

# GRAPH INFERENCE ACCELERATION BY BRIDGING GNNs AND MLPs WITH SELF-SUPERVISED LEARNING

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Graph Neural Networks (GNNs) have demonstrated their effectiveness in a variety of graph learning tasks such as node classification and link prediction. However, GNN inference mainly relies on neighborhood aggregation, which limits the deployment in latency-sensitive (i.e., real-time) applications such as financial fraud detection. To solve this problem, recent works have proposed to distill knowledge from teacher GNNs to student Multi-Layer Perceptrons (MLPs) trained on node content for inference acceleration. Despite the progress, these studies still suffer insufficient exploration of structural information when inferring unseen nodes. To address this issue, we propose a new method (namely **SSL-GM**) to fully integrate rich structural information into MLPs by bridging GNNs and MLPs with Self-Supervised Learning (SSL) for graph inference acceleration while improving model generalization capability. A key new insight of SSL-GM is that, without fetching their neighborhoods, the structural information of unseen nodes can be inferred solely from the nodes themselves with SSL. Specifically, SSL-GM employs self-supervised contrastive learning to align the representations encoded by graph context-aware GNNs and neighborhood dependency-free MLPs, fully integrating the structural information into MLPs. In particular, SSL-GM approximates the representations of GNNs using a non-parametric aggregator to avoid potential model collapse and exploits augmentation to facilitate the training; additionally, SSL-GM further incorporates reconstruction regulation to prevent representation shift caused by augmentation. Theoretically, we interpret our proposed SSL-GM through the principle of information bottleneck, demonstrating its generalization capability; we also analyze model capacity in incorporating structural information from the perspective of mutual information maximization and graph smoothness. Empirically, we demonstrate the superiority of SSL-GM over existing state-of-the-art models in both efficiency and effectiveness. In particular, SSL-GM obtains significant performance gains (7~26%) in comparison to MLPs, and a remarkable acceleration of GNNs (90~126×) on large-scale graph datasets.

## 1 INTRODUCTION

Due to the ubiquity of graph-structured data (e.g., computing networks, recommender systems in e-commerce, citation networks, and social networks), Graph Neural Networks (GNNs) have drawn significant attention in recent years. Generally, GNNs are based on neighborhood aggregation (Gilmer et al., 2017) to learn representations of given graphs. Despite their effectiveness in various graph learning tasks such as node classification and link prediction, due to the cost of neighborhood fetching during the testing stage, GNNs still face limitations of the deployment in latency-sensitive (i.e., real-time) applications such as financial fraud detection (Zhang et al., 2022; Wang et al., 2021). To solve this problem, existing works mainly adopt quantization (Ding et al., 2021), pruning (Zhou et al., 2021), and knowledge distillation (Yan et al., 2020) for graph inference acceleration. However, these improvements are limited as they still rely on neighborhood dependency (Zhang et al., 2022). To address this issue, as Multi-Layer Perceptrons (MLPs) have no dependency on graph data and can be efficiently deployed in latency-sensitive applications, researchers have explored distilling knowledge from pre-trained GNNs into MLPs (Zhang et al., 2022; Tian et al., 2023; Wang et al., 2023). Despite the progress, these methods inevitably sacrifice the model generalization capacity, as they cannot fully leverage structural information when inferring testing nodes. Given these challenges, we naturally ask: *how to bridge graph context-aware GNNs and neighborhood dependency-free MLPs for graph inference acceleration while improving model generalization capability?*

To answer the above question, we bring a key new insight different from existing works on graph inference acceleration: without fetching their neighborhoods, the structural information of unseen nodes can be inferred solely from nodes themselves with Self-Supervised Learning (SSL) (Chen et al., 2020b). Accordingly, we propose a new method (namely **SSL-GM**) to fully integrate rich structural information into MLPs by bridging GNNs and MLPs with Self-Supervised Learning (SSL) for graph inference acceleration while improving model generalization capability (Huang et al., 2023; Cabannes et al., 2023). More specifically, our proposed SSL-GM applies self-supervised contrastive learning (He et al., 2020) to align the consistency between GNNs and MLPs in the representation space. In particular, SSL-GM approximates the representations of GNNs using a non-parametric aggregator to avoid potential model collapse (Grill et al., 2020) and exploits augmentation to further enhance generalization (Zhao et al., 2021); additionally, SSL-GM further incorporates reconstruction regulation to prevent representation shift caused by augmentation. Theoretically, we have demonstrated that minimizing the objective function of SSL-GM is equivalent to optimizing the information bottleneck, thereby assuring model generalization (Alemi et al., 2017); in addition, we also analyze model capacity in incorporating structural information from the perspective of mutual information maximization and graph smoothness. Empirically, through extensive experiments over large-scale graph datasets, SSL-GM shows state-of-the-art performance on node classification tasks across transductive, inductive (Zhang et al., 2022), and cold-start settings. In terms of inference efficiency, SSL-GM exhibits remarkable acceleration compared to GNNs ( $90\sim 126\times$ ) and other acceleration techniques ( $5\sim 90\times$ ). The main contributions of our work are summarized below:

**(Methodology)** We observe that existing graph inference acceleration methods using low-latency MLPs are incapable of acquiring generalizable structure-aware representations. To this end, we propose SSL-GM to bridge graph context-aware GNNs and neighborhood dependency-free MLPs with SSL *at the first attempt* for graph inference acceleration while improving model generalization.

**(Theory)** We establish the theoretical equivalence between our objective and information bottleneck, proving the generalization capability of SSL-GM. Moreover, we analyze the capability of SSL-GM on encoding structural knowledge through mutual information maximization and graph smoothness.

**(Experiments)** Our SSL-GM achieves state-of-the-art performance over ten graph benchmarks on node classification tasks under transductive, inductive, and cold-start settings. It also shows a remarkable graph inference acceleration compared to GNNs ( $90\sim 126\times$ ) and exhibits significant performance improvements over vanilla MLPs ( $7\sim 26\%$ ).

## 2 RELATED WORK

**Graph Neural Networks** learn node representations by passing and aggregating messages from neighboring nodes. For example, GCN (Kipf & Welling, 2017) employs the normalized Laplacian matrix to guide message passing, GraphSAGE (Hamilton et al., 2017) utilizes neighborhood sampling, and GAT (Veličković et al., 2018) applies attention mechanisms. More recently, certain studies (Wu et al., 2019; Gasteiger et al., 2019; Han et al., 2023) decompose feature transformation from message passing, demonstrating that the effectiveness of GNNs stems from message propagation (Yang et al., 2023a). Despite their success, the inference speed on testing nodes remains a limitation due to neighborhood dependencies, which will be addressed in this paper.

**Self-Supervised Learning** (SSL) (Chen et al., 2020b; He et al., 2020) is a pre-training strategy that cultivates discriminative representations without supervision. Numerous studies (Veličković et al., 2019; Hassani & Khasahmadi, 2020; Zhu et al., 2020; 2021; Thakoor et al., 2022) have introduced methods for acquiring knowledge from graphs, prompting the downstream tasks (Sun et al., 2023). Particularly, BGRL (Thakoor et al., 2022) employs bootstrapping (Grill et al., 2020) to minimize the distance between node representations in two augmented views with an efficient approach. Despite potential improvements in generalization (Jiang et al., 2019; Huang et al., 2023), the dependency on neighborhood information continues to constrain inference speed.

**Inference Acceleration** encompasses quantization (Gupta et al., 2015; Jacob et al., 2018), pruning (Han et al., 2015; Frankle & Carbin, 2019), and knowledge distillation (KD) (Hinton et al., 2015). Quantization (Ding et al., 2021) approximates continuous data with limited discrete values, pruning (Zhou et al., 2021) involves removing connections within neural networks, and KD distills knowledge from large GNNs to small GNNs (Yan et al., 2020). However, these methods still fail to eliminate neighborhood dependencies, resulting in constrained inference acceleration. In light of this, GLNN (Zhang et al., 2022) distills knowledge from teacher GNNs to student MLPs, bypassing

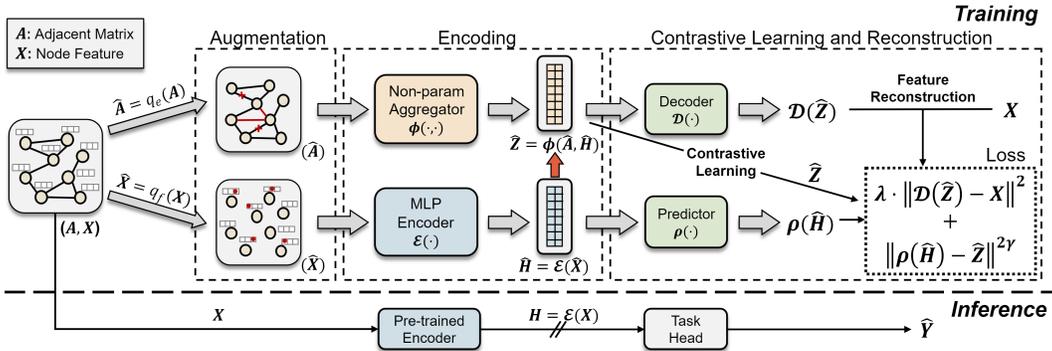


Figure 1: In training, SSL-GM augments the original graph and encodes representations through MLPs and GNNs, where the GNN representations are approximated by a non-parametric aggregator. Following, SSL-GM employs contrastive learning to maximize the alignment between these two representations and applies reconstruction based on GNN representations. In inference, only the MLP is leveraged to encode graph representations, leading to a substantial inference acceleration.

neighbor-fetching latency, but cannot fully model the structural information. Following studies integrate structural knowledge in pre-processing, e.g., appending positional embeddings (Perozzi et al., 2014) to node content (Tian et al., 2023; Wang et al., 2023), conducting label propagation (Yang et al., 2021), acquiring a motif cookbook (Yang et al., 2023b). However, these methods introduce additional time overhead and are impractical for out-of-distribution and cold-start generalization. In parallel, GraphMLP (Hu et al., 2021) and following works (Dong et al., 2022; Liu et al., 2022) employ neighborhood-aware distillation to train MLPs, albeit limited to the transductive setting. Unlike these methods, SSL-GM applies self-supervised learning to bridge GNNs and MLPs, achieving significant inference acceleration meanwhile enhancing model generalization.

### 3 BRIDGING GNNs AND MLPs WITH SELF-SUPERVISED LEARNING

In this section, we provide a detailed description of SSL-GM. Figure 1 illustrates an overview of the SSL-GM framework. The key insight is that, without fetching node neighborhoods, the potential structural distribution of unseen nodes can be inferred solely with SSL. To achieve the goal, SSL-GM employs self-supervised contrastive learning to align the consistency between the representations encoded by GNNs and MLPs, integrating fine-grained structural information into MLPs to learn generalizable graph representations (Hendrycks et al., 2019; Cabannes et al., 2023).

**Problem Statement.** We consider a graph  $\mathcal{G} = (\mathbf{A}, \mathbf{X})$  consisting of node set  $\mathcal{V}$  and edge set  $E$ , with  $N$  nodes in total. We have node features  $\mathbf{X} \in \mathbb{R}^{N \times d}$  with dimension  $d$ , and adjacent matrix  $\mathbf{A}$ , where  $A_{ij} = 1$  if node  $i$  and  $j$  are connected, and  $A_{ij} = 0$  otherwise. Our focus lies in training an MLP encoder  $\mathcal{E}(\cdot)$  without supervision, which receives node content  $\mathbf{X}$  as input and generates node representations  $\mathbf{H} \in \mathbb{R}^{N \times d'}$  preserving the semantics of both  $\mathbf{A}$  and  $\mathbf{X}$ .

#### 3.1 STRUCTURE-AWARE MLPs WITH SELF-SUPERVISED LEARNING

The effectiveness of GNNs stems from their ability to learn graph contextual information. Although some methods propose to distill knowledge from GNNs to MLPs, they inevitably sacrifice model capability (Tian et al., 2020) and generalization (Tian et al., 2023), as they only align the predictions of MLPs and GNNs in the label space, falling short to fully explore graph structural knowledge. To this end, we employ self-supervised contrastive learning to comprehensively incorporate structural information into MLPs by aligning GNNs and MLPs in the representation space. Specifically, we treat the representations encoded by MLPs as  $\mathbf{H}$  and GNNs as  $\mathbf{Z}$ , which corresponds to the encodings on node view  $\mathcal{G}_1 = (\emptyset, \mathbf{X})$  and graph view  $\mathcal{G}_2 = (\mathbf{A}, \mathbf{X})$ , respectively. The objective is to optimize the consistency between these two representations, thereby encoding structural knowledge into MLPs. We employ the Bootstrap loss (Grill et al., 2020) as objective function, defined as

$$\mathcal{L}_{cont} = \mathbb{E} \|\rho(\mathbf{H}) - \mathbf{Z}\|^{2\gamma}, \quad (1)$$

where  $\gamma \geq 1$  serves as a scaling term, akin to an adaptive sample reweighing technique (Hou et al., 2022; Lin et al., 2017). The projector  $\rho(\cdot)$  can either be identity or learnable. Here we opt for a non-linear MLP to enhance the expressiveness in estimating instance distances (Chen et al., 2020b).

**Non-Parametric Aggregator for Approximating GNN Representations.** Directly applying GNNs may lead to model collapse, as evidenced in Appendix D.3. We suppose it derives from the inconsistency between representations learned by MLPs and GNNs (He et al., 2020; Grill et al., 2020). To mitigate the issue, we propose to propagate the representations learned by MLPs to approximate the representations of GNNs, preserving their inherent consistency. Specifically, we employ a non-parametric aggregator  $\phi(\cdot, \cdot)$  to conduct message passing as follows

$$\text{GNN} : \mathbf{Z} = \text{GNN}(\mathbf{A}, \mathbf{X}; \Theta) \implies \text{SSL-GM} : \mathbf{Z} = \phi(\mathbf{A}, \mathbf{H}), \mathbf{H} = \mathcal{E}(\mathbf{X}; \Theta), \quad (2)$$

where  $\Theta$  represents the parameters of the model. Unlike existing GNNs, our SSL-GM decomposes the linear transformation and message passing by applying MLP encoder to transform node features  $\mathbf{H} = \mathcal{E}(\mathbf{X})$  and then employing the non-parametric aggregator to perform message passing  $\mathbf{Z} = \phi(\mathbf{A}, \mathbf{H})$  without additional transformation. Further details are presented in Appendix D.6. This approach has been empirically (Wu et al., 2019; Gasteiger et al., 2019) and theoretically (Han et al., 2023; Yang et al., 2023a) shown to be expressive. The choice of aggregation type can be arbitrary. We employ a GCN-like (Kipf & Welling, 2017) framework for neighborhood aggregation.

### 3.2 AUGMENTATION TO FACILITATE TRAINING

The fundamental assumption behind training MLPs on graphs is that nodes with similar features have similar surrounding ego-graphs (Chen et al., 2021). This aligns with our insight that the contextual neighborhoods can be inferred based on the target nodes themselves. However, the assumption indicates the necessity of high-quality node features (Zhang et al., 2022; Guo et al., 2023) and inherently requires the consistency in structural distribution between training and testing graphs (Luan et al., 2022; 2023). This prevents MLP-based graph learning algorithms from generalizing to out-of-distribution settings. To solve the issue, we augment the node and ego-graph pairs to increase the quality and diversity of training data (Feng et al., 2020; You et al., 2020; Zhao et al., 2021). This approach will enhance model generalization and robustness for graphs that originate from distributions different from the training set. The augmentation is applied in each training epoch as

$$\hat{\mathcal{G}} = (\hat{\mathbf{A}}, \hat{\mathbf{X}}), \hat{\mathbf{A}} \sim q_e(\mathbf{A}), \hat{\mathbf{X}} \sim q_f(\mathbf{X}), \quad (3)$$

where  $q_e(\cdot)$  and  $q_f(\cdot)$  are two random augmentation methods for graph structures and node features, respectively. This augmentation aligns with the objective of optimal contrastive learning (Xu et al., 2021) that aims to train an augmentation-invariant encoder, formulated as

$$\mathcal{E}^* = \arg \min_{\mathcal{E}} I(\mathcal{G}_1, \mathcal{G}_2) - I(\hat{\mathbf{H}}, \hat{\mathbf{Z}}), \quad (4)$$

where  $\mathcal{G}_1 = (\emptyset, \hat{\mathbf{X}})$  and  $\mathcal{G}_2 = (\hat{\mathbf{A}}, \hat{\mathbf{X}})$ , with  $\hat{\mathbf{H}}$  and  $\hat{\mathbf{Z}}$  are representations encoded by MLPs and GNNs on augmented graphs, respectively. This process will facilitate the training of the MLP encoder  $\mathcal{E}$  by incorporating more structure-relevant information into SSL-GM (Xu et al., 2021).

