20±0.45	6.64
ACMGNN	kNN	$\{G, G'\}$	$\theta_1 = \theta_2$	52.74±5.22	51.18±2.12	33.11±1.38	69.06±4.65	62.30±3.23	31.58±4.39	38.79±4.73	29.06±2.60	54.10±7.59	59.19±8.87	70.39 ± 9.58	59.74±1.87	65.17±1.94	79.53±1.69	9.86
ACMGNN	kNN	$\{G, G'\}$	$\theta_1 \neq \theta_2$	87.59±0.88	73.21±0.63	49.06±0.53	81.34±0.85	73.95±1.35	39.18±2.18	41.70±3.71	34.67±1.11	68.48±5.78	68.92±5.87	80.20±3.13	80.46±2.26	73.14±1.31	88.87±0.51	4.07
MLP	-	-	-	79.55±1.23	65.45±0.99	46.65±0.83	75.94±1.38	74.92±1.39	39.29±2.22	43.57±4.18	35.40±1.38	80.46±6.44	73.78±7.34	85.88±7.78	87.97±1.80	76.68±2.10	87.39±2.18	2.29
MixHop	-	-		90.10±5.59	81.70±0.89	50.95±0.71	84.56±1.19	77.66±1.24	41.22±2.66	43.11±4.73	33.59±1.23	72.54±8.98	62.43±9.54	75.88±8.27	87.76±1.94	76.51±1.93	89.42±0.81	1.86
MixHop	cos-graph	$\{G'\}$		64.75±4.59	51.83±0.53	41.47±2.00	68.78±1.94	71.45±1.38	37.75±2.41	37.79±2.10	31.77±1.75	55.72±6.39	60.27±5.85	70.20 ± 4.60	84.42±1.35	74.20 ± 0.83	88.74±0.29	8.21
MixHop	cos-graph	$\{\mathcal{G}, \mathcal{G}'\}$	$\theta_1 = \theta_2$	54.22±10.75	63.50±0.86	44.21±1.36	74.22±2.21	70.64±1.32	37.16±1.34	39.06±3.08		58.16±9.18	66.22±5.59	73.73±7.80	65.14±2.62	68.66±1.24	86.63±0.51	7.54
MixHop	cos-graph	$\{G, G'\}$	$\theta_1 \neq \theta_2$	84.71±1.19	55.41±1.63	43.37±0.75	74.41±1.33	69.63±2.03	37.64±2.19			61.13±7.96		75.29 ± 6.00	85.42±1.21	74.57±1.34	88.16±0.46	6.50
MixHop	cos-node	$\{G'\}$	-	60.56±7.08	51.74±0.68	42.71±0.97	74.27 ± 1.84	72.83 ± 1.12	38.35±1.99	38.88±3.00	33.05±1.04	58.42±6.52	60.27±5.98	71.57±4.91	83.22±1.16	74.11±1.12	88.23±0.45	6.71
MixHop	cos-node	$\{\mathcal{G}, \mathcal{G}'\}$	$\theta_1 = \theta_2$	54.22±10.75	63.50±0.86	44.21±1.36	74.22±2.21	70.64±1.32	37.16±1.34	39.06±3.08	32.24±1.33	58.16±9.18	66.22±5.59	73.73±7.80	65.14±2.62	68.66±1.24	86.63±0.51	7.64
MixHop	cos-node	$\{G, G'\}$	$\theta_1 \neq \theta_2$	85.43±0.57	55.95±2.35	44.15±0.59	76.54±0.91	72.03±2.45	37.47±2.07	39.52±3.33	32.50±1.10	60.61±8.73	62.97±6.75	75.10±6.20	85.36±0.89	74.68±1.13	88.18±0.52	4.79
MixHop	kNN	$\{G'\}$		59.50±6.26	50.39±0.72	42.07±0.93	70.49±1.70	69.57±1.32	38.07±1.72	38.76±2.91	33.23±1.30	59.25±4.49	57.30±6.96	69.22±7.22	83.99±1.28	74.96±1.18	87.99±0.40	8.00
MixHop	kNN	$\{G, G'\}$	$\theta_1 = \theta_2$	54.22±10.75	63.50±0.86	44.21±1.36	74.22±2.21	70.64±1.32	37.16±1.34	39.06±3.08	32.24±1.33	58.16±9.18	66.22±5.59	73.73±7.80	65.14±2.62	68.66±1.24	86.63±0.51	7.54
MixHop	kNN	$\{G, G'\}$	$\theta_1 \neq \theta_2$	85.53±0.50	57.48±1.98	43.28 ± 0.68	77.24±1.61	70.34±1.76	38.15±2.01	40.12±3.76	32.30±1.53	60.05±9.45	63.51±7.56	74.90±8.21	85.18±1.26	74.59±1.19	88.20±0.57	4.93

