# Contrastive Representation Learning Based on Multiple Node-centered Subgraphs

Dong Li
Tianjin University
Tianjin, China
ld2022244154@tju.edu.cn

Wenjun Wang
Tianjin University
Tianjin, China
wjwang@tju.edu.cn

Minglai Shao*
Tianjin University
Tianjin, China
shaoml@tju.edu.cn

Chen Zhao
Baylor University
Waco, Texas , USA
chen_zhao@baylor.edu

## ABSTRACT

As the basic element of graph-structured data, node has been recognized as the main object of study in graph representation learning. A single node intuitively has multiple node-centered subgraphs from the whole graph (e.g., one person in a social network has multiple social circles based on his different relationships). We study this intuition under the framework of graph contrastive learning, and propose a multiple node-centered subgraphs contrastive representation learning method to learn node representation on graphs in a self-supervised way. Specifically, we carefully design a series of node-centered regional subgraphs of the central node. Then, the mutual information between different subgraphs of the same node is maximized by contrastive loss. Experiments on various real-world datasets and different downstream tasks demonstrate that our model has achieved state-of-the-art results.

## CCS CONCEPTS

• **Computing methodologies** → **Unsupervised learning**; *Neural networks*; **Learning latent representations**.

## KEYWORDS

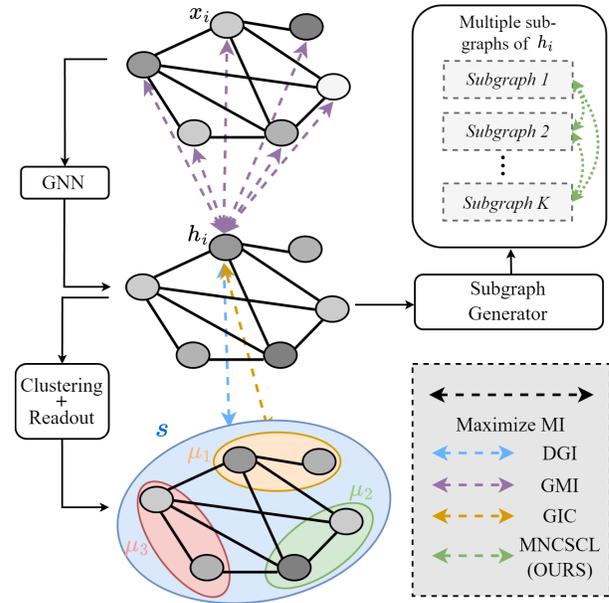Contrastive Learning, Node-centered Subgraph, Graph Representation Learning

Figure 1: A descriptive illustration of different proxy tasks among DGI, GMI, GIC and our proposed MNCSCL. The two-way arrows represent MI maximization, and the different colors represent different models. The subgraph generator and multiple subgraphs of $h_i$ are described in details in Section 3.1 and Section 3.2.

## 1 INTRODUCTION

Graph representation learning has received increasing attention recently [5], which aims to transform high-dimensional graph-structured data into low-dimensional dense vectorized representations. As the basic elements of graph-structured data, node representation has been the main object of graph representation learning.

A comprehensive node representation can be well used for a variety of downstream tasks, such as node classification [14] and link prediction [3].

A widespread graph representation learning method is the application of graph neural networks (GNN) [4, 14, 29, 32, 40]. But most of such methods focus on supervised learning, relying on supervised signals in the graph. For real-world networks, the acquisition of these supervised signals is often cumbersome and expensive. Self-supervised learning [11] is a popular research area in recent years, which designs proxy tasks for unlabeled data to mine the representational properties of the data itself as supervised information.

As one of the representative methods of self-supervised learning, the proxy task of contrastive learning is to maximize the Mutual Information (MI) [25] between the input and related content [34]. For example, Deep Graph Infomax (DGI) [30] maximizes MI between local view and global view of the input graph and its corresponding
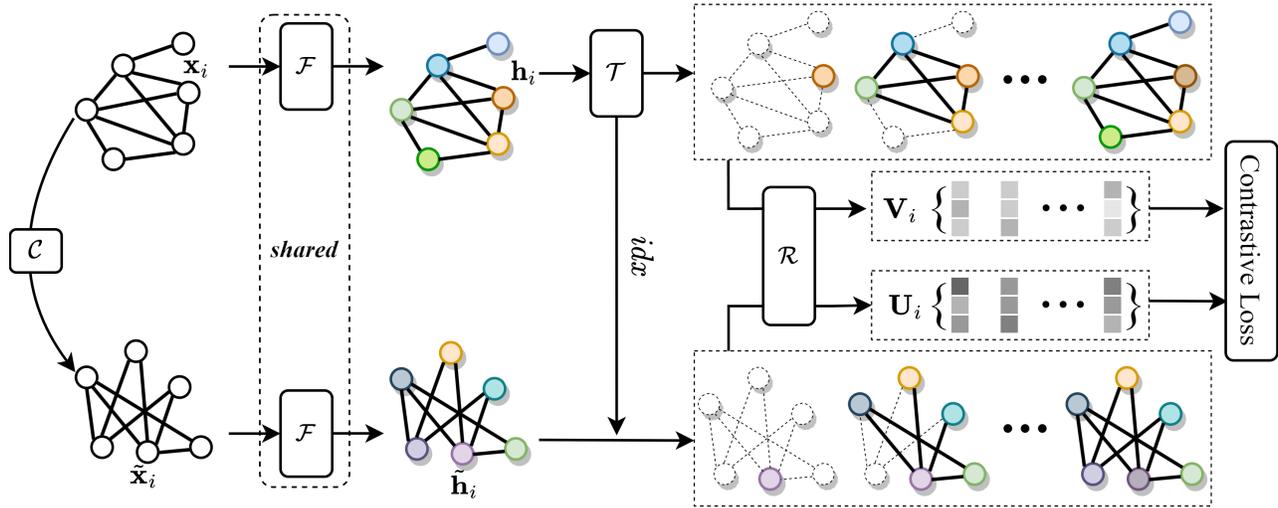
---