### 3.3 RECONSTRUCTION FOR MITIGATING REPRESENTATION SHIFT

While augmentation can facilitate the training process, it may impact the distribution of encoded representations. **This representation shift, particularly pronounced in structural augmentation, can significantly alter the local structure of the target node. For instance, as depicted in Figure 2, simple edge permutation can dramatically change the 2-hop neighborhoods of a target node, leading to a substantial shift in the representations.** Although this shift can potentially benefit existing graph contrastive learning methods by providing adversarial samples (Suresh et al., 2021; Kong et al., 2022; Feng et al., 2022), it may result in a mismatch between the augmented node and ego-graph pairs, thereby impairing the quality of representations learned by MLPs. **To counter this problem, we hypothesize that if GNN representations can preserve localized information, the impact of representation shift can be minimized. Based on this, we introduce a reconstruction regularizer that reconstructs the raw node features  $\mathbf{X}$  based on the GNN representations  $\hat{\mathbf{Z}}$  on augmented graphs  $\hat{\mathcal{G}} = (\hat{\mathbf{A}}, \hat{\mathbf{X}})$ .** The reconstruction term is defined as

$$\mathcal{L}_{rec} = \mathbb{E} \|\mathcal{D}(\hat{\mathbf{Z}}) - \mathbf{X}\|^2, \hat{\mathbf{Z}} = \phi(\mathbf{A}, \hat{\mathbf{H}}), \hat{\mathbf{H}} = \mathcal{E}(\hat{\mathbf{X}}), \quad (5)$$

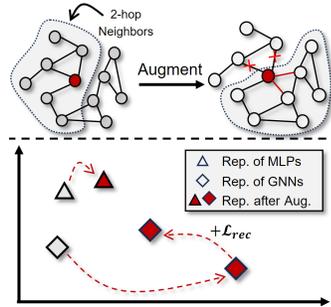


Figure 2: Augmentation leads to severe representation shift.

where  $\mathcal{D}(\cdot)$  represents the decoder, while  $\hat{\mathbf{Z}}$  and  $\hat{\mathbf{H}}$  correspond to representations of GNNs and MLPs on the augmented graphs, respectively. This term encourages GNNs to preserve more localized information, ensuring GNN representations do not significantly shift from MLP representations, shown in Figure 2. However, when no augmentations are applied, the term may even lead to model performance degradation, evaluated in Appendix D.5. Additionally, the reconstruction regularizer also provides denoising capability (Batson & Royer, 2019) to some extent, enhancing the robustness of SSL-GM to noisy data.

### 3.4 OBJECTIVE FUNCTION

Considering the above three modules, we define the overall objective of SSL-GM as

$$\mathcal{L} = \mathcal{L}_{cont} + \lambda \cdot \mathcal{L}_{rec} = \mathbb{E}[\underbrace{\|\rho(\hat{\mathbf{H}}) - \hat{\mathbf{Z}}\|^{2\gamma}}_{invariance} + \lambda \cdot \underbrace{\|\mathcal{D}(\hat{\mathbf{Z}}) - \mathbf{X}\|^2}_{reconstruct}], \quad (6)$$

where  $\hat{\mathbf{H}} = \mathcal{E}(\hat{\mathbf{X}})$  and  $\hat{\mathbf{Z}} = \phi(\hat{\mathbf{A}}, \hat{\mathbf{H}})$  are representations of MLPs and GNNs on the augmented graph  $\hat{\mathcal{G}} = (\hat{\mathbf{A}}, \hat{\mathbf{X}})$ .  $\rho$  and  $\mathcal{D}$  indicate projector head and decoder.  $\gamma$  and  $\lambda$  denote scaling factor and trade-off coefficient. We dubbed the contrastive loss as invariance term and the reconstruction loss as reconstruction term. While seemingly distinct, these two terms collaboratively reciprocate each other. The invariance term produces graph context-aware MLP representations, akin to positional encoding (Rampášek et al., 2022), improving the quality of GNN representations (Dwivedi & Bresson, 2020; Kreuzer et al., 2021; Dwivedi et al., 2022; Wang et al., 2022). This enhanced discrimination capability facilitates node feature reconstruction. The reconstruction term, in turn, prevents potential representation shifts by preserving additional localized information in representations learned by GNNs, thus providing better high-order structural signals for training MLPs.

**Theorem 1** Suppose  $\mathcal{G} = (\mathbf{A}, \mathbf{X})$  is sampled from a latent graph  $\mathcal{G}_{\mathcal{I}} = (\mathbf{A}, \mathbf{F})$ ,  $\mathcal{G} \sim P(\mathcal{G}_{\mathcal{I}})$  (Xie et al., 2022), and  $\mathbf{F}^*$  is the lossless compression of  $\mathbf{F}$  that  $\mathbb{E}[\mathbf{X}|\mathbf{A}, \mathbf{F}^*] = \mathbf{F}$ . Let  $\mathcal{E}$  be a  $l$ -Lipschitz continuous function respect to  $l_2$ -norm,  $\rho$  be an identity projector, and  $\lambda = 1, \gamma = 1$ . Optimizing Eq. 6 equals to finding the optimal compression  $\mathbf{T}^*$  with minimal sufficient information  $\mathbf{C}$  where

$$\mathbf{T}^* = \arg \min_{\mathbf{T}} I(\mathcal{G}; \mathbf{T}) - \beta I(\mathbf{T}; \mathcal{G}_{\mathcal{I}}), \text{ s.t.}, I(\mathbf{T}, \mathcal{G}_{\mathcal{I}}) \geq \mathbf{C}, \mathbf{T} = (\mathbf{H}, \mathbf{Z}). \quad (7)$$

Theorem 1 reveals the equivalence of optimization objectives between our loss function in Eq. 6 and information bottleneck (Tishby et al., 2000; Tishby & Zaslavsky, 2015), ensuring our SSL-GM to learn informative and generalizable representations (Alemi et al., 2017) for downstream tasks.

## 4 EXPERIMENTS

### 4.1 EXPERIMENTAL SETTINGS

**Datasets.** We evaluate our SSL-GM model on ten public benchmark datasets for node classification task, including Cora, Citeseer, Pubmed, Amazon-CS, Amazon-Photo, Coauthor-CS, Coauthor-Physics, Wiki-CS, Flickr, and Arxiv. Dataset details are presented in Appendix B.1. We also evaluate the model performance on graph classification datasets, shown in Appendix B.1 and D.1.

**Baselines.** We compare SSL-GM to a variety of baselines, including supervised GNNs, e.g., SAGE (Hamilton et al., 2017), GAT (Veličković et al., 2018), SGC (Wu et al., 2019), APPNP (Gasteiger et al., 2019), self-supervised methods, including DGI (Veličković et al., 2019), MVGRL (Hassani & Khasahmadi, 2020), GRACE (Zhu et al., 2020), BGRL (Thakoor et al., 2022), GCA (Zhu et al., 2021), and MLP-based methods, such as basic MLP, GraphMLP (Hu et al., 2021), GLNN (Zhang et al., 2022), GENN (Wang et al., 2023), and NOSMOG (Tian et al., 2023). Appendix B.2 presents the details.

**Evaluation Protocol.** We report the mean and standard deviation of ten runs using different random seeds to assess model performance in node classification. Accuracy serves as the metric. We conduct comparison for SSL-GM against baseline methods in transductive, inductive (production) (Zhang et al., 2022), and cold-start settings. Appendix B.3 and B.4 show additional information.

### 4.2 NODE CLASSIFICATION RESULTS

**Transductive Setting.** Table 1 demonstrates the node classification results of SSL-GM and baselines in transductive setting. Our SSL-GM consistently outperforms MLP-based methods across all datasets, which highlights the expressiveness of incorporating fine-grained structural knowledge

		Cora	Citeseer	PubMed	Amazon-CS	Amazon-Photo	Co-CS	Co-Phys	Wiki-CS	Flickr	Arxiv
GNN	SAGE	81.44 $\pm$ 0.91	70.44 $\pm$ 1.39	85.94 $\pm$ 0.43	88.88 $\pm$ 0.27	93.81 $\pm$ 0.42	93.41 $\pm$ 0.16	95.72 $\pm$ 0.05	80.87 $\pm$ 0.63	48.45 $\pm$ 0.78	72.05 $\pm$ 0.25
	GAT	82.33 $\pm$ 1.17	68.89 $\pm$ 1.47	84.73 $\pm$ 0.43	89.92 $\pm$ 0.48	91.94 $\pm$ 0.44	91.98 $\pm$ 0.32	95.08 $\pm$ 0.15	79.97 $\pm$ 0.56	51.38 $\pm$ 0.16	71.79 $\pm$ 0.40
	APPNP	75.48 $\pm$ 1.57	68.09 $\pm$ 1.20	84.60 $\pm$ 0.30	87.41 $\pm$ 0.27	93.37 $\pm$ 0.47	94.62 $\pm$ 0.21	95.44 $\pm$ 0.12	79.10 $\pm$ 0.31	47.53 $\pm$ 0.29	71.01 $\pm$ 0.24
	SGC	81.80 $\pm$ 0.88	68.96 $\pm$ 1.63	85.28 $\pm$ 0.32	89.31 $\pm$ 0.62	92.74 $\pm$ 0.43	94.02 $\pm$ 0.20	94.78 $\pm$ 0.20	81.06 $\pm$ 0.55	51.75 $\pm$ 0.24	69.95 $\pm$ 0.35
GCL	DGI	82.31 $\pm$ 0.60	71.81 $\pm$ 0.73	76.78 $\pm$ 0.70	79.98 $\pm$ 0.19	91.60 $\pm$ 0.21	92.22 $\pm$ 0.53	94.50 $\pm$ 0.04	76.42 $\pm$ 0.55	46.88 $\pm$ 0.13	70.13 $\pm$ 0.15
	MVGRL	83.89 $\pm$ 0.50	72.14 $\pm$ 1.25	86.33 $\pm$ 0.59	87.85 $\pm$ 0.31	91.88 $\pm$ 0.15	92.15 $\pm$ 0.07	95.30 $\pm$ 0.04	77.64 $\pm$ 0.09	49.32 $\pm$ 0.11	70.88 $\pm$ 0.10
	GRACE	80.50 $\pm$ 1.03	65.52 $\pm$ 2.06	84.64 $\pm$ 0.50	88.44 $\pm$ 0.33	92.83 $\pm$ 0.56	93.01 $\pm$ 0.30	95.43 $\pm$ 0.06	78.59 $\pm$ 0.47	49.33 $\pm$ 0.11	70.96 $\pm$ 0.13
	GCA	83.53 $\pm$ 0.49	71.33 $\pm$ 0.15	86.03 $\pm$ 0.37	87.42 $\pm$ 0.30	92.61 $\pm$ 0.21	93.06 $\pm$ 0.03	95.72 $\pm$ 0.03	78.35 $\pm$ 0.05	49.03 $\pm$ 0.07	70.90 $\pm$ 0.08
MLP	BGRL	81.30 $\pm$ 0.59	66.90 $\pm$ 0.58	84.92 $\pm$ 0.24	88.19 $\pm$ 0.21	92.54 $\pm$ 0.11	92.11 $\pm$ 0.12	95.21 $\pm$ 0.07	77.54 $\pm$ 0.79	49.67 $\pm$ 0.06	70.84 $\pm$ 0.12
	MLP	64.49 $\pm$ 1.90	64.01 $\pm$ 1.26	80.69 $\pm$ 0.28	80.79 $\pm$ 0.33	87.77 $\pm$ 0.49	91.65 $\pm$ 0.32	95.11 $\pm$ 0.12	75.16 $\pm$ 0.46	46.21 $\pm$ 0.07	56.44 $\pm$ 0.30
	GraphMLP	79.50 $\pm$ 0.81	72.10 $\pm$ 0.48	84.27 $\pm$ 0.23	84.01 $\pm$ 0.58	90.90 $\pm$ 1.03	90.36 $\pm$ 0.64	93.51 $\pm$ 0.15	76.39 $\pm$ 0.53	46.25 $\pm$ 0.21	63.36 $\pm$ 0.18
	GLNN	81.32 $\pm$ 1.15	71.15 $\pm$ 0.71	86.34 $\pm$ 0.46	87.47 $\pm$ 0.60	93.87 $\pm$ 0.31	94.16 $\pm$ 0.21	95.40 $\pm$ 0.07	80.66 $\pm$ 0.74	46.18 $\pm$ 0.19	64.03 $\pm$ 0.51
NOSMOG	GENN	82.13 $\pm$ 0.77	71.42 $\pm$ 1.31	86.28 $\pm$ 0.31	87.12 $\pm$ 0.55	93.64 $\pm$ 0.65	93.82 $\pm$ 0.29	95.45 $\pm$ 0.05	80.48 $\pm$ 0.74	46.35 $\pm$ 0.34	70.13 $\pm$ 0.60
	NOSMOG	82.27 $\pm$ 1.13	72.39 $\pm$ 1.27	86.18 $\pm$ 0.33	87.64 $\pm$ 1.14	93.94 $\pm$ 0.47	93.83 $\pm$ 0.23	95.74 $\pm$ 0.12	80.53 $\pm$ 0.77	46.69 $\pm$ 0.25	70.84 $\pm$ 0.44
SSL-GM	SSL-GM	84.60 $\pm$ 0.24	73.52 $\pm$ 0.53	86.99 $\pm$ 0.09	88.46 $\pm$ 0.16	94.28 $\pm$ 0.08	94.87 $\pm$ 0.07	96.17 $\pm$ 0.03	81.21 $\pm$ 0.13	49.85 $\pm$ 0.09	71.12 $\pm$ 0.10
	$\Delta_{BGRL}$	$\uparrow 4.06\%$	$\uparrow 9.90\%$	$\uparrow 2.44\%$	$\uparrow 0.31\%$	$\uparrow 1.88\%$	$\uparrow 3.00\%$	$\uparrow 1.01\%$	$\uparrow 4.73\%$	$\uparrow 0.36\%$	$\uparrow 0.40\%$
	$\Delta_{MLP}$	$\uparrow 31.18\%$	$\uparrow 14.86\%$	$\uparrow 7.81\%$	$\uparrow 9.49\%$	$\uparrow 7.42\%$	$\uparrow 3.51\%$	$\uparrow 1.11\%$	$\uparrow 8.05\%$	$\uparrow 7.88\%$	$\uparrow 26.01\%$
	$\Delta_{NOSMOG}$	$\uparrow 2.83\%$	$\uparrow 1.56\%$	$\uparrow 0.94\%$	$\uparrow 0.94\%$	$\uparrow 0.36\%$	$\uparrow 1.11\%$	$\uparrow 0.45\%$	$\uparrow 0.84\%$	$\uparrow 6.77\%$	$\uparrow 0.40\%$
w/o Aggr.	55.91 $\pm$ 0.66	57.36 $\pm$ 0.33	79.93 $\pm$ 0.32	72.76 $\pm$ 0.71	77.05 $\pm$ 0.18	91.19 $\pm$ 0.13	93.35 $\pm$ 0.12	73.87 $\pm$ 0.26	45.82 $\pm$ 0.07	54.83 $\pm$ 0.41	
w/o Pred.	81.78 $\pm$ 0.30	73.09 $\pm$ 0.25	85.33 $\pm$ 0.10	83.12 $\pm$ 0.25	91.25 $\pm$ 0.27	93.32 $\pm$ 0.08	94.98 $\pm$ 0.06	76.13 $\pm$ 0.22	48.31 $\pm$ 0.14	67.48 $\pm$ 0.44	
w/o Aug.	82.10 $\pm$ 0.45	71.83 $\pm$ 0.43	86.89 $\pm$ 0.13	87.12 $\pm$ 0.15	93.52 $\pm$ 0.20	93.10 $\pm$ 0.05	94.56 $\pm$ 0.06	80.98 $\pm$ 0.13	48.21 $\pm$ 0.10	70.58 $\pm$ 0.20	
w/o Rec.	84.37 $\pm$ 0.27	73.18 $\pm$ 0.24	86.86 $\pm$ 0.10	88.25 $\pm$ 0.07	94.15 $\pm$ 0.07	94.64 $\pm$ 0.06	96.01 $\pm$ 0.07	81.10 $\pm$ 0.13	49.60 $\pm$ 0.11	70.38 $\pm$ 0.22	

Table 1: Node classification accuracy (%) under transductive setting.  $\Delta_{BGRL}$ ,  $\Delta_{MLP}$ , and  $\Delta_{NOSMOG}$  represent the performance gap (%) between our methods and MLP, BGRL, and NOSMOG, where green indicates the improvement over 4%.