performance remains comparable to that of the GSL-enhanced counterparts. This suggests that GSL does not consistently enhance GNN performance in graph classification. Due to page limitations, we only tested a few methods in this paper. We believe it would be valuable to explore additional state-of-the-art methods, datasets, and theoretical justifications for the effectiveness of GSL in graph classification in future work.

F.8 Robustness of GSL

We investigate the robustness of GSL with GNNs using 3 types of graph perturbation strategies:

- Additive Feature Noise: We randomly inject noise into node features, where the noise follows a normal distribution $N(0, \sigma^2)$. The level of noise is controlled by σ , taking values from the set [0, 0.01, 0.02, 0.05, 0.1, 0.2, 0.5, 1.0, 2.0, 3.0, 10.0].
- Edge Addition: We randomly add edges to the graph structure, with the ratio of added edges proportional to the original number of edges, ranging from $r \in [0.0, 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9]$.
- Edge Removal: We randomly remove edges from the graph structure, with the ratio of removed edges also proportional to the original number of edges, following the same range as in edge addition.

We then measure model performance through accuracy or AUC-ROC in node classification.

Figure 12 illustrates the differences in model performance between GNN baselines and their GSL-enhanced counterparts across additional datasets beyond those shown in Figure 4. Generally, the performance of GSL is comparable to or even worse than that of the GNN baselines for all three types of perturbed graphs. Notably, model performance is not consistently stable for structural

Table 7: Performance of GNNs with their counterparts of trainable GSL.

Model	GSL Type	Mines.	Roman.	Amazon.	Tolokers	Questions	Cora	CiteSeer	PubMed	Rank
GCN	No GSL Trainable GSL Non-trainable GSL	90.07±5.79 90.07±0.58 89.17±0.68	81.46±1.25 78.76±0.46 72.63±1.45	50.89±1.16 50.89±0.65 48.31±0.96	84.61 ± 0.99 84.61 ± 0.65 82.91±0.97	77.68±1.10 OOM 75.56±1.05	87.97±1.51 84.92±1.51 85.69±1.73	76.75 ± 2.30 74.89±1.13 75.49±1.42	89.47 ± 0.64 88.66±0.45 88.72±0.71	1.19 2.31 2.50
SGC	No GSL Trainable GSL Non-trainable GSL	83.45±4.47 83.45±1.03 79.03±3.76	78.04±0.69 74.74±0.57 67.84±1.87	51.38±0.68 51.38±0.57 47.93±0.94	84.88±1.13 84.88±0.65 78.09±1.84	77.39±1.23 OOM 75.46±1.43	88.10±1.89 86.99±1.64 87.47±1.86	77.52±2.20 75.13±1.26 76.36±1.27	89.39 ± 0.62 88.94±0.31 89.37±0.41	1.19 2.31 2.50
GraphSAGE	No GSL Trainable GSL Non-trainable GSL	$\frac{90.66 \pm 0.88}{90.66 \pm 0.58}$ 90.67 ± 0.66	85.02±0.97 82.54±0.60 79.02±1.21	52.93±0.83 52.93±0.59 52.10±0.84	83.31±1.12 83.31±0.50 82.17±0.89	75.95±1.41 OOM 75.38±0.96	88.13±1.77 83.48±1.69 83.60±1.78	76.65 ± 2.00 74.18±1.02 74.39±1.35	89.18 ± 0.65 88.67±0.39 88.88±0.50	1.31 2.44 2.25
GAT	No GSL Trainable GSL Non-trainable GSL	90.41±1.34 90.41±0.61 89.96±0.79	84.51±0.84 83.10±0.58 77.23±1.63	52.00±2.84 52.10±0.62 49.79±0.72	84.37±0.96 84.35±0.56 82.78±0.95	77.78±1.27 OOM 76.67±1.13	88.02±1.92 86.23±1.58 86.97±1.75	76.77 ± 2.02 74.39±1.14 75.20±1.55	89.21±0.67 88.13±0.56 87.97±0.51	1.19 2.19 2.62
ACMGNN	No GSL Trainable GSL Non-trainable GSL	90.56±1.03 90.56±0.63 87.46±1.02	84.86±0.73 81.90±0.71 74.63±0.76	52.07±1.72 51.87±0.44 49.35±0.58	84.41±1.12 84.40±0.79 81.63±0.87	77.72±1.59 OOM 73.84±1.41	$\begin{array}{c} \textbf{88.23} {\pm} \textbf{1.81} \\ \underline{81.16} {\pm} 1.81 \\ 80.83 {\pm} 1.84 \end{array}$	76.63±2.34 73.91±1.16 73.43±1.47	89.37 ± 0.56 88.55±0.39 88.98±0.47	1.06 2.19 2.75
МіхНор	No GSL Trainable GSL Non-trainable GSL	90.10±5.59 90.10±0.52 85.43±0.57	81.70±0.89 79.07±0.75 55.95±2.35	50.95±0.71 50.95±0.71 44.15±0.59	84.56±1.19 84.55±0.67 76.54±0.91	77.66±1.24 OOM 72.03±2.45	87.76 ± 1.94 84.84±1.28 85.36±0.89	76.51 ± 1.93 74.45±1.11 74.68±1.13	89.42±0.81 88.48±0.62 88.18±0.52	1.12 2.25 2.62