		Cora	Citeseer	PubMed	Amazon-CS	Amazon-Photo	Co-CS	Co-Phys	Wiki-CS	Flickr	Arxiv
SAGE	SAGE	77.51 $\pm$ 1.77	68.40 $\pm$ 1.61	85.04 $\pm$ 0.44	87.24 $\pm$ 0.43	93.20 $\pm$ 0.45	92.88 $\pm$ 0.40	95.74 $\pm$ 0.12	79.26 $\pm$ 0.65	47.17 $\pm$ 0.73	68.52 $\pm$ 0.56
	BGRL	77.73 $\pm$ 1.07	64.33 $\pm$ 1.56	83.97 $\pm$ 0.48	87.33 $\pm$ 0.48	91.47 $\pm$ 0.62	91.26 $\pm$ 0.35	94.38 $\pm$ 0.29	76.25 $\pm$ 1.09	49.12 $\pm$ 0.31	69.29 $\pm$ 0.38
MLP	MLP	63.76 $\pm$ 1.65	63.98 $\pm$ 1.22	80.91 $\pm$ 0.45	81.00 $\pm$ 0.54	87.73 $\pm$ 0.88	91.68 $\pm$ 0.59	95.18 $\pm$ 0.13	75.08 $\pm$ 0.71	46.14 $\pm$ 0.22	55.89 $\pm$ 0.51
	GLNN	78.34 $\pm$ 1.04	69.61 $\pm$ 1.13	85.44 $\pm$ 0.48	87.04 $\pm$ 0.50	93.28 $\pm$ 0.43	93.72 $\pm$ 0.35	95.76 $\pm$ 0.09	78.39 $\pm$ 0.54	46.11 $\pm$ 0.27	63.53 $\pm$ 0.48
	GENN	77.83 $\pm$ 1.57	67.30 $\pm$ 1.48	84.34 $\pm$ 0.47	85.75 $\pm$ 1.20	92.09 $\pm$ 0.96	93.57 $\pm$ 0.37	95.67 $\pm$ 0.06	78.27 $\pm$ 1.01	45.56 $\pm$ 0.51	68.52 $\pm$ 0.54
NOSMOG	NOSMOG	77.83 $\pm$ 1.94	68.58 $\pm$ 1.41	83.84 $\pm$ 0.45	86.61 $\pm$ 1.22	92.52 $\pm$ 0.68	93.45 $\pm$ 0.44	95.78 $\pm$ 0.10	78.35 $\pm$ 0.70	46.05 $\pm$ 0.55	69.10 $\pm$ 0.80
	SSL-GM	81.37 $\pm$ 1.20	72.33 $\pm$ 0.90	86.47 $\pm$ 0.28	87.65 $\pm$ 0.40	93.87 $\pm$ 0.32	94.63 $\pm$ 0.16	96.04 $\pm$ 0.12	79.26 $\pm$ 0.83	49.27 $\pm$ 0.18	70.23 $\pm$ 0.47
$\Delta_{BGRL}$	$\uparrow 4.68\%$	$\uparrow 12.44\%$	$\uparrow 2.98\%$	$\uparrow 0.37\%$	$\uparrow 2.62\%$	$\uparrow 3.69\%$	$\uparrow 1.76\%$	$\uparrow 3.95\%$	$\uparrow 0.31\%$	$\uparrow 1.36\%$	
$\Delta_{MLP}$	$\uparrow 27.62\%$	$\uparrow 13.05\%$	$\uparrow 6.87\%$	$\uparrow 8.21\%$	$\uparrow 7.00\%$	$\uparrow 3.22\%$	$\uparrow 0.90\%$	$\uparrow 5.57\%$	$\uparrow 6.78\%$	$\uparrow 25.66\%$	
$\Delta_{NOSMOG}$	$\uparrow 4.55\%$	$\uparrow 5.47\%$	$\uparrow 3.14\%$	$\uparrow 1.20\%$	$\uparrow 1.46\%$	$\uparrow 1.26\%$	$\uparrow 0.27\%$	$\uparrow 1.16\%$	$\uparrow 6.99\%$	$\uparrow 1.64\%$	

Table 2: Node classification accuracy (%) under inductive (production) settings.

		Cora	Citeseer	PubMed	Amazon-CS	Amazon-Photo	Co-CS	Co-Phys	Wiki-CS	Flickr	Arxiv
SAGE	SAGE	60.23 $\pm$ 5.03	56.62 $\pm$ 5.10	77.98 $\pm$ 1.53	61.01 $\pm$ 4.51	59.52 $\pm$ 8.02	91.30 $\pm$ 0.84	94.64 $\pm$ 0.88	52.73 $\pm$ 7.93	41.06 $\pm$ 2.25	43.47 $\pm$ 2.53
	BGRL	78.80 $\pm$ 1.14	65.10 $\pm$ 2.08	84.18 $\pm$ 0.80	86.13 $\pm$ 0.76	90.39 $\pm$ 0.30	90.23 $\pm$ 0.48	94.06 $\pm$ 0.28	78.15 $\pm$ 1.17	48.73 $\pm$ 0.13	64.11 $\pm$ 0.20
MLP	MLP	64.15 $\pm$ 2.11	64.43 $\pm$ 1.76	80.90 $\pm$ 0.72	80.80 $\pm$ 0.91	87.88 $\pm$ 0.96	91.78 $\pm$ 0.81	95.16 $\pm$ 0.18	74.94 $\pm$ 1.81	46.09 $\pm$ 0.50	55.91 $\pm$ 0.69
	GLNN	71.96 $\pm$ 1.68	69.14 $\pm$ 2.58	84.42 $\pm$ 0.87	83.98 $\pm$ 0.70	91.05 $\pm$ 0.49	93.34 $\pm$ 0.47	95.70 $\pm$ 0.09	77.64 $\pm$ 1.42	46.05 $\pm$ 0.43	60.55 $\pm$ 0.55
	GENN	69.06 $\pm$ 4.80	65.44 $\pm$ 2.33	78.19 $\pm$ 2.08	79.44 $\pm$ 1.66	90.18 $\pm$ 0.62	93.54 $\pm$ 0.55	95.55 $\pm$ 0.25	67.31 $\pm$ 1.66	45.24 $\pm$ 0.72	61.30 $\pm$ 0.59
NOSMOG	NOSMOG	70.69 $\pm$ 2.45	68.03 $\pm$ 2.79	81.48 $\pm$ 1.30	81.95 $\pm$ 1.04	91.15 $\pm$ 0.88	93.63 $\pm$ 0.42	95.54 $\pm$ 0.40	68.49 $\pm$ 3.61	46.07 $\pm$ 0.30	61.64 $\pm$ 0.93
	SSL-GM	80.48 $\pm$ 2.15	72.81 $\pm$ 1.61	86.44 $\pm$ 0.51	87.58 $\pm$ 0.99	93.91 $\pm$ 0.58	94.51 $\pm$ 0.15	95.97 $\pm$ 0.24	78.46 $\pm$ 1.48	49.41 $\pm$ 0.46	66.13 $\pm$ 1.05
$\Delta_{BGRL}$	$\uparrow 2.13\%$	$\uparrow 11.84\%$	$\uparrow 2.68\%$	$\uparrow 1.68\%$	$\uparrow 3.89\%$	$\uparrow 4.74\%$	$\uparrow 2.03\%$	$\uparrow 0.40\%$	$\uparrow 1.40\%$	$\uparrow 3.15\%$	
$\Delta_{MLP}$	$\uparrow 25.46\%$	$\uparrow 13.01\%$	$\uparrow 6.85\%$	$\uparrow 8.39\%$	$\uparrow 6.86\%$	$\uparrow 2.97\%$	$\uparrow 0.85\%$	$\uparrow 4.70\%$	$\uparrow 7.20\%$	$\uparrow 18.28\%$	
$\Delta_{NOSMOG}$	$\uparrow 13.85\%$	$\uparrow 7.03\%$	$\uparrow 6.09\%$	$\uparrow 6.87\%$	$\uparrow 3.03\%$	$\uparrow 0.94\%$	$\uparrow 0.45\%$	$\uparrow 14.56\%$	$\uparrow 7.25\%$	$\uparrow 7.28\%$	

Table 3: Node classification accuracy (%) under cold-start setting.

into MLPs. More specifically, SSL-GM surpasses MLP, GraphMLP, and GLNN on large-scale dataset Arxiv by 26%, 12%, and 11% improvements. Furthermore, SSL-GM outperforms other self-supervised methods in all datasets, and fully-supervised GNNs in 7 out of 10 datasets. In particular, SSL-GM achieves superior performance compared to SGC and APPNP, both of which employ propagators to encode node representations in a manner similar to our aggregator. To further analyze the expressiveness of SSL-GM, we conduct graph classification in Appendix D.1. Experimental results show our SSL-GM achieves the best or sub-best performance across 6 out of 7 datasets.

**Inductive (Production) Setting.** Table 2 presents the node classification results of SSL-GM and baseline methods in inductive (production) setting. We partition the original graph into two non-overlapping sets, namely the transductive set  $\mathcal{G}^T = (\mathcal{V}^T, E^T)$  and the inductive set  $\mathcal{G}^I = (\mathcal{V}^I, E^I)$ , each with distinct structural distributions. We evaluate transductive and inductive results on these sets and then interpolate them to derive the production results. Unlike Tian et al. (2023), which establish connections between nodes in  $\mathcal{V}^I$  and  $\mathcal{V}^T$  during inductive inference, our setting is more challenging as it treats these two sets as independent from each other, simulating an *out-of-distribution* scenario. Appendix B.4 includes more details. We report the production results in Table 2 and provide the comprehensive results in Appendix D.2. SSL-GM outperforms all baselines in all datasets, demonstrating the effectiveness of SSL-GM in real-world settings. Surprisingly, GLNN outper-

forms the advanced NOSMOG on 8 out of 10 datasets. We suppose that the reliance on additional positional embeddings hinders the generalization of NOSMOG on graphs with distinct structures.

**Cold-start Setting.** In latency-constrained systems, newly emerged nodes may become isolated (Hao et al., 2021; Zheng et al., 2022), which is commonly referred to as the cold-start issue. Our SSL-GM, which infers potential structural information of newly emerged nodes solely based on node content, provides a promising solution to address this issue. To simulate the cold-start scenario, we follow the inductive (production) setting while removing connections within the inductive set. This ensures that all inductive nodes remain isolated, forming a cold-start set  $\mathcal{G}^C = (\mathcal{V}^I, \emptyset)$ . Table 3 presents the cold-start performance results for SSL-GM and baselines. SSL-GM achieves notable improvements compared to all baselines. Specifically, SSL-GM demonstrates performance improvements of 7% and 18% over MLP on the Flickr and Arxiv datasets and achieves 7% and 7% enhancements over NOSMOG. We suppose the deficiency of NOSMOG derives from the absence of positional embedding. Note that BGRL, which employs augmentation in model training, also achieves exceptional performance, even if it highly relies on structural knowledge.

### 4.3 INFERENCE ACCELERATION

Our primary aim of SSL-GM is to accelerate inference. In this section, we demonstrate the capability of SSL-GM by illustrating the trade-off between prediction accuracy and model inference time using the Arxiv dataset under a cold-start setting, as depicted in Figure 3. We observe that SSL-GM attains the best trade-off between accuracy and inference time. Compared to baselines with similar inference times, SSL-GM outperforms them significantly, achieving an accuracy of 66%, whereas NOSMOG and MLPs only reach 62% and 56% accuracy, respectively. In contrast, methods that achieve performance similar to SSL-GM demand a substantial amount of inference time. For example, the 2-layer BGRL (BGRL-L2) requires 314.7ms, and the 3-layer BGRL (BGRL-L3) requires 635.9ms, while SSL-GM only needs 2.5ms, resulting in an acceleration of  $125\times$  and  $254\times$ , respectively. In the case of NOSMOG, increasing the hidden dimension (Zhang et al., 2022) can enhance performance by preserving more task-relevant information. We conduct a comparison between SSL-GM and NOSMOGw4 (four times wider than NOSMOG) and NOSMOGw8 (eight times wider than NOSMOG). We find that the wider models perform even worse than SSL-GM and demand more inference time. Consequently, we conclude that our SSL-GM surpasses existing baselines in accuracy while maintaining competitive inference speed. Additionally, we demonstrate that SSL-GM surpasses other acceleration methods in inference time on Flickr and Arxiv datasets ( $5\sim 90\times$ ) in Appendix C.

### 4.4 ABLATION STUDY

**Model Component.** In this section, we examine the impact of each component in the model. We provide a detailed analysis in Appendix D.3, D.4, D.5, and D.6. Table 1 shows the ablation results for ten benchmark datasets. The abbreviations *Aggr.*, *Pred.*, *Aug.*, and *Rec.* stand for aggregator, predictor, augmentation, and reconstruction. We observe the aggregator significantly enhances the model performance by injecting fine-grained structural information into the MLP encoder. The projector, which facilitates distance measurements between node representations, indeed improves the model performance, consistent with the findings in (Chen et al., 2020b). The augmentations improve the training of SSL-GM by synthesizing diverse node and ego-graph pairs during training. Moreover, the reconstruction term improves model performance by preventing the representation shifts of GNNs. Additionally, we implement the supervised version of SSL-GM in Appendix E.

**Robustness of SSL-GM to Noisy Data and Label Sparsity.** The quality of MLP encoder depends on both graph structure and node content. In this section, we introduce noise into the original graph to assess the robustness of SSL-GM on noisy data. Furthermore, we analyze the model performance

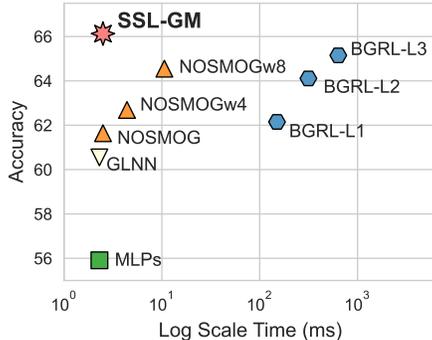


Figure 3: Accuracy vs. Inference Time on Arxiv dataset under cold-start setting.

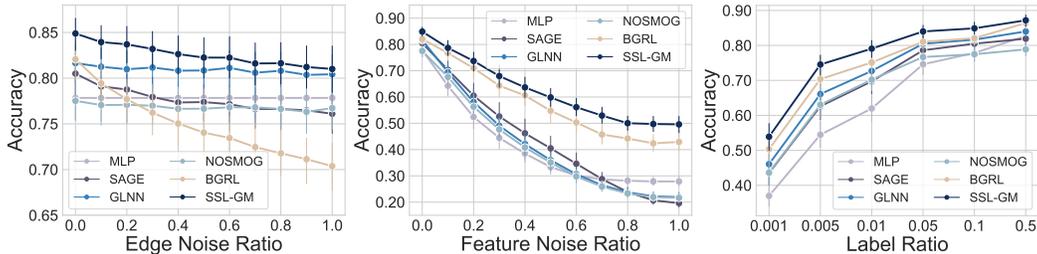


Figure 4: **Left:** Edge Noise. SSL-GM consistently demonstrates robustness against edge noise, even in scenarios with exceptionally high noise ratios. **Middle:** Feature Noise. SSL-GM demonstrates robustness to feature noise, whereas other MLP-based methods are susceptible to it. **Right:** Label Sparsity. SSL-GM significantly outperforms other baselines with a low label ratio.

with sparse labels. We report results on inductive set of production setting, averaged over seven datasets, including Cora, Citeseer, PubMed, Amazon-CS, Amazon-Photo, Co-CS, and Wiki-CS.

**(Noisy Topology)** To introduce structural noise, we randomly flip edges within the graph. Specifically, we replace  $\mathbf{A}$  with  $\tilde{\mathbf{A}} = \mathbf{M} \odot (1 - \mathbf{A}) + (1 - \mathbf{M}) \odot \mathbf{A}$ ,  $M_{ij} \sim \mathcal{B}(p)$ , where  $\mathcal{B}(p)$  is a Bernoulli distribution with probability  $p$ . The results under various noise levels are depicted in Figure 4 (**Left**). Our SSL-GM consistently outperforms others, demonstrating its robustness. The increase of noise level leads to substantial performance degradation for GNNs, particularly for self-supervised BGRL, but imposes minimal impact on MLPs, even  $\tilde{\mathbf{A}}$  becomes independent to  $\mathbf{A}$ .

**(Noisy Node Features)** Following, we examine the impact of node feature noise by introducing random Gaussian noise. We replace  $\mathbf{X}$  with  $\tilde{\mathbf{X}} = (1 - \alpha)\mathbf{X} + \alpha x$ , where  $x$  represents random noise agnostic to  $\mathbf{X}$ , and noise level  $\alpha \in [0, 1]$ . As depicted in Figure 4 (**Middle**), SSL-GM surpasses all baselines in all settings, even though node content quality is a crucial factor for MLP-based methods (Zhang et al., 2022; Guo et al., 2023). We attribute this robustness to the augmentation that synthesizes additional high-quality node and ego-graph pairs, thereby aiding MLP training. This augmentation also contributes to the robustness of BGRL. However, we observe that the performance of other MLP-based methods degrades rapidly with the increase in noise levels.

**(Label Sparsity)** In this section, we examine the robustness of SSL-GM under label sparsity. Figure 4 (**Right**) shows the model performance under various label ratios for node classification. We observe our method consistently outperforms all other baselines, even with extremely limited training data (0.001). This highlights the robustness of SSL-GM to label sparsity. Furthermore, we also observe that self-supervised methods exhibit better robustness Huang et al. (2023) compared to supervised methods, raised from the ability to leverage unlabeled data during training.

## 5 HOW SSL-GM LEARN FROM STRUCTURAL KNOWLEDGE?

**Mutual Information Maximization.** We interpret the SSL-GM from the perspective of information theory to theoretically analyze how MLPs encode structural information. Given a graph  $\mathcal{G} = (\mathbf{X}, \mathbf{A}, \mathbf{Y})$  with node features  $\mathbf{X}$ , graph structure  $\mathbf{A}$  and labels  $\mathbf{Y}$ , single MLP aims to minimize the cross-entropy between  $\mathbf{X}$  and  $\mathbf{Y}$ , which corresponds to maximize the mutual information  $\sum_{i \in \mathcal{V}} I(\mathbf{y}_i; \mathbf{x}_i)$  (Boudiaf et al., 2020) while disregarding the impact of structure. GNNs follows message passing framework that leverages subgraphs surrounding the target nodes to make prediction. We define subgraph around node  $i$  as  $\mathcal{S}_i = (\mathbf{X}^{[i]}, \mathbf{A}^{[i]})$ , where  $\mathbf{X}^{[i]}$  is the node features for neighborhoods of node  $i$ , and  $\mathbf{A}^{[i]}$  is the adjacent matrix that describes the aggregation rule. Thus, optimizing GNNs equals to maximize  $\sum_{i \in \mathcal{V}} I(\mathbf{y}_i; \mathcal{S}_i) = \sum_{i \in \mathcal{V}} I(\mathbf{y}_i; \mathbf{X}^{[i]}) + \sum_{i \in \mathcal{V}} I(\mathbf{y}_i; \mathbf{A}^{[i]} | \mathbf{X}^{[i]})$ , which models the correlation between label  $\mathbf{y}$  and both node feature  $\mathbf{X}$  and graph structure  $\mathbf{A}$ . The objective of GLNN is to maximize  $\sum_{i \in \mathcal{V}} I(\mathbf{x}_i; \mathbf{y}_i | \mathcal{S}_i)$ , where  $\mathbf{y}_i | \mathcal{S}_i$  denotes the soft labels given by GNNs. However, the approach cannot directly model the correlation between subgraph  $\mathcal{S}$  and label  $\mathbf{y}$ , preventing to acquire structural knowledge. Some models, e.g., GENN and NOSMOG, utilize positional encoding to incorporate structural knowledge based on GLNN, whose objective is to maximize  $\sum_{i \in \mathcal{V}} I(\mathbf{x}_i; \mathbf{y}_i | \mathcal{S}_i) + I(\mathbf{y}_i; \mathbf{A}^{[i]})$ . Although they capture the

structural knowledge to some extent, they only model the correlation between label  $\mathbf{y}$  and the graph structure  $\mathbf{A}$  instead of the subgraph  $\mathcal{S}$ , failing to model fine-grained structural knowledge.

Unlike these models, the aim of our SSL-GM is to maximize  $\sum_{i \in \mathcal{V}} I(\mathbf{y}_i; \mathbf{x}_i | \mathcal{S}_i) + I(\mathbf{x}_i; \mathcal{S}_i)$ . The first term optimizes the model on downstream tasks and the second term is the objective of SSL-GM. We argue that when the second term is maximized, the objective technically equals to maximize  $\sum_{i \in \mathcal{V}} I(\mathbf{y}_i; \mathcal{S}_i)$ , which corresponds to the objective of GNNs. The objective ensures SSL-GM comprehensively leverages graph information  $\mathcal{S}$ , including node content and graph structure, in downstream tasks. Our analysis also aligns with the findings in Chen et al. (2021) and Zhang et al. (2022) that the expressiveness of GNNs and MLPs are theoretically bounded by the equivalence classes of induced rooted graphs, which corresponds to  $\mathcal{S}$  in our case. Furthermore, SSL-GM augments node and ego-graph pairs to enhance the diversity of the training data, which inherently improves the quality of  $\mathcal{S}$  and thereby improve the optimizing process of our SSL-GM.

**Graph Smoothness.** In addition to theoretical analysis, we empirically demonstrate the capability of SSL-GM in learning structural information as an inductive bias, which can be measured by graph smoothness. A low smoothness value indicates that representations of closely connected nodes are similar, allowing the model to extract more information from the graph data (Hou et al., 2019). We employ the Mean Average Distance (MAD)  $\mathcal{L}_{MAD} = \frac{\sum_{v \in \mathcal{V}} \sum_{j \in \mathcal{N}(i)} (\mathbf{H}_i - \mathbf{H}_j)^2}{\sum_{v \in \mathcal{V}} \sum_{j \in \mathcal{N}(i)} \mathbf{1}}$  (Chen et al., 2020a) to quantify graph smoothness. For supervised meth-

	Cora	Citeseer	PubMed	Amazon-CS	Amazon-Photo	Average
Raw Feat.	0.8221	0.7825	0.7342	0.5393	0.5399	0.6836
SAGE	0.1132	0.1835	0.1426	0.1564	0.1089	0.1409
BGRL	0.1553	0.1023	0.3326	0.2509	0.2031	0.2088
MLP	0.4633	0.4442	0.4853	0.4557	0.4317	0.4560
GLNN	0.2818	0.2684	0.4208	0.3549	0.3976	0.3447
NOSMOG	0.2672	0.2301	0.3942	0.3056	<b>0.2773</b>	0.2949
SSL-GM	<b>0.1964</b>	<b>0.1703</b>	<b>0.3604</b>	<b>0.2986</b>	0.2878	<b>0.2627</b>

Table 4: The graph smoothness value measures the representation similarity between nodes and their neighborhoods. Lower values indicate a smoother graph, signifying that the model can capture more structural knowledge.

ods, the representations are extracted from the output of the final layer before the prediction head. Table 4 presents the graph smoothness values for SSL-GM and the baseline methods. We observe that message passing methods inherently yield lower smoothness values than MLP-based methods. SSL-GM, which pulls close the outputs of MLPs and GNNs in the representation space, achieves an average smoothness value of 0.2627. This value is lower than that of MLP (0.4560), GLNN (0.3447), and NOSMOG (0.2949), indicating the superiority of learning fine-grained structural information. Note that a low smoothness value might potentially lead to over-smoothing (Li et al., 2019), but this topic falls outside the scope of this paper. We utilize smoothness to indicate whether the encoded representations can accurately reflect the graph structure.

**Normalized Cut.** Graph smoothness evaluates the consistency between representations and graph structure. To further analyze the alignment between predictions and structure, we adopt normalized cut, which approximates the min-cut problem. This problem involves minimizing the number of removed edges to partition nodes  $\mathcal{V}$  into  $K$  disjoint subsets. The min-cut problem can be formulated

as  $\mathcal{L}_{cut} = \frac{\text{tr}(\hat{\mathbf{Y}}^T \mathbf{A} \hat{\mathbf{Y}})}{\text{tr}(\hat{\mathbf{Y}}^T \mathbf{D} \hat{\mathbf{Y}})}$  where  $\hat{\mathbf{Y}}$  represents the model predictions. A higher  $\mathcal{L}_{cut}$  value indicates a better alignment between model predictions and the underlying structure. Table 5 presents the evaluation of min-cut scores for SSL-GM and baselines across five datasets. We observe message passing methods generally outperform MLPs due to their explicit encoding of structure. Remarkably, SSL-GM surpasses existing MLP-based methods and is even competitive with message passing models, demonstrating its superior ability to capture structural information.

	Cora	Citeseer	PubMed	Amazon-CS	Amazon-Photo	Average
SAGE	0.9243	0.9426	0.9177	0.8541	0.8721	0.9022
BGRL	0.8847	0.9346	0.8556	0.8336	0.8493	0.8716
MLP	0.6663	0.8035	0.8625	0.7183	0.7467	0.7595
GLNN	0.8863	0.9162	0.7934	0.8038	0.8113	0.8422
NOSMOG	0.9023	0.9317	0.8337	0.8384	0.8226	0.8657
SSL-GM	<b>0.9335</b>	<b>0.9575</b>	<b>0.8863</b>	<b>0.9014</b>	<b>0.8604</b>	<b>0.9078</b>

Table 5: Our model predictions are more consistent with the graph structure in terms of normalized cut value, although the model is trained in an unsupervised manner.

## 6 CONCLUSION

We analyze that existing MLP-based graph inference acceleration methods cannot learn generalizable structure-aware representations, and propose a novel framework SSL-GM to solve it. The key insight is that the potential structural information of unseen nodes can be inferred solely based on the node content with SSL. Specifically, we employ self-supervised contrastive learning to align GNNs and MLPs in the representation space to integrate rich structural information into MLPs. Additionally, we apply non-parametric aggregator, graph augmentation, and feature reconstruction to improve model performance. We also empirically and theoretically demonstrate the expressiveness, generalization, and robustness of SSL-GM. In the future, we will extend SSL-GM to real-world applications, such as financial fraud detection. We hope our work can inspire researchers seeking to develop novel learning algorithms for graphs that go beyond the scope of GNNs.

## REFERENCES

- Alexander A. Alemi, Ian Fischer, Joshua V. Dillon, and Kevin Murphy. Deep variational information bottleneck. In *ICLR*, 2017.
- Joshua Batson and Loic Royer. Noise2self: Blind denoising by self-supervision. In *ICML*, 2019.
- Malik Boudiaf, Jérôme Rony, Imtiaz Masud Ziko, Eric Granger, Marco Pedersoli, Pablo Piantanida, and Ismail Ben Ayed. A unifying mutual information view of metric learning: cross-entropy vs. pairwise losses. In *ECCV*, 2020.
- Vivien Cabannes, Bobak Kiani, Randall Balestriero, Yann LeCun, and Alberto Bietti. The ssl interplay: Augmentations, inductive bias, and generalization. In *ICML*, 2023.
- Deli Chen, Yankai Lin, Wei Li, Peng Li, Jie Zhou, and Xu Sun. Measuring and relieving the over-smoothing problem for graph neural networks from the topological view. In *AAAI*, 2020a.
- Lei Chen, Zhengdao Chen, and Joan Bruna. On graph neural networks versus graph-augmented MLPs. In *ICLR*, 2021.
- Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *ICML*, 2020b.
- Mucong Ding, Kezhi Kong, Jingling Li, Chen Zhu, John Dickerson, Furong Huang, and Tom Goldstein. VQ-GNN: A universal framework to scale up graph neural networks using vector quantization. In *NeurIPS*, 2021.
- Wei Dong, Junsheng Wu, Yi Luo, Zongyuan Ge, and Peng Wang. Node representation learning in graph via node-to-neighbourhood mutual information maximization. In *CVPR*, 2022.
- Vijay Prakash Dwivedi and Xavier Bresson. A generalization of transformer networks to graphs. *arXiv*, 2020.
- Vijay Prakash Dwivedi, Anh Tuan Luu, Thomas Laurent, Yoshua Bengio, and Xavier Bresson. Graph neural networks with learnable structural and positional representations. In *ICLR*, 2022.
- Shengyu Feng, Baoyu Jing, Yada Zhu, and Hanghang Tong. Adversarial graph contrastive learning with information regularization. In *WWW*, 2022.
- Wenzheng Feng, Jie Zhang, Yuxiao Dong, Yu Han, Huanbo Luan, Qian Xu, Qiang Yang, Evgeny Kharlamov, and Jie Tang. Graph random neural networks for semi-supervised learning on graphs. In *NeurIPS*, 2020.
- Jonathan Frankle and Michael Carbin. The lottery ticket hypothesis: Finding sparse, trainable neural networks. In *ICLR*, 2019.
- Johannes Gasteiger, Aleksandar Bojchevski, and Stephan Günnemann. Combining neural networks with personalized pagerank for classification on graphs. In *ICLR*, 2019.

- Justin Gilmer, Samuel S Schoenholz, Patrick F Riley, Oriol Vinyals, and George E Dahl. Neural message passing for quantum chemistry. In *ICML*, 2017.
- Jean-Bastien Grill, Florian Strub, Florent Alché, Corentin Tallec, Pierre Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Guo, Mohammad Gheshlaghi Azar, et al. Bootstrap your own latent-a new approach to self-supervised learning. In *NeurIPS*, 2020.
- Zhichun Guo, William Shiao, Shichang Zhang, Yozen Liu, Nitesh V Chawla, Neil Shah, and Tong Zhao. Linkless link prediction via relational distillation. In *ICML*, 2023.
- Suyog Gupta, Ankur Agrawal, Kailash Gopalakrishnan, and Pritish Narayanan. Deep learning with limited numerical precision. In *ICML*, 2015.
- Will Hamilton, Zhitao Ying, and Jure Leskovec. Inductive representation learning on large graphs. In *NeurIPS*, 2017.
- Song Han, Jeff Pool, John Tran, and William Dally. Learning both weights and connections for efficient neural network. In *NeurIPS*, 2015.
- Xiaotian Han, Tong Zhao, Yozen Liu, Xia Hu, and Neil Shah. MLPInit: Embarrassingly simple GNN training acceleration with MLP initialization. In *ICLR*, 2023.
- Bowen Hao, Jing Zhang, Hongzhi Yin, Cuiping Li, and Hong Chen. Pre-training graph neural networks for cold-start users and items representation. In *WSDM*, 2021.
- Kaveh Hassani and Amir Hosein Khasahmadi. Contrastive multi-view representation learning on graphs. In *ICML*, 2020.
- Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *CVPR*, 2020.
- Dan Hendrycks, Mantas Mazeika, Saurav Kadavath, and Dawn Song. Using self-supervised learning can improve model robustness and uncertainty. In *NeurIPS*, 2019.
- Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv*, 2015.
- Yifan Hou, Jian Zhang, James Cheng, Kaili Ma, Richard TB Ma, Hongzhi Chen, and Ming-Chang Yang. Measuring and improving the use of graph information in graph neural networks. In *ICLR*, 2019.
- Zhenyu Hou, Xiao Liu, Yukuo Cen, Yuxiao Dong, Hongxia Yang, Chunjie Wang, and Jie Tang. Graphmae: Self-supervised masked graph autoencoders. In *KDD*, 2022.
- Weihua Hu, Matthias Fey, Marinka Zitnik, Yuxiao Dong, Hongyu Ren, Bowen Liu, Michele Catasta, and Jure Leskovec. Open graph benchmark: Datasets for machine learning on graphs. In *NeurIPS*, 2020.
- Yang Hu, Haoxuan You, Zhecan Wang, Zhicheng Wang, Erjin Zhou, and Yue Gao. Graph-mlp: Node classification without message passing in graph. *arXiv*, 2021.
- Weiran Huang, Mingyang Yi, Xuyang Zhao, and Zihao Jiang. Towards the generalization of contrastive self-supervised learning. In *ICLR*, 2023.
- Benoit Jacob, Skirmantas Kligys, Bo Chen, Menglong Zhu, Matthew Tang, Andrew Howard, Hartwig Adam, and Dmitry Kalenichenko. Quantization and training of neural networks for efficient integer-arithmic-only inference. In *CVPR*, 2018.
- Huajie Jiang, Ruiping Wang, Shiguang Shan, and Xilin Chen. Transferable contrastive network for generalized zero-shot learning. In *ICCV*, 2019.
- Thomas N. Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. In *ICLR*, 2017.

- Kezhi Kong, Guohao Li, Mucong Ding, Zuxuan Wu, Chen Zhu, Bernard Ghanem, Gavin Taylor, and Tom Goldstein. Robust optimization as data augmentation for large-scale graphs. In *CVPR*, 2022.
- Devin Kreuzer, Dominique Beaini, Will Hamilton, Vincent Létourneau, and Prudencio Tossou. Rethinking graph transformers with spectral attention. In *NeurIPS*, 2021.
- Guohao Li, Matthias Muller, Ali Thabet, and Bernard Ghanem. Deepgcns: Can gcns go as deep as cnns? In *ICCV*, 2019.
- Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *ICCV*, 2017.
- Siwei Liu, Iadh Ounis, and Craig Macdonald. An MLP-based algorithm for efficient contrastive graph recommendations. In *SIGIR*, 2022.
- Sitao Luan, Chenqing Hua, Qincheng Lu, Jiaqi Zhu, Xiao-Wen Chang, and Doina Precup. When do we need gnn for node classification? *arXiv*, 2022.
- Sitao Luan, Chenqing Hua, Minkai Xu, Qincheng Lu, Jiaqi Zhu, Xiao-Wen Chang, Jie Fu, Jure Leskovec, and Doina Precup. When do graph neural networks help with node classification: Investigating the homophily principle on node distinguishability. *arXiv*, 2023.
- Péter Mernyei and Cătălina Cangea. Wiki-cs: A wikipedia-based benchmark for graph neural networks. *arXiv*, 2020.
- Christopher Morris, Nils M Kriege, Franka Bause, Kristian Kersting, Petra Mutzel, and Marion Neumann. Tudataset: A collection of benchmark datasets for learning with graphs. *arXiv*, 2020.
- Annamalai Narayanan, Mahinthan Chandramohan, Rajasekar Venkatesan, Lihui Chen, Yang Liu, and Shantanu Jaiswal. graph2vec: Learning distributed representations of graphs. *arXiv*, 2017.
- Bryan Perozzi, Rami Al-Rfou, and Steven Skiena. Deepwalk: Online learning of social representations. In *KDD*, 2014.
- Ladislav Rampášek, Michael Galkin, Vijay Prakash Dwivedi, Anh Tuan Luu, Guy Wolf, and Dominique Beaini. Recipe for a general, powerful, scalable graph transformer. In *NeurIPS*, 2022.
- Oleksandr Shchur, Maximilian Mumme, Aleksandar Bojchevski, and Stephan Günnemann. Pitfalls of graph neural network evaluation. *arXiv*, 2018.
- Nino Shervashidze, Pascal Schweitzer, Erik Jan Van Leeuwen, Kurt Mehlhorn, and Karsten M Borgwardt. Weisfeiler-lehman graph kernels. *JMLR*, 2011.
- Ravid Shwartz-Ziv and Naftali Tishby. Opening the black box of deep neural networks via information. *arXiv*, 2017.
- Fan-Yun Sun, Jordan Hoffman, Vikas Verma, and Jian Tang. Infograph: Unsupervised and semi-supervised graph-level representation learning via mutual information maximization. In *ICLR*, 2020.
- Xiangguo Sun, Hong Cheng, Jia Li, Bo Liu, and Jihong Guan. All in one: Multi-task prompting for graph neural networks. In *KDD*, 2023.
- Susheel Suresh, Pan Li, Cong Hao, and Jennifer Neville. Adversarial graph augmentation to improve graph contrastive learning. In *NeurIPS*, 2021.
- Shantanu Thakoor, Corentin Tallec, Mohammad Gheshlaghi Azar, Mehdi Azabou, Eva L Dyer, Remi Munos, Petar Veličković, and Michal Valko. Large-scale representation learning on graphs via bootstrapping. In *ICLR*, 2022.
- Yijun Tian, Chuxu Zhang, Zhichun Guo, Xiangliang Zhang, and Nitesh Chawla. Learning MLPs on graphs: A unified view of effectiveness, robustness, and efficiency. In *ICLR*, 2023.