Table 8: Ablation study of GSL-enhanced methods for graph classification.

Model	Cora		PubMed	l	CiteSeer		
1,10,001	AUC	Time	AUC	Time	Acc	Time	
ProGNN	76.28±0.52	959s	OOM	-	67.14±0.23	1776s	
ProGNN,w/o. GSL	78.96 ± 0.64	30s	75.80 ± 0.95	326s	67.24 ± 1.48	44s	
GEN	79.88 ± 0.93	219s	OOM	-	66.98 ± 1.28	320s	
GEN,w/o. GSL	78.32 ± 1.21	3s	76.94 ± 0.40	47s	64.66 ± 1.46	3s	
GRCN	83.04 ± 0.33	56s	74.55 ± 0.96	249s	70.85 ± 0.87	113s	
GRCN,w/o. GSL	71.82 ± 0.61	9s	74.18 ± 0.63	28s	58.33 ± 0.17	24s	
IDGL	83.32 ± 0.59	144s	OOM	-	70.57 ± 0.26	330s	
IDGL,w/o. GSL	83.32 ± 0.59	129s	OOM	-	71.12 ± 0.31	401s	

perturbations in heterophilous graphs. We attribute this inconsistency to the non-informative nature of the structural information in these graphs, which leads to diminished responses to edge addition or removal. Despite this, GSL still fails to consistently outperform GNN baselines.

1310

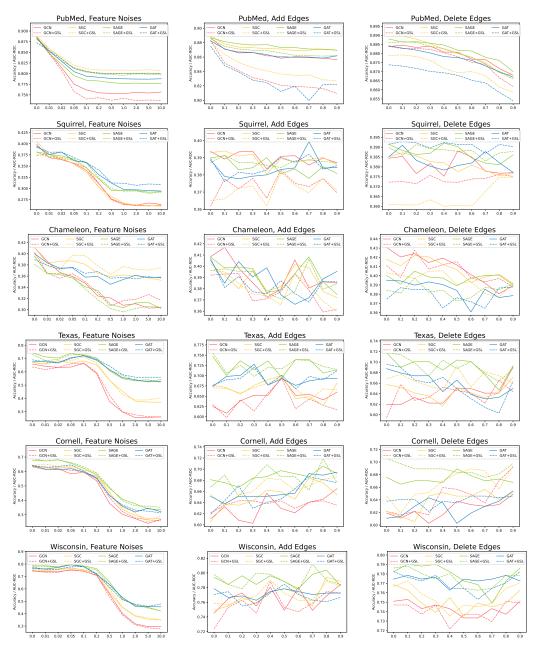


Figure 12: Response to feature noise, edge additions, and edge removals in GNN baselines and their GSL-enhanced counterparts.