- Yonglong Tian, Dilip Krishnan, and Phillip Isola. Contrastive representation distillation. In *ICLR*, 2020.
- Naftali Tishby and Noga Zaslavsky. Deep learning and the information bottleneck principle. In *ITW*, 2015.
- Naftali Tishby, Fernando C Pereira, and William Bialek. The information bottleneck method. *arXiv*, 2000.
- Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio. Graph attention networks. In *ICLR*, 2018.
- Petar Veličković, William Fedus, William L. Hamilton, Pietro Liò, Yoshua Bengio, and R Devon Hjelm. Deep graph infomax. In *ICLR*, 2019.
- Haorui Wang, Haoteng Yin, Muhan Zhang, and Pan Li. Equivariant and stable positional encoding for more powerful graph neural networks. In *ICLR*, 2022.
- Xuhong Wang, Ding Lyu, Mengjian Li, Yang Xia, Qi Yang, Xinwen Wang, Xinguang Wang, Ping Cui, Yupu Yang, Bowen Sun, et al. Apan: Asynchronous propagation attention network for real-time temporal graph embedding. In *SIGMOD*, 2021.
- Yiwei Wang, Bryan Hooi, Yozen Liu, and Neil Shah. Graph explicit neural networks: Explicitly encoding graphs for efficient and accurate inference. In *WSDM*, 2023.
- Felix Wu, Amauri Souza, Tianyi Zhang, Christopher Fifty, Tao Yu, and Kilian Weinberger. Simplifying graph convolutional networks. In *ICML*, 2019.
- Yaochen Xie, Zhao Xu, and Shuiwang Ji. Self-supervised representation learning via latent graph prediction. In *ICML*, 2022.
- Dongkuan Xu, Wei Cheng, Dongsheng Luo, Haifeng Chen, and Xiang Zhang. Infogcl: Information-aware graph contrastive learning. In *NeurIPS*, 2021.
- Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural networks? In *ICLR*, 2019.
- Bencheng Yan, Chaokun Wang, Gaoyang Guo, and Yunkai Lou. Tinygcn: Learning efficient graph neural networks. In *KDD*, 2020.
- Pinar Yanardag and SVN Vishwanathan. Deep graph kernels. In *KDD*, 2015.
- Cheng Yang, Jiawei Liu, and Chuan Shi. Extract the knowledge of graph neural networks and go beyond it: An effective knowledge distillation framework. In *WWW*, 2021.
- Chenxiao Yang, Qitian Wu, Jiahua Wang, and Junchi Yan. Graph neural networks are inherently good generalizers: Insights by bridging GNNs and MLPs. In *ICLR*, 2023a.
- Ling Yang, Ye Tian, Minkai Xu, Zhongyi Liu, Shenda Hong, Wei Qu, Wentao Zhang, Bin Cui, Muhan Zhang, and Jure Leskovec. Vqgraph: Graph vector-quantization for bridging gnn and mlps. *arXiv*, 2023b.
- Zhilin Yang, William Cohen, and Ruslan Salakhudinov. Revisiting semi-supervised learning with graph embeddings. In *ICML*, 2016.
- Yuning You, Tianlong Chen, Yongduo Sui, Ting Chen, Zhangyang Wang, and Yang Shen. Graph contrastive learning with augmentations. In *NeurIPS*, 2020.
- Yuning You, Tianlong Chen, Yang Shen, and Zhangyang Wang. Graph contrastive learning automated. In *ICML*, 2021.
- Hanqing Zeng, Hongkuan Zhou, Ajitesh Srivastava, Rajgopal Kannan, and Viktor Prasanna. Graph-saint: Graph sampling based inductive learning method. In *ICLR*, 2020.

Shichang Zhang, Yozen Liu, Yizhou Sun, and Neil Shah. Graph-less neural networks: Teaching old MLPs new tricks via distillation. In *ICLR*, 2022.

Tong Zhao, Yozen Liu, Leonardo Neves, Oliver Woodford, Meng Jiang, and Neil Shah. Data augmentation for graph neural networks. In *AAAI*, 2021.

Wenqing Zheng, Edward W Huang, Nikhil Rao, Sumeet Katariya, Zhangyang Wang, and Karthik Subbian. Cold brew: Distilling graph node representations with incomplete or missing neighborhoods. In *ICLR*, 2022.

Hongkuan Zhou, Ajitesh Srivastava, Hanqing Zeng, Rajgopal Kannan, and Viktor Prasanna. Accelerating large scale real-time gnn inference using channel pruning. In *VLDB*, 2021.

Yanqiao Zhu, Yichen Xu, Feng Yu, Qiang Liu, Shu Wu, and Liang Wang. Deep graph contrastive representation learning. *arXiv*, 2020.

Yanqiao Zhu, Yichen Xu, Feng Yu, Qiang Liu, Shu Wu, and Liang Wang. Graph contrastive learning with adaptive augmentation. In *WWW*, 2021.

## A PROOF OF THEOREM 1

**Theorem 1** Suppose  $\mathcal{G} = (\mathbf{A}, \mathbf{X})$  is sampled from a latent graph  $\mathcal{G}_{\mathcal{I}} = (\mathbf{A}, \mathbf{F})$ ,  $\mathcal{G} \sim P(\mathcal{G}_{\mathcal{I}})$  (Xie et al., 2022), and  $\mathbf{F}^*$  is the lossless compression of  $\mathbf{F}$  that  $\mathbb{E}[\mathbf{X}|\mathbf{A}, \mathbf{F}^*] = \mathbf{F}$ . Let  $\mathcal{E}$  be a  $l$ -Lipschitz continuous function respect to  $l_2$ -norm,  $\rho$  be an identity projector, and  $\lambda = 1, \gamma = 1$ . Optimizing Eq. 6 equals to finding the optimal compression  $\mathbf{T}^*$  with minimal sufficient information  $\mathbf{C}$  where

$$\mathbf{T}^* = \arg \min_{\mathbf{T}} I(\mathcal{G}; \mathbf{T}) - \beta I(\mathbf{T}; \mathcal{G}_{\mathcal{I}}), \text{ s.t.}, I(\mathbf{T}, \mathcal{G}_{\mathcal{I}}) \geq \mathbf{C}, \mathbf{T} = (\mathbf{H}, \mathbf{Z})$$

*Proof.* Consider  $\mathcal{G} = (\mathbf{A}, \mathbf{X})$ , where  $\mathbf{X} \in \mathbb{R}^{N \times d}$  is derived from a latent graph  $\mathcal{G}_{\mathcal{I}} = (\mathbf{A}, \mathbf{F})$  following a distribution  $\mathcal{G} \sim P(\mathcal{G}_{\mathcal{I}})$ , with  $\mathbf{F} \in \mathbb{R}^{N \times d}$  representing the latent node semantics. Let  $\hat{\mathcal{G}} = (\hat{\mathbf{X}}, \hat{\mathbf{A}})$  denote the randomly augmented version of  $\mathcal{G}$ , achieved by applying augmentations to both node features and graph structure. Additionally, consider an encoder  $\mathcal{E}$  and a decoder  $\mathcal{D}$  implemented as fully-connected layers, ensuring  $l$ -Lipschitz continuity with respect to the  $l_2$ -norm, and a non-parametric aggregator  $\phi$ . This yields  $\hat{\mathbf{H}} = \mathcal{E}(\hat{\mathbf{X}})$  with  $\hat{\mathbf{H}} \in \mathbb{R}^{N \times d'}$ , and  $\hat{\mathbf{Z}} = \phi(\hat{\mathbf{H}}, \hat{\mathbf{A}})$  with  $\hat{\mathbf{Z}} \in \mathbb{R}^{N \times d'}$ . Furthermore,  $\mathbf{F}^* \in \mathbb{R}^{N \times d'}$  denotes the lossless compression of  $\mathbf{F}$  that  $\mathbb{E}[\mathbf{X}|\mathbf{A}, \mathbf{F}^*] = \mathbf{F}$ . The Eq. 6 can be rewritten as:

$$\mathcal{E}^* = \arg \min_{\mathcal{E}} \mathbb{E}_{\hat{\mathbf{A}}, \hat{\mathbf{X}}} \left[ \|\hat{\mathbf{H}} - \hat{\mathbf{Z}}\|^2 + \|\mathcal{D}(\hat{\mathbf{Z}}) - \mathbf{X}\|^2 | \hat{\mathbf{A}}, \hat{\mathbf{X}} \right] \quad (8)$$

$$= \arg \min_{\mathcal{E}} \mathbb{E}_{\hat{\mathbf{A}}, \hat{\mathbf{X}}} \left[ \|(\hat{\mathbf{H}} - \mathbf{F}^*) - (\hat{\mathbf{Z}} - \mathbf{F}^*)\|^2 + \|\mathcal{D}(\hat{\mathbf{Z}}) - \mathbf{X}\|^2 | \hat{\mathbf{A}}, \hat{\mathbf{X}} \right] \quad (9)$$

$$= \arg \min_{\mathcal{E}} \mathbb{E}_{\hat{\mathbf{A}}, \hat{\mathbf{X}}} \left[ \|\hat{\mathbf{H}} - \mathbf{F}^*\|^2 + \|\hat{\mathbf{Z}} - \mathbf{F}^*\|^2 + \|\mathcal{D}(\hat{\mathbf{Z}}) - \mathbf{X}\|^2 | \hat{\mathbf{A}}, \hat{\mathbf{X}} \right] \\ - 2\mathbb{E}_{\hat{\mathbf{A}}, \hat{\mathbf{X}}} \left[ \langle \hat{\mathbf{H}} - \mathbf{F}^*, \hat{\mathbf{Z}} - \mathbf{F}^* \rangle | \hat{\mathbf{A}}, \hat{\mathbf{X}} \right] \quad (10)$$

$$= \arg \min_{\mathcal{E}} \mathbb{E}_{\hat{\mathbf{A}}, \hat{\mathbf{X}}} \left[ \|\hat{\mathbf{H}} - \mathbf{F}^*\|^2 + \|\hat{\mathbf{Z}} - \mathbf{F}^*\|^2 + \|\mathcal{D}(\hat{\mathbf{Z}}) - \mathbf{X}\|^2 | \hat{\mathbf{A}}, \hat{\mathbf{X}} \right] \\ - 2\mathbb{E}_{\hat{\mathbf{A}}, \hat{\mathbf{X}}, \mathbf{F}^*} \left[ \sum_i (\hat{H}_i - \mathbf{F}_i^*)(\hat{Z}_i - \mathbf{F}_i^*) | \hat{\mathbf{A}}, \hat{\mathbf{X}}, \mathbf{F}^* \right] \quad (11)$$

$$= \arg \min_{\mathcal{E}} \mathbb{E}_{\hat{\mathbf{A}}, \hat{\mathbf{X}}} \left[ \|\hat{\mathbf{H}} - \mathbf{F}^*\|^2 + \|\hat{\mathbf{Z}} - \mathbf{F}^*\|^2 + \|\mathcal{D}(\hat{\mathbf{Z}}) - \mathbf{X}\|^2 | \hat{\mathbf{A}}, \hat{\mathbf{X}} \right] \\ - 2\mathbb{E}_{\hat{\mathbf{A}}, \hat{\mathbf{X}}, \mathbf{F}^*} \left[ \sum_i \text{Cov}(\hat{H}_i - \mathbf{F}_i^*, \hat{Z}_i - \mathbf{F}_i^*) | \hat{\mathbf{A}}, \hat{\mathbf{X}}, \mathbf{F}^* \right] \quad (12)$$

$$= \arg \min_{\mathcal{E}} \mathbb{E}_{\hat{\mathbf{A}}, \hat{\mathbf{X}}} \left[ \|\hat{\mathbf{H}} - \mathbf{F}^*\|^2 + \|\hat{\mathbf{Z}} - \mathbf{F}^*\|^2 + \|\mathcal{D}(\hat{\mathbf{Z}}) - \mathbf{X}\|^2 | \hat{\mathbf{A}}, \hat{\mathbf{X}} \right] \\ - 2\mathbb{E}_{\hat{\mathbf{A}}, \hat{\mathbf{X}}, \mathbf{F}^*} \left[ \sum_i \text{Cov}(\hat{H}_i, \hat{Z}_i) | \hat{\mathbf{A}}, \hat{\mathbf{X}}, \mathbf{F}^* \right]. \quad (13)$$

The terms  $\|\hat{\mathbf{H}} - \mathbf{F}^*\|^2$  and  $\|\hat{\mathbf{Z}} - \mathbf{F}^*\|^2$  represent the reconstruction errors of MLP representations  $\hat{\mathbf{H}}$  and GNN representations  $\hat{\mathbf{Z}}$  on the latent graph  $\mathcal{G}_{\mathcal{I}}$ . These errors ensure the invariance of these representations with respect to latent semantics. The term  $\|\mathcal{D}(\hat{\mathbf{Z}}) - \mathbf{X}\|^2$  optimizes the distance between the reconstructed GNN representations and the node features. This optimization prevents representation shifts caused by augmentations. The final term,  $-\sum_i \text{Cov}(\hat{H}_i, \hat{Z}_i)$ , quantifies the covariance between the representations encoded by MLPs and GNNs. Maximizing this covariance ensures that the encoder  $\mathcal{E}$  effectively captures structural knowledge.

We will now interpret SSL-GM from the perspective of the information bottleneck. The fundamental idea behind the information bottleneck (Tishby et al., 2000; Tishby & Zaslavsky, 2015; Shwartz-Ziv & Tishby, 2017) is to compress the original information while retaining the latent information. In our case, the information to be compressed is the original graph  $\mathcal{G}$ , while we consider  $\mathcal{G}_{\mathcal{I}}$  as the latent information. We regard the compressed information as the combination of MLP and GNN representations, denoted as  $\mathbf{T} = (\mathbf{H}, \mathbf{Z})$ . Thus, we define the information bottleneck as  $\mathbf{T}^* = \arg \min_{\mathbf{T}} I(\mathcal{G}; \mathbf{T}) - \beta I(\mathbf{T}; \mathcal{G}_{\mathcal{I}})$ . To present the information bottleneck in a more accessible manner, we perform the following transformation:

$$\mathbf{T}^* = \arg \min_{\mathbf{T}} I(\mathcal{G}; \mathbf{T}) - \beta I(\mathbf{T}; \mathcal{G}_{\mathcal{I}}) \quad (14)$$

$$= \arg \min_{\mathbf{T}} (1 - \beta)H(\mathbf{T}) + \beta H(\mathbf{T} | \mathcal{G}_{\mathcal{I}}) - H(\mathbf{T} | \mathcal{G}) \quad (15)$$

$$= \arg \min_{\mathbf{T}} H(\mathbf{T}) + \lambda H(\mathbf{T} | \mathcal{G}_{\mathcal{I}}) \quad (16)$$

$$= \arg \min_{\mathbf{H}, \mathbf{Z}} \lambda H(\mathbf{H} | \mathcal{G}_{\mathcal{I}}) + \lambda H(\mathbf{Z} | \mathbf{H}, \mathcal{G}_{\mathcal{I}}) + H(\mathbf{Z}) + H(\mathbf{H} | \mathbf{Z}), \quad (17)$$

where  $\lambda = \frac{\beta}{1-\beta} > 0$ . We consider that the four terms in Eq. 17 correspond to the four terms in Eq. 13. Specifically, minimizing the conditional entropy  $H(\mathbf{H} | \mathcal{G}_{\mathcal{I}})$  and  $H(\mathbf{Z} | \mathbf{H}, \mathcal{G}_{\mathcal{I}})$  in Eq. 17 can reduce the uncertainty of  $\mathbf{H}$  and  $\mathbf{Z}$  with respect to the latent graph  $\mathcal{G}_{\mathcal{I}}$ , thereby ensuring the invariance of  $\mathbf{H}$  and  $\mathbf{Z}$  concerning  $\mathcal{G}_{\mathcal{I}}$ . Consequently, these two terms correspond to the reconstruction terms  $\|\hat{\mathbf{H}} - \mathbf{F}^*\|^2$  and  $\|\hat{\mathbf{Z}} - \mathbf{F}^*\|^2$  in Eq. 13, respectively. Additionally, we posit that optimizing the third term  $H(\mathbf{Z})$  in Eq. 17 reduces the uncertainty of  $\mathbf{Z}$ . This aligns with the third term  $\|\mathcal{D}(\hat{\mathbf{Z}}) - \mathbf{X}\|^2$

in Eq. 13, which constrains  $\mathbf{Z}$  to be invariant to  $\mathbf{X}$ . Furthermore, concerning the final term in Eq. 17, which minimizes the conditional entropy  $H(\mathbf{H}|\mathbf{Z})$ , it guarantees the consistency between MLP representations  $\mathbf{H}$  and GNN representations  $\mathbf{Z}$ . This alignment can be attained by maximizing the covariance term  $\sum_i \text{Cov}(\hat{\mathbf{H}}_i, \hat{\mathbf{Z}}_i)$ , corresponding to the last term in Eq. 13. This analysis reveals that minimizing our objective is equivalent to optimizing the information bottleneck.

## B EXPERIMENT SETUP DETAILS

In this section, we provide a comprehensive description of the experimental setup for both node classification and graph classification. The experiments are conducted on Nvidia A100 (80GB) for the Arxiv dataset, and Nvidia GeForce RTX 3090 (24GB) for the remaining datasets.

### B.1 DATASET STATISTICS

Dataset	Task	# Graphs	# Nodes	# Edges	# Features	# Classes	Split
Cora	Node-level	1	2,708	10,556	1,433	7	10%/10%/80%
Citeseer	Node-level	1	3,327	9,104	3,703	6	10%/10%/80%
PubMed	Node-level	1	19,717	88,648	500	3	10%/10%/80%
Amazon-CS	Node-level	1	13,752	491,722	767	10	10%/10%/80%
Amazon-Photo	Node-level	1	7,650	238,162	745	8	10%/10%/80%
Co-CS	Node-level	1	18,333	163,788	6,805	15	10%/10%/80%
Co-Phys	Node-level	1	34,493	495,924	8,415	5	10%/10%/80%
Wiki-CS	Node-level	1	11,701	432,246	300	10	10%/10%/80%
Flickr	Node-level	1	89,250	899,756	500	7	10%/10%/80%
Arxiv	Node-level	1	169,343	1,166,243	128	40	Public Split

Dataset	Task	# Graphs	# Nodes	# Edges	# Features	# Classes	Split
IMDB-B	Graph-level	1,000	~19.8	~193.1	-	2	10-fold CV
IMDB-M	Graph-level	1,500	~13.0	~65.9	-	3	10-fold CV
COLLAB	Graph-level	5,000	~74.5	~4,914.4	-	3	10-fold CV
PTC-MR	Graph-level	344	~14.3	~14.7	18	2	10-fold CV
MUTAG	Graph-level	118	~17.9	~39.6	7	2	10-fold CV
DD	Graph-level	1,178	~284.3	~715.6	89	2	10-fold CV
PROTEINS	Graph-level	1,113	~39.1	~145.6	3	2	10-fold CV

Table 6: The statistics of datasets for node-level and graph-level tasks.

We select 17 benchmark datasets, with 10 designated for node classification and 7 for graph classification, to evaluate the performance of SSL-GM and other approaches. These datasets are collected from diverse domains, encompassing citation networks, social networks, molecule networks, etc. We present the statistics of these datasets in Table 6.

**Node Classification:** Specifically, **Cora**, **Citeseer**, **PubMed** (Yang et al., 2016) are three citation networks, in which nodes denote papers and edges represent citations. The node features are represented as bag-of-words based on paper keywords. **Amazon-CS** and **Amazon-Photo** (Shchur et al., 2018) are two co-purchase networks that describe the frequent co-purchases of items (nodes). **Co-CS (Coauthor-CS)** and **Co-Phys (Coauthor-physics)** (Shchur et al., 2018) consist of nodes representing authors and edges indicating collaborations between authors. **Wiki-CS** (Mernyei & Cangea, 2020) is extracted from Wikipedia, comprising computer science articles (nodes) connected by hyperlinks (edges). **Flickr** (Zeng et al., 2020) consists online images, with the goal of categorizing images based on their descriptions and common properties. All these datasets are available through PyG (Pytorch Geometric), and we partition them randomly into training, validation, and testing sets with a split ratio of 10%/10%/80%. Additionally, we employ **Arxiv** dataset from OGB benchmarks (Hu et al., 2020) to evaluate model performance on large-scale datasets. We process the dataset in PyG using OGB public interfaces with standard public split setting.

**Graph Classification:** All graph classification datasets are sourced from TU datasets (Morris et al., 2020). We employ several datasets, including biochemical molecule datasets (**PTC-MR**, **MUTAG**, **DD**, **PROTEINS**) and social networks (**IMDB-B**, **IMDB-M**, **COLLAB**), to access whether SSL-GM can acquire generalizable and global information. In the case of the PTC-MR and DD datasets, we utilize the original node features, whereas for other datasets lacking rich node features, we generate one-hot features based on node degrees. These datasets are available in PyG library following a 10-fold cross validation data split.

## B.2 SUMMARY OF BASELINES

We compare SSL-GM against a range of baselines, encompassing supervised GNNs, self-supervised graph contrastive learning (GCL) methods, and MLP-based graph learning methods.

**Supervised GNNs.** Our primary node classification baselines include **GraphSAGE** (Hamilton et al., 2017) and **GAT** (Veličković et al., 2018), while for graph classification, we utilize **GIN** (Xu et al., 2019). Furthermore, we also incorporate **SGC** (Wu et al., 2019) and **APPNP** (Gasteiger et al., 2019) as additional node classification baselines.

**Self-supervised GNNs.** We compare SSL-GM to self-supervised graph learning methods. **DGI** (Veličković et al., 2019) and **MVGRL** (Hassani & Khasahmadi, 2020) conduct contrastive learning between graph patches and graph summaries to integrate knowledge into node representations. **GRACE** (Zhu et al., 2020) and subsequent **GCA** (Zhu et al., 2021) perform contrast between nodes in two corrupted views to acquire augmentation-invariant representations. **BGRL** (Thakoor et al., 2022) utilizes predictive objective for node-level contrastive learning to achieve efficient training. For graph-level tasks, we explore traditional graph kernels for classification, including **WL kernel** (Shervashidze et al., 2011) and **DGK** (Yanardag & Vishwanathan, 2015). Furthermore, we include contrastive learning approaches, such as **graph2vec** (Narayanan et al., 2017), **MVGRL** (Hassani & Khasahmadi, 2020), **InfoGraph** (Sun et al., 2020), **GraphCL** (You et al., 2020), and **JOAO** (You et al., 2021), which conduct contrastive learning between representations of two augmented graphs.

**MLPs on Graphs.** In node classification, we employ basic **MLP** that considers only node content as baseline. Furthermore, we incorporate **GraphMLP** (Hu et al., 2021) that trains an MLP by emphasizing consistency between target nodes and their direct neighborhoods. We exclude the following works (Dong et al., 2022; Liu et al., 2022) as baselines since they are high-order versions of GraphMLP. To achieve this, we slightly modify the original GraphMLP to enable the ability in learning high-order information, and search the number of layers within  $\{1, 2, 3\}$ . **GLNN** (Zhang et al., 2022) employs knowledge distillation to transfer knowledge from GNNs to MLPs, **GENN** leverages positional encoding to acquire structural knowledge, while **NOSMOG** (Tian et al., 2023) jointly integrates positional information and robust training strategies based on GLNN. Note that the public code of GENN is not available, thus we implement GENN based on the code of NOSMOG. In graph classification, we employ a pooling function to generate graph-level representations for training an MLP. For other baselines, they cannot be readily applied to graph-level tasks. Therefore, we implement an MLP with graph-level knowledge distillation as another baseline to access the role of knowledge distillation on graph classification.

## B.3 HYPER-PARAMETER SETTING

Hyper-parameters	Node Classification									
	Cora	Citeseer	PubMed	Amazon-CS	Amazon-Photo	Co-CS	Co-Phys	Wiki-CS	Flickr	Arxiv
Epochs	1000	1000	1000	1000	1000	2000	1000	2000	2000	5000
Optimizer	AdamW used for all datasets									
Learning Rate	1e-3	5e-4	5e-4	1e-3	1e-3	1e-4	1e-3	5e-4	1e-3	1e-3
Weight Decay	0	5e-5	1e-5	-	1e-4	-	1e-4	1e-5	5e-4	-
Activation	ReLU used for all datasets									
Hidden Dimension	512	512	512	512	512	512	512	512	1024	1024
Normalization	Batchnorm used for all datasets									
# Encoder Layers	2	2	2	3	2	2	2	2	2	8
# Aggregator Layers	2	3	3	2	1	1	1	2	3	3
Feature Mask Ratio	0.50	0.75	0.25	0.25	0.25	0.50	0.75	0.00	0.25	0.00
Edge Mask Ratio	0.25	0.50	0.25	0.25	0.50	0.75	0.50	0.25	0.50	0.25

Table 7: Hyper-parameters used for SSL-GM for node-level task.

We perform hyper-parameter tuning for each approach using a grid search strategy. Specifically, we set the number of epochs to 1,000, the hidden dimension to 512, and employ ReLU as the activation function. We explore various learning rates  $\{5e-4, 1e-4, 5e-4, 1e-3, 5e-3, 1e-2\}$ , weight decay values  $\{5e-5, 1e-5, 5e-3, 1e-4, 0\}$ , and the number of layers  $\{1, 2, 3\}$ . In self-supervised learning methods, we employ a 2-layer GCN (Kipf & Welling, 2017) as the encoder for node-level tasks. For graph-level tasks, we utilize a 5-layer GIN (Xu et al., 2019) model concatenated with a readout function, which is selected from  $\{\text{MEAN}, \text{SUM}, \text{MAX}\}$ . Subsequently, we assess the quality of the acquired representations by training a Logistic regression function on downstream

Hyper-parameters	Graph Classification						
	IMDB-B	IMDB-M	COLLAB	PTC-MR	MUTAG	DD	PROTEINS
Epochs	200	100	30	100	100	100	500
Optimizer	AdamW used for all datasets						
Learning Rate	1e-2	1e-2	5e-4	1e-2	1e-2	1e-3	1e-3
Weight Decay	0 used for all datasets						
Activation	PReLU used for all datasets						
Batch Size	64	128	32	64	64	32	64
Raw Feature	N	N	N	Y	N	Y	N
Deg4Feature	Y	Y	Y	N	Y	N	Y
Pooling	MEAN	MEAN	MEAN	SUM	SUM	MEAN	SUM
Hidden Dimension	512 used for all datasets						
Normalization	Batchnorm used for all datasets						
# Encoder Layers	2 used for all datasets						
# Aggregator Layers	2	2	2	2	1	2	1
Feature Mask Ratio	0.50	0.25	0.75	0.25	0.5	0.00	0.00
Edge Mask Ratio	0.75	0.50	0.75	0.00	0.25	0.00	0.50

Table 8: Hyper-parameters used for SSL-GM for graph-level task.

tasks (Zhu et al., 2020). For other settings, we follow the settings reported in the original papers. Regarding SSL-GM, we provide a comprehensive overview of the hyper-parameter settings for both node-level and graph-level tasks in Tables 7 and 8, respectively.

#### B.4 DETAILED EVALUATION SETTING

**Transductive Setting.** We consider a graph  $\mathcal{G} = (\mathcal{V}, E)$  in which all nodes are visible during the training stage. We partition the nodes into three non-overlapping sets:  $\mathcal{V} = \mathcal{V}_{train} \sqcup \mathcal{V}_{val} \sqcup \mathcal{V}_{test}$ . In the supervised setting, we train the encoder using  $\mathcal{V}_{train}$  and evaluate its performance on  $\mathcal{V}_{val}$  and  $\mathcal{V}_{test}$ . In the self-supervised setting, we train the encoder on  $\mathcal{V}$  and utilize  $\mathcal{V}_{trans}$  for training the downstream head, while using  $\mathcal{V}_{val}$  and  $\mathcal{V}_{test}$  for evaluation. We perform the split 10 times to evaluate the quality of the learned representations to alleviate the impact of randomness.

**Inductive (Production) Setting.** In the production setting, we partition a graph  $\mathcal{G} = (\mathcal{V}, E)$  into transductive set  $\mathcal{G}^T$  and inductive set  $\mathcal{G}^I$ , where  $\mathcal{G} = \mathcal{G}^T \sqcup \mathcal{G}^I$  and  $\emptyset = \mathcal{G}^T \cap \mathcal{G}^I$ .  $\mathcal{G}^T = (\mathcal{V}^T, E^T)$ , containing 80% of the nodes, is used for training, while  $\mathcal{G}^I = (\mathcal{V}^I, E^I)$ , containing the remaining 20% of the nodes, remains unseen during training. We further partition the nodes in  $\mathcal{G}^T$  into non-overlapping training, validation, and testing sets, denoted as  $\mathcal{V}^T = \mathcal{V}_{train}^T \sqcup \mathcal{V}_{valid}^T \sqcup \mathcal{V}_{test}^T$ . The model is trained on  $\mathcal{G}^T$ , and we report the transductive results on  $\mathcal{V}_{test}^T$  and inductive results on  $\mathcal{V}^I$ . We interpolate the results of these two settings to represent production results. Our approach differs from that of Zhang et al. (2022) and Tian et al. (2023). We treat nodes in the inductive set as disconnected from nodes in the transductive set, even during inference, creating a more challenging out-of-distribution setting.

**Cold-start Setting.** The cold-start setting is closely similar to inductive (production) setting, but we assume the nodes in the inductive set are isolated,  $\mathcal{G}^C = (\mathcal{V}^I, \emptyset)$ . This is a more challenging yet practical setting, as new agents often emerge independently in real-world systems. We follow the production setting to train the model on transductive set  $\mathcal{G}^T$  but only report the model performance on cold-start set  $\mathcal{G}^C$  to assess cold-start performance. For models that rely on positional encoding, such as GENN and NOSMOG, we set the positional embeddings as zero vectors.

## C EMPIRICALLY RUNNING TIME AND MEMORY USAGE

Table 9 presents a comparison of the running time and memory usage between our SSL-GM and other baseline methods, namely GAT (Veličković et al., 2018), GRACE (Zhu et al., 2020), and BGRL (Thakoor et al., 2022). Despite our primary focus on accelerating inference speed, we have observed that SSL-GM outperforms these approaches in both training time and memory utilization. In particular, GAT, which employs 4 attention heads, imposes a substantial computational burden during training due to attention score computation, resulting in significant memory consumption. Considering self-supervised methods, GRACE utilizes the InfoNCE loss for model training, involving computation of node pair similarities. This operation results in considerable time and memory overhead. In comparison to GRACE, our SSL-GM demonstrates improvements in terms of memory

usage ( $3.8 \sim 6.8\times$ ) and training time efficiency ( $4.8 \sim 8.3\times$ ). BGRL employs bootstrap loss (Grill et al., 2020) to predict node representations from different views, remarkably enhancing training and memory efficiency. However, our SSL-GM remains more efficient than BGRL.

Additionally, we conducted comparison between SSL-GM and other acceleration techniques on Flickr and Arxiv datasets, as summarized in Table 10 and 11. Note that the inference time for self-supervised approaches in this paper aligns with that of SAGE (Hamilton et al., 2017). Therefore, we establish SAGE as our baseline and apply acceleration techniques based on that, including quantization (QSAGE), pruning (PSAGE), and neighbor sampling (Neighbor Sample), to facilitate the inference. Furthermore, we include an evaluation of the inference time of SGC (Wu et al., 2019), which comprises an MLP and a one-layer message passing. It is evident that even the most efficient methods achieve only modest acceleration ( $3.2 \sim 4.0\times$ ), and SGC, which employs only a single-layer aggregation, achieves a mere acceleration ( $1.1 \sim 1.2\times$ ). This observation demonstrates that the aggregation process leads to significant time consumption during inference. In contrast, our SSL-GM achieves remarkable inference acceleration by disregarding neighborhood dependency, which is faster than SAGE ( $89.7 \sim 125.9\times$ ). Compared to other methods employing MLPs as encoders like GLNN and NOSMOG, we do not observe significant distinctions in terms of inference speed. However, methods like NOSMOG (Tian et al., 2023) and GENN (Wang et al., 2023) utilize additional positional embeddings, introducing significant time consumption in learning positional embedding. For example, NOSMOG remains fast inference under 5 milliseconds, while the positional encoding takes more than 5 seconds on the Arxiv dataset, even with a minimal epoch setting of 1. Considering the time consumption on encoding positional embeddings, the overall inference consumption of GENN and NOSMOG far surpasses that of SAGE.

Dataset	Amazon-CS		Amazon-Photo		Coauthor-CS		Coauthor-Phys		Wiki-CS	
	Memory	Training Time	Memory	Training Time	Memory	Training Time	Memory	Training Time	Memory	Training Time
GAT	5239 MB	73.8 (s)	2571 MB	41.9 (s)	2539 MB	60.4 (s)	13199 MB	265.2 (s)	4568 MB	74.4 (s)
GRACE	8142 MB	349.5 (s)	2755 MB	138.4 (s)	11643 MB	261.4 (s)	16294 MB	573.2 (s)	5966 MB	290.9 (s)
BGRL	2196 MB	96.8 (s)	1088 MB	64.1 (s)	2513 MB	129.9 (s)	5556 MB	273.8 (s)	1899 MB	108.8 (s)
<b>SSL-GM</b>	<b>1969 MB</b>	<b>53.4 (s)</b>	<b>694 MB</b>	<b>27.0 (s)</b>	<b>1716 MB</b>	<b>54.8 (s)</b>	<b>3920 MB</b>	<b>110.7 (s)</b>	<b>1590 MB</b>	<b>35.5 (s)</b>

Table 9: Computational requirements of different baseline methods on a set of standard benchmark graphs. The experiments are performed on a 24GB Nvidia GeForce RTX 3090.

Datasets	SAGE	BGRL	SGC	APPNP	QSAGE	PSAGE	Neighbor Sample	<b>SSL-GM</b>
Flickr	Time (ms)	80.7	80.7 (1.00 $\times$ )	76.9 (1.05 $\times$ )	78.1 (1.03 $\times$ )	70.6 (1.14 $\times$ )	67.4 (1.20 $\times$ )	25.5 (3.16 $\times$ )
	Acc (%)	47.17	49.12	47.35	47.53	47.22	47.25	47.01
Arxiv	Time (ms)	314.7	314.7 (1.00 $\times$ )	265.9 (1.18 $\times$ )	284.1 (1.11 $\times$ )	289.5 (1.09 $\times$ )	297.5 (1.06 $\times$ )	78.3 (4.02 $\times$ )
	Acc (%)	68.52	69.29	68.93	69.10	68.48	68.55	68.35

Table 10: The inference time and accuracy of different acceleration methods.

Dataset	Models	Trans		Ind		Cold-start	
		Time (ms)	Acc (%)	Time (ms)	Acc (%)	Time (ms)	Acc (%)
Pubmed	SAGE	73	85.94	15	85.04	15	77.98
	SGC	64	85.28	14	85.22	13	76.10
	NOSMOG	5	86.18	3	83.84	3	81.48
	<b>SSL-GM</b>	<b>3</b>	<b>86.99</b>	<b>3</b>	<b>86.47</b>	<b>3</b>	<b>86.44</b>
Amazon-CS	SAGE	103	88.88	31	87.24	25	61.01
	SGC	89	<b>89.31</b>	26	87.12	24	63.08
	NOSMOG	5	87.64	4	86.61	4	81.95
	<b>SSL-GM</b>	<b>4</b>	88.46	<b>3</b>	<b>87.65</b>	<b>3</b>	<b>87.58</b>
Arxiv	SAGE	485	<b>72.05</b>	315	68.52	305	43.47
	SGC	410	69.95	266	68.93	250	42.08
	NOSMOG	6	70.84	4	69.10	4	61.64
	<b>SSL-GM</b>	<b>4</b>	71.12	<b>3</b>	<b>70.23</b>	<b>3</b>	<b>66.13</b>

Table 11: Inference acceleration across different settings.

		IMDB-B	IMDB-M	COLLAB	PTC-MR	MUTAG	DD	PROTEINS
Supervised	GIN	75.10 $\pm$ 5.10	52.30 $\pm$ 2.80	80.20 $\pm$ 1.90	64.60 $\pm$ 1.70	89.40 $\pm$ 5.60	74.88 $\pm$ 3.12	76.20 $\pm$ 2.80
Graph Kernel	WL	72.30 $\pm$ 3.44	46.95 $\pm$ 0.46	-	57.97 $\pm$ 0.49	80.72 $\pm$ 3.00	-	72.92 $\pm$ 0.56
	DGK	66.96 $\pm$ 0.56	44.55 $\pm$ 0.52	-	60.08 $\pm$ 2.55	87.44 $\pm$ 2.72	-	73.30 $\pm$ 0.82
GCL	graph2vec	71.10 $\pm$ 0.54	50.44 $\pm$ 0.87	-	60.17 $\pm$ 6.86	83.15 $\pm$ 9.25	-	73.30 $\pm$ 2.05
	MVGRL	71.84 $\pm$ 0.78	50.84 $\pm$ 0.92	73.10 $\pm$ 0.56	-	89.24 $\pm$ 1.31	75.20 $\pm$ 0.55	74.02 $\pm$ 0.32
	InfoGraph	73.03 $\pm$ 0.87	49.69 $\pm$ 0.53	70.65 $\pm$ 1.13	61.65 $\pm$ 1.43	89.01 $\pm$ 1.13	72.85 $\pm$ 1.78	74.44 $\pm$ 0.31
	GraphCL	71.14 $\pm$ 0.44	48.58 $\pm$ 0.67	71.36 $\pm$ 1.15	-	86.80 $\pm$ 1.34	78.62 $\pm$ 0.40	74.39 $\pm$ 0.45
	JOAO	70.21 $\pm$ 3.08	49.20 $\pm$ 0.77	69.50 $\pm$ 0.36	-	87.35 $\pm$ 1.02	-	74.55 $\pm$ 0.41
MLP	MLP*	49.50 $\pm$ 1.66	33.11 $\pm$ 1.59	51.90 $\pm$ 0.95	54.39 $\pm$ 1.41	67.22 $\pm$ 0.99	58.56 $\pm$ 1.40	59.20 $\pm$ 1.00
	MLP + KD*	72.85 $\pm$ 1.04	48.14 $\pm$ 0.52	75.38 $\pm$ 1.53	59.38 $\pm$ 1.38	87.44 $\pm$ 0.67	73.59 $\pm$ 1.69	73.54 $\pm$ 1.78
<b>SSL-GM</b>		<b>74.06<math>\pm</math>0.22</b>	<b>51.41<math>\pm</math>0.52</b>	<b>81.04<math>\pm</math>0.11</b>	<b>60.28<math>\pm</math>1.07</b>	<b>87.67<math>\pm</math>0.24</b>	<b>78.44<math>\pm</math>0.47</b>	<b>75.31<math>\pm</math>0.13</b>
$\Delta_{\text{GraphCL}}$		$\uparrow$ 4.10%	$\uparrow$ 5.83%	$\uparrow$ 13.57%	-	$\uparrow$ 1.00%	$\downarrow$ 0.23%	$\uparrow$ 1.24%
$\Delta_{\text{MLP}}$		$\uparrow$ 49.62%	$\uparrow$ 55.27%	$\uparrow$ 56.15%	$\uparrow$ 10.83%	$\uparrow$ 30.42%	$\uparrow$ 33.95%	$\uparrow$ 27.21%
$\Delta_{\text{MLP+KD}}$		$\uparrow$ 1.66%	$\uparrow$ 6.79%	$\uparrow$ 7.51%	$\uparrow$ 1.52%	$\uparrow$ 0.26%	$\uparrow$ 6.59%	$\uparrow$ 2.41%

The reported results of baselines are from previous papers if available (You et al., 2020; 2021; Hou et al., 2022). \* indicates the results are from our implementation.

Table 12: Graph classification accuracy (%).  $\Delta_{\text{MLP}}$ ,  $\Delta_{\text{GraphCL}}$ ,  $\Delta_{\text{MLP+KD}}$  represents the performance gap (%) between our methods and GraphCL, MLP, and knowledge distillation-enhanced MLP, where green indicates the improvement over 4% and red indicates the degradation.

## D ADDITIONAL EXPERIMENTAL RESULTS

### D.1 EFFECTIVENESS OF SSL-GM ON GRAPH CLASSIFICATION

In this section, we present the experimental results of SSL-GM alongside state-of-the-art baselines for graph classification. The datasets encompass a variety of graph types, including biomolecular graphs and social networks (as described in Appendix B.1), and the baselines comprise supervised GNNs, traditional graph kernels, and self-supervised graph learning approaches (as detailed in Appendix B.2). Note that existing MLP-based methods, such as GLNN and NOSMOG, are specifically designed for node-level tasks, which hinder their direct applicability to graph classification tasks. To assess the impact of knowledge distillation on graph-level tasks, we also implement a combination of MLP and KD, which is similar to GLNN but conducts knowledge distillation at the graph level. Due to the generalizability inherent in self-supervised learning (Sun et al., 2023), SSL-GM can be readily extended to various downstream tasks. Our results demonstrate that SSL-GM outperforms other MLP-based methods, achieving the best or sub-best performance on 6 out of 7 baselines. In contrast, MLP-based methods generally underperform when compared to self-supervised GCL methods and even traditional graph kernel methods. Notably, SSL-GM surpasses other baselines, even on the large-scale dataset COLLAB, highlighting the potential of MLP-based methods in graph-level tasks. Our findings expand the usability of MLPs from node-level tasks to graph-level tasks.

### D.2 COMPREHENSIVE PREDICTION RESULTS UNDER INDUCTIVE (PRODUCTION) SETTING

Table 13 presents the comprehensive experimental results of our inductive (production) setting. We report the performance on both transductive and inductive sets, along with the interpolated production results. Note that the results of GLNN (Zhang et al., 2022) and NOSMOG (Tian et al., 2023) reported in this paper differ from their respective original papers. This discrepancy arises because our experimental setting presents a more challenging task where the inductive set is disconnected from the transductive set in inference. In this table, we present the performance of six baseline methods, which encompass supervised GNNs, self-supervised approaches, and MLP-based techniques. We can observe that our SSL-GM attains state-of-the-art performance in the majority of settings. Among GNN methods, we note that SAGE generally outperforms BGRL in transductive settings, but underperforms it in inductive settings. The robustness of BGRL stems from the augmentation, which aids in learning augmentation-invariant representations, enabling BGRL to work on inductive sets, even though the distributions of these two sets are potentially different. Regarding MLP-based methods, we observe that NOSMOG outperforms GLNN in transductive settings, particularly on large-scale graphs, while GLNN significantly outperforms NOSMOG in inductive settings. We assume that the learned positional embeddings of NOSMOG on the inductive set differ from those on the transductive set, rendering them untrustworthy and less meaningful. For our SSL-GM, the

model maximizes consistency of MLPs and GNNs in the representation space, thereby preserving more fine-grained structural knowledge, which ensures the model generalization.

	Setting	Cora	Citeseer	PubMed	Amazon-CS	Amazon-Photo	Co-CS	Co-Phys	Wiki-CS	Flickr	Arxiv
SAGE	<i>prod</i>	77.51 $\pm$ 1.77	68.40 $\pm$ 1.61	85.04 $\pm$ 0.44	87.24 $\pm$ 0.43	93.20 $\pm$ 0.45	92.88 $\pm$ 0.40	95.74 $\pm$ 0.12	<b>79.26</b> $\pm$ 0.65	47.17 $\pm$ 0.73	68.52 $\pm$ 0.56
	<i>trans</i>	79.46 $\pm$ 1.49	68.73 $\pm$ 1.37	85.57 $\pm$ 0.30	<b>87.98</b> $\pm$ 0.30	93.70 $\pm$ 0.42	93.13 $\pm$ 0.33	95.77 $\pm$ 0.04	80.01 $\pm$ 0.41	48.15 $\pm$ 0.63	<b>71.79</b> $\pm$ 0.50
	<i>ind</i>	69.70 $\pm$ 2.89	67.11 $\pm$ 2.57	82.90 $\pm$ 0.98	84.45 $\pm$ 0.94	91.18 $\pm$ 0.56	91.87 $\pm$ 0.68	95.63 $\pm$ 0.05	76.27 $\pm$ 1.63	43.25 $\pm$ 1.14	55.45 $\pm$ 0.78
BGRL	<i>prod</i>	77.73 $\pm$ 1.07	64.33 $\pm$ 1.56	83.97 $\pm$ 0.48	<b>87.33</b> $\pm$ 0.48	91.47 $\pm$ 0.62	91.26 $\pm$ 0.35	94.38 $\pm$ 0.29	76.25 $\pm$ 1.09	<b>49.12</b> $\pm$ 0.31	<b>69.29</b> $\pm$ 0.38
	<i>trans</i>	77.32 $\pm$ 0.90	64.15 $\pm$ 1.40	83.97 $\pm$ 0.34	87.27 $\pm$ 0.42	91.47 $\pm$ 0.51	91.31 $\pm$ 0.33	94.40 $\pm$ 0.25	76.32 $\pm$ 0.97	<b>49.09</b> $\pm$ 0.24	70.36 $\pm$ 0.35
	<i>ind</i>	<b>79.38</b> $\pm$ 1.74	65.03 $\pm$ 2.19	83.98 $\pm$ 1.02	<b>87.59</b> $\pm$ 0.75	<b>91.46</b> $\pm$ 1.05	91.09 $\pm$ 0.45	94.33 $\pm$ 0.46	75.96 $\pm$ 1.57	<b>49.26</b> $\pm$ 0.60	<b>65.03</b> $\pm$ 0.50
MLP	<i>prod</i>	63.76 $\pm$ 1.65	63.98 $\pm$ 1.22	80.91 $\pm$ 0.45	81.00 $\pm$ 0.54	87.73 $\pm$ 0.88	91.68 $\pm$ 0.59	95.18 $\pm$ 0.13	75.08 $\pm$ 0.71	46.14 $\pm$ 0.22	55.89 $\pm$ 0.51
	<i>trans</i>	63.66 $\pm$ 1.53	63.86 $\pm$ 1.09	80.92 $\pm$ 0.38	81.05 $\pm$ 0.45	87.69 $\pm$ 0.86	91.66 $\pm$ 0.54	95.18 $\pm$ 0.12	75.12 $\pm$ 0.43	46.16 $\pm$ 0.15	55.89 $\pm$ 0.46
	<i>ind</i>	64.15 $\pm$ 2.11	64.43 $\pm$ 1.76	80.90 $\pm$ 0.72	80.80 $\pm$ 0.91	87.88 $\pm$ 0.96	91.78 $\pm$ 0.81	95.16 $\pm$ 0.18	74.94 $\pm$ 1.81	46.09 $\pm$ 0.50	55.91 $\pm$ 0.69
GLNN	<i>prod</i>	<b>78.34</b> $\pm$ 1.04	<b>69.61</b> $\pm$ 1.13	<b>85.44</b> $\pm$ 0.48	87.04 $\pm$ 0.50	<b>93.28</b> $\pm$ 0.43	<b>93.72</b> $\pm$ 0.35	95.76 $\pm$ 0.09	<b>78.39</b> $\pm$ 0.54	46.11 $\pm$ 0.27	63.53 $\pm$ 0.48
	<i>trans</i>	79.93 $\pm$ 0.87	<b>69.73</b> $\pm$ 0.77	85.70 $\pm$ 0.38	87.80 $\pm$ 0.45	93.84 $\pm$ 0.42	<b>93.82</b> $\pm$ 0.32	95.78 $\pm$ 0.04	78.58 $\pm$ 0.32	46.13 $\pm$ 0.22	64.27 $\pm$ 0.46
	<i>ind</i>	71.96 $\pm$ 1.68	<b>69.14</b> $\pm$ 2.58	<b>84.42</b> $\pm$ 0.87	83.98 $\pm$ 0.70	91.05 $\pm$ 0.49	<b>93.34</b> $\pm$ 0.47	<b>95.70</b> $\pm$ 0.09	<b>77.64</b> $\pm$ 1.42	46.05 $\pm$ 0.43	60.55 $\pm$ 0.55
GENN	<i>prod</i>	77.83 $\pm$ 1.57	67.30 $\pm$ 1.48	84.34 $\pm$ 0.47	85.75 $\pm$ 1.20	92.09 $\pm$ 0.96	93.57 $\pm$ 0.37	95.67 $\pm$ 0.06	78.27 $\pm$ 1.01	45.56 $\pm$ 0.51	68.52 $\pm$ 0.54
	<i>trans</i>	<b>80.27</b> $\pm$ 1.41	67.86 $\pm$ 1.16	<b>85.81</b> $\pm$ 0.38	87.42 $\pm$ 1.04	93.35 $\pm$ 0.60	93.79 $\pm$ 0.35	95.78 $\pm$ 0.05	<b>80.31</b> $\pm$ 0.85	45.68 $\pm$ 0.45	70.01 $\pm$ 0.50
	<i>ind</i>	68.09 $\pm$ 2.21	65.07 $\pm$ 2.78	78.44 $\pm$ 0.84	79.09 $\pm$ 1.84	87.05 $\pm$ 2.42	92.68 $\pm$ 0.45	95.23 $\pm$ 0.08	70.13 $\pm$ 1.65	45.08 $\pm$ 0.74	62.58 $\pm$ 0.69
NOSMOG	<i>prod</i>	77.83 $\pm$ 1.94	68.58 $\pm$ 1.41	83.84 $\pm$ 0.45	86.61 $\pm$ 1.22	92.52 $\pm$ 0.68	93.45 $\pm$ 0.44	<b>95.78</b> $\pm$ 0.10	78.35 $\pm$ 0.70	46.05 $\pm$ 0.55	69.10 $\pm$ 0.80
	<i>trans</i>	<b>80.27</b> $\pm$ 1.69	68.95 $\pm$ 1.24	85.43 $\pm$ 0.37	<b>88.30</b> $\pm$ 1.14	<b>93.88</b> $\pm$ 0.47	93.68 $\pm$ 0.38	<b>95.85</b> $\pm$ 0.10	<b>80.35</b> $\pm$ 0.58	46.24 $\pm$ 0.51	70.50 $\pm$ 0.79
	<i>ind</i>	68.11 $\pm$ 2.95	67.07 $\pm$ 2.10	77.44 $\pm$ 0.80	79.83 $\pm$ 1.52	87.08 $\pm$ 1.52	92.55 $\pm$ 0.69	95.52 $\pm$ 0.10	70.36 $\pm$ 1.18	45.27 $\pm$ 0.72	63.49 $\pm$ 0.83
SSL-GM	<i>prod</i>	<b>81.37</b> $\pm$ 1.20	<b>72.33</b> $\pm$ 0.90	<b>86.47</b> $\pm$ 0.28	<b>87.65</b> $\pm$ 0.40	<b>93.87</b> $\pm$ 0.32	<b>94.63</b> $\pm$ 0.16	<b>96.04</b> $\pm$ 0.12	<b>79.26</b> $\pm$ 0.83	<b>49.27</b> $\pm$ 0.18	<b>70.23</b> $\pm$ 0.47
	<i>trans</i>	<b>81.60</b> $\pm$ 0.96	<b>72.21</b> $\pm$ 0.73	<b>86.48</b> $\pm$ 0.23	87.67 $\pm$ 0.25	<b>93.86</b> $\pm$ 0.26	<b>94.66</b> $\pm$ 0.16	<b>96.06</b> $\pm$ 0.09	79.46 $\pm$ 0.66	<b>49.23</b> $\pm$ 0.11	<b>71.25</b> $\pm$ 0.33
	<i>ind</i>	<b>80.48</b> $\pm$ 2.15	<b>72.81</b> $\pm$ 1.61	<b>86.44</b> $\pm$ 0.51	<b>87.58</b> $\pm$ 0.99	<b>93.91</b> $\pm$ 0.58	<b>94.51</b> $\pm$ 0.15	<b>95.97</b> $\pm$ 0.24	<b>78.46</b> $\pm$ 1.48	<b>49.41</b> $\pm$ 0.46	<b>66.13</b> $\pm$ 1.05

Table 13: Node classification accuracy (%) under inductive (production) scenario for both transductive and inductive settings. *ind* represents the accuracy on  $\mathcal{V}^I$ , *trans* represents the accuracy on  $\mathcal{V}_{test}^T$ , and *prod* is the interpolated accuracy of both *ind* and *trans*.

### D.3 EMPLOYING STANDARD GNNs WILL LEAD TO MODEL COLLAPSE

In this section, we attempt to answer why we approximate the representations of GNNs by a non-parametric aggregator instead of standard GNNs. We empirically demonstrate that the use of GNNs will lead to model collapse, which is a common issue in contrastive learning. To solve the issue, existing works in computer vision employ moving average network (Grill et al., 2020) to preserve the consistency between representations learned from two networks. Inspired by this, we propose a non-parametric aggregator to approximate the GNN representations based on MLP representations, thereby preserving the inherent consistency between these two representations. In Figure 5, we compare the accuracy and loss curves between two versions of SSL-GM: one utilizes the non-parametric aggregator (MLP + PROP.), and the other directly employs GNNs. The experiments are conducted on five benchmark datasets in a transductive setting. We observe that MLP + PROP. achieves superior accuracy and higher losses than the GNN version. Specifically, we suppose the model collapse stems from the lack of consistency between the outputs of GNNs and MLPs in the representation space, which is a significant problem in contrastive learning (He et al., 2020). In our SSL-GM, we utilize the non-parametric aggregator to approximate GNN representations based on MLP representations, naturally preserving the consistency between MLPs and GNNs in the representation space. Our observation aligns with the findings in (He et al., 2020; Thakoor et al., 2022).

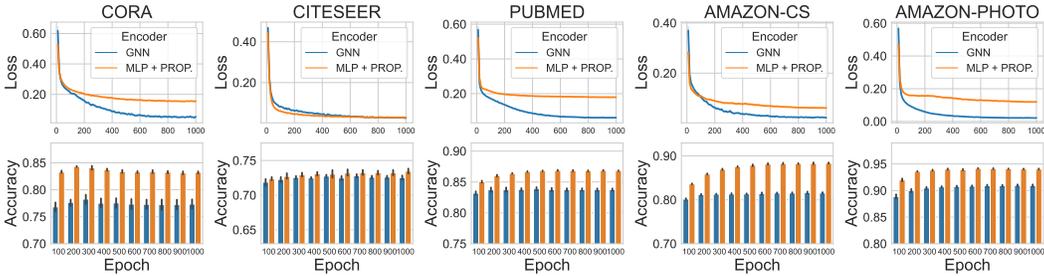


Figure 5: The training curves and model accuracies on five benchmark datasets. We compare the method applied in our SSL-GM (MLP + PROP.) with standard GNNs in learning node representations. We observe that employing GNNs will lead to model collapse.

D.4 HOW AUGMENTATIONS IMPROVE THE MODEL CAPABILITY?

In this section, we evaluate the influence of augmentations on model performance. We conduct an ablation study to assess the model performance without specific augmentations, demonstrating both semantic and structural augmentations enhance model performance. Table 14 presents a comprehensive ablation study of edge masking and node feature masking for node classification on ten benchmark datasets within a transductive setting. We observe that these two types of augmentations significantly enhance model performance by improving different aspects of the datasets. Furthermore, the combination of these two techniques further enhance the performance of SSL-GM, indicating that our model can benefit from both augmentations simultaneously.

Additionally, we conduct a detailed analysis of how augmentations impact model performance. Figure 6 illustrates the model performance at different augmentation probabilities on Cora, Citeseer, PubMed, Amazon-CS, and Amazon-Photo datasets within a transductive setting. The augmentation ratio is searched among {0.0, 0.25, 0.5, 0.75}. These figures enable us to gain insight into the specific effects of augmentations on model performance.

Feature Masking	Edge Masking	Cora	Citeseer	PubMed	Amazon-CS	Amazon-Photo	Co-CS	Co-Phys	Wiki-CS	Flickr	Arxiv
-	-	82.10 <sub>±0.45</sub>	71.83 <sub>±0.43</sub>	86.89 <sub>±0.13</sub>	87.12 <sub>±0.15</sub>	93.52 <sub>±0.20</sub>	93.10 <sub>±0.05</sub>	94.56 <sub>±0.06</sub>	80.98 <sub>±0.13</sub>	48.21 <sub>±0.10</sub>	70.58 <sub>±0.20</sub>
✓	-	84.78 <sub>±0.25</sub>	73.00 <sub>±0.63</sub>	86.98 <sub>±0.09</sub>	88.27 <sub>±0.18</sub>	94.19 <sub>±0.14</sub>	94.50 <sub>±0.10</sub>	96.12 <sub>±0.06</sub>	81.03 <sub>±0.11</sub>	49.55 <sub>±0.11</sub>	70.03 <sub>±0.23</sub>
-	✓	82.33 <sub>±0.61</sub>	71.78 <sub>±0.77</sub>	86.98 <sub>±0.13</sub>	87.35 <sub>±0.29</sub>	93.69 <sub>±0.07</sub>	94.35 <sub>±0.08</sub>	95.88 <sub>±0.06</sub>	81.04 <sub>±0.22</sub>	49.33 <sub>±0.07</sub>	<b>71.12</b> <sub>±0.10</sub>
✓	✓	<b>84.60</b> <sub>±0.24</sub>	<b>73.52</b> <sub>±0.53</sub>	<b>86.99</b> <sub>±0.09</sub>	<b>88.46</b> <sub>±0.16</sub>	<b>94.28</b> <sub>±0.08</sub>	<b>94.87</b> <sub>±0.07</sub>	<b>96.17</b> <sub>±0.03</sub>	<b>81.21</b> <sub>±0.13</sub>	<b>49.85</b> <sub>±0.09</sub>	<b>71.12</b> <sub>±0.10</sub>

Table 14: Ablation study of augmentation methods used in SSL-GM.

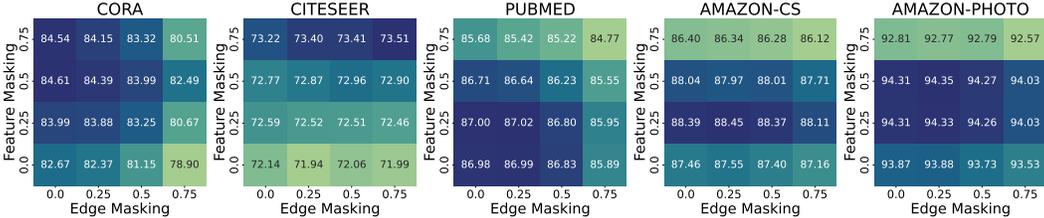


Figure 6: Node classification accuracy on transductive setting with different augmentation ratios.

D.5 DOES THE RECONSTRUCTION TERM PREVENT REPRESENTATION SHIFTS?

In this section, we evaluate the role of the reconstruction term in SSL-GM. It is important to note that the reconstruction term serves to mitigate representation shifts caused by augmentation. The reconstruction term enables the outputs of GNNs to remain invariant to node features, thereby preserving more localized knowledge. Intuitively, if we do not apply augmentations, the incorporation of the reconstruction term will hinder the process of contrastive learning. To evaluate this, we assess the impact of the reconstruction term with and without applying augmentations, shown in Table 15. Intriguingly, the model performance demonstrates that the reconstruction term is effective only when augmentation is utilized. Without augmentation, the model performance is adversely affected. When both augmentation and reconstruction are applied, our SSL-GM achieves the highest performance, demonstrating the importance of employing reconstruction to prevent representation shifts brought by augmentation.

Aug.	Rec.	Cora	Citeseer	PubMed	Amazon-CS	Amazon-Photo	Co-CS	Co-Phys	Wiki-CS	Flickr	Arxiv
-	-	<b>82.80</b> <sub>±0.65</sub>	<b>72.03</b> <sub>±0.83</sub>	<b>86.89</b> <sub>±0.07</sub>	<b>87.42</b> <sub>±0.17</sub>	<b>93.87</b> <sub>±0.04</sub>	<b>93.12</b> <sub>±0.03</sub>	94.55 <sub>±0.04</sub>	<b>81.06</b> <sub>±0.13</sub>	<b>48.61</b> <sub>±0.10</sub>	70.38 <sub>±0.32</sub>
-	✓	82.10 <sub>±0.45</sub>	71.83 <sub>±0.43</sub>	<b>86.89</b> <sub>±0.13</sub>	87.12 <sub>±0.15</sub>	93.52 <sub>±0.20</sub>	93.10 <sub>±0.05</sub>	<b>94.56</b> <sub>±0.06</sub>	80.98 <sub>±0.13</sub>	48.21 <sub>±0.10</sub>	<b>70.58</b> <sub>±0.20</sub>
✓	✓	84.60 <sub>±0.24</sub>	73.52 <sub>±0.53</sub>	86.99 <sub>±0.09</sub>	88.46 <sub>±0.16</sub>	94.28 <sub>±0.08</sub>	94.87 <sub>±0.07</sub>	96.17 <sub>±0.03</sub>	81.21 <sub>±0.13</sub>	49.85 <sub>±0.09</sub>	71.12 <sub>±0.10</sub>

Table 15: Reconstruction term leads to performance drop without augmentation.

## D.6 INFLUENCE OF AGGREGATOR TYPES AND AGGREGATION LAYERS

In this section, we examine the role of the aggregator of SSL-GM. Unlike SGC or APPNP, which employ high-order adjacent matrices to guide the message passing, we utilize a GNN-like architecture to guide aggregation, preserving non-linearity between layers. While various aggregation methods can be employed, we adopt the normalized Laplacian matrix, which is similar to GCN (Kipf & Welling, 2017), to aggregate high-order representations, as formulated as:

$$\mathbf{H}^{(l)} = \phi^{(l)}(\mathbf{A}, \mathbf{H}^{(l-1)}) = \sigma(\hat{\mathbf{D}}^{-\frac{1}{2}} \hat{\mathbf{A}} \hat{\mathbf{D}}^{-\frac{1}{2}} \mathbf{H}^{(l-1)}) \quad (18)$$

where  $\mathbf{H}^{(l)}$  represents the high-order representation at  $l$ -th aggregator layer,  $\phi^{(l)}(\cdot)$  denotes the  $l$ -th aggregator,  $\hat{\mathbf{A}} = \mathbf{A} + \mathbf{I}$  is adjacent matrix with self-loop,  $\hat{\mathbf{D}}$  denotes the degree matrix of  $\hat{\mathbf{A}}$ , and  $\sigma$  is the activation function. The aggregation process is iterated  $L$  times to yield the final high-order representations  $\mathbf{Z} = \phi^{(L)}(\mathbf{A}, \mathbf{H}^{(L-1)})$ . Note that the aggregation is non-parametric.

Additionally, we explore two other aggregation types, the row-normalized Laplacian matrix  $\tilde{\mathbf{A}}\tilde{\mathbf{D}}^{-1}$  and column-normalized Laplacian matrix  $\tilde{\mathbf{D}}^{-1}\tilde{\mathbf{A}}$ . Table 16 presents the results of these three aggregation types. In this table, we observe that there is no significant difference in performance among the various aggregation methods. All of these methods can achieve desirable performance. Nevertheless, the GCN-like aggregation method consistently outperforms the others. We consider that if the encoder is non-parametric, then different aggregation methods will not bring significant differences in inductive bias, aligning with the findings in Yang et al. (2023a).

	Cora	Citeseer	PubMed	Amazon-CS	Amazon-Photo	Co-CS	Co-Phys	Wiki-CS	Flickr	Arxiv
GCN (Ours)	<b>84.60</b> $\pm 0.24$	<b>73.52</b> $\pm 0.53$	<b>86.99</b> $\pm 0.09$	<b>88.46</b> $\pm 0.16$	<b>94.28</b> $\pm 0.08$	<b>94.87</b> $\pm 0.07$	<b>96.17</b> $\pm 0.03$	<b>81.21</b> $\pm 0.13$	<b>49.85</b> $\pm 0.09$	<b>71.12</b> $\pm 0.10$
Col	84.14 $\pm 0.34$	73.48 $\pm 0.53$	86.92 $\pm 0.08$	87.93 $\pm 0.27$	93.11 $\pm 0.15$	94.81 $\pm 0.06$	96.09 $\pm 0.03$	80.62 $\pm 0.30$	49.15 $\pm 0.16$	71.03 $\pm 0.09$
Row	84.09 $\pm 0.32$	73.49 $\pm 0.54$	86.92 $\pm 0.08$	87.96 $\pm 0.27$	93.07 $\pm 0.15$	94.82 $\pm 0.06$	96.07 $\pm 0.04$	80.63 $\pm 0.25$	49.18 $\pm 0.10$	71.04 $\pm 0.09$

Table 16: Ablation study on node aggregation type. *GCN* indicates bi-normalized Laplacian aggregation matrix  $\tilde{\mathbf{D}}^{-1/2}\tilde{\mathbf{A}}\tilde{\mathbf{D}}^{-1/2}$ , *Col* indicates column-normalized Laplacian aggregation matrix  $\tilde{\mathbf{D}}^{-1}\tilde{\mathbf{A}}$ , and *Row* indicates row-normalized Laplacian aggregation matrix  $\tilde{\mathbf{A}}\tilde{\mathbf{D}}^{-1}$ .

In addition to analyzing the aggregation type, we also assess the performance of SSL-GM with varying numbers of aggregation layers, as depicted in Figure 7. In this figure, we illustrate the performance of SSL-GM on five benchmark datasets for both validation and testing sets under the transductive setting. We can see that the model attains the optimal performance with 2 or 3 layers, consistent with prior research on GNNs (Kipf & Welling, 2017; Hamilton et al., 2017). We consider that our model might encounter over-smoothing issues with a high number of layers (Li et al., 2019).

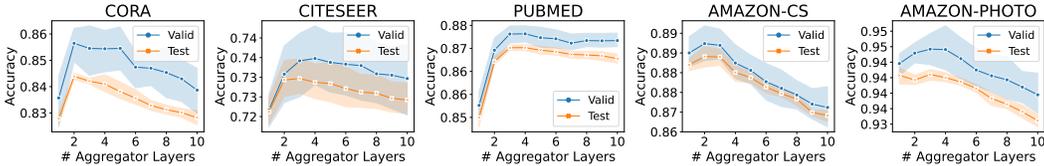


Figure 7: Node classification accuracy on transductive setting with different aggregation layers.

## E EXTENSION TO SUPERVISED SSL-GM

Our SSL-GM is a self-supervised method that follows the training then evaluation paradigm. In particular, the model initially generates representations for each node and subsequently trains classification heads for downstream tasks. Despite the generalization capability of self-supervised learning techniques, their performance typically lags behind that of supervised methods (Chen et al., 2020b). To bridge the gap, we implement a supervised version of SSL-GM by jointly optimizing the objective function and downstream task, which can be formulated as

$$\mathcal{L} = \mathcal{L}_{cont} + \lambda \cdot \mathcal{L}_{rec} + \beta \cdot \mathcal{L}_{sup}. \quad (19)$$

We evaluate the performance of supervised SSL-GM on transductive node classification and search the value of  $\beta$  within the range  $\{0.1, 1, 10\}$ . Table 17 reports the experimental results. It is counter-intuitive that our self-supervised version outperforms the supervised version on 7 out of 10 datasets. We suppose that the integration of cross-entropy loss may impact the learning process of contrastive learning by introducing additional gradients on model parameters. This could lead the model to acquire pseudo-knowledge unrelated to our prediction target. In other words, the model parameters might be optimized by the annotations rather than the outputs of GNNs.

	Cora	Citeseer	PubMed	Amazon-CS	Amazon-Photo	Co-CS	Co-Phys	Wiki-CS	Flickr	Arxiv
w/o sup.	<b>84.60</b> $\pm 0.24$	73.52 $\pm 0.53$	86.99 $\pm 0.09$	88.46 $\pm 0.16$	<b>94.28</b> $\pm 0.08$	<b>94.87</b> $\pm 0.07$	<b>96.17</b> $\pm 0.03$	<b>81.21</b> $\pm 0.13$	<b>49.85</b> $\pm 0.09$	<b>71.12</b> $\pm 0.10$
w. sup.	77.23 $\pm 1.41$	<b>75.20</b> $\pm 0.55$	<b>88.64</b> $\pm 0.27$	<b>88.68</b> $\pm 0.25$	94.17 $\pm 0.18$	94.60 $\pm 0.19$	95.89 $\pm 0.10$	80.64 $\pm 0.34$	46.39 $\pm 1.73$	70.60 $\pm 0.43$

Table 17: Ablation study on the model with or without supervisions.