*Corresponding author

**Figure 2: The pipelines of MNCSCL. For a specific node $v_i$ with attribute $\mathbf{x}_i$, we first get its negative example $\tilde{\mathbf{x}}_i$ through perturbing the structure and attributes of the input graph $\mathcal{G}$ with a corruption function $C$. Then we use a subgraph generator $\mathcal{T}$ to get a series of node-centered subgraphs $\mathbf{G}_i$ and their corresponding negative set $\tilde{\mathbf{G}}_i$ from $\mathbf{h}_i$ and $\tilde{\mathbf{h}}_i$ (obtained by a shared encoder $\mathcal{F}$). Finally, the mutual information between $\mathbf{G}_i$ and $\tilde{\mathbf{G}}$ is maximized in the latent space $\mathbf{V}_i$ and $\mathbf{U}_i$ (obtained by a readout function $\mathcal{R}$) by contrastive loss. For more details, refer to the "overall framework" subsection in Section 3.**

corrupted graph. Graphical Mutual Information (GMI) [23] doesn't use corruption function, instead, it maximizes the MI between the hidden representation of nodes and their original local structure. Graph InfoClust (GIC) [19], on the other hand, maximizes the MI between the node representation and its corresponding cluster representation on the basis of DGI. Although these methods have achieved many advances, they all focus on the MI between node embeddings and only one related graph structure, as shown in Figure 1.

In reality, we can look at a specific thing from multiple perspectives. For graph data, we can observe individual nodes in a graph from multiple perspectives, yet little literature has focused on this. Intuitively, for an individual in a social network, there may be a social circle of relatives based on blood relations, a social circle of colleagues based on work relations, and other social circles of friends with many different interests. If we analyze this individual from these different social circles, it is actually equivalent to learning from multiple perspectives on the nodes in this network.

Based on this intuition, we propose Multiple Node-centered Subgraphs Contrastive Representation Learning (MNCSCL). MNC-SCL takes each node in the network as the center and samples its node-centered regional subgraphs under different semantics, thus forming several different perspectives of the corresponding node, as shown in Figure 1. More specifically, we first generate a negative example through the corrupt function, then generate a series of node-centered subgraphs of the original graph by the view generator, and sample the corresponding subgraphs on the negative example. Then, these subgraphs are fed into graph neural network

encoders to obtain the representations of central nodes and its subgraphs after pooling. Finally, the mutual information between different subgraphs of the same node is maximized in the latent space by contrastive learning objective function. Experimental results on a variety of datasets demonstrate the superb performance of our design. The major contributions of this paper are as follows:

- We propose a novel framework to learn node representation through multiple node-centered subgraphs of nodes, which is a novel idea in current work to observe a single node from multiple perspectives.
- We carefully design five node-centered subgraphs and analyze the influence of different subgraphs on the learning quality of node representation through extensive experiments, which is of reference significance.
- We evaluated MNCSCL on five standard datasets and two downstream tasks to validate the effectiveness of the proposed method. Our experiments show that the contrastive learning of multiple subgraphs outperforms the above mentioned single-subgraph contrastive learning in terms of results.

## 2 RELATED WORK

### 2.1 Graph Representation Learning Based on Contrastive Learning

Inspired by the advances of contrastive learning in fields such as CV and NLP, some work has started to apply contrastive learning on graph data for graph representation learning. The objective of

graph contrastive learning is to maximize the MI between similar instances in a graph [33], and its model design focuses on three modules: data augmentation, proxy tasks, and contrastive objectives. Among them, the most important modules is the proxy task, which describes the definition of similar instances (i.e. positive example) and dissimilar instances (i.e. negative example). DGI [30] extends the idea of DIM [8] to graph data to learn node representations by maximizing the MI between local node representation and global graph representation. GMI [23] takes nodes and their neighbors as objects of study and maximizes the MI between hidden representation of each node and the original features of its neighboring nodes. GIC [19] clusters the nodes in the graph by a differentiable version of K-means clustering, and then maximizes the MI between the node representation and its corresponding cluster summaries. SUBG-CON [10] obtains the context subgraph of each node by subgraph sampling based data augmentation, and then maximizes the consistency between them. Despite the good results achieved, these works perform graph contrastive learning only on a single perspective for nodes.

## 2.2 Multi-view Contrastive Learning

Recently, multi-view representation learning has become a rapidly growing direction in machine learning and data mining areas [17, 22, 37–40]. It has had a lot of success in areas such as computer vision. For instance, Contrastive Multiview Coding (CMC) [28] uses contrastive learning to maximize the mutual information between multiple views of a dataset to perform representation learning of images. MVGRL [6] obtains multiple views of graph through data augmentation, and they find out that unlike visual representation learning, increasing the number of views of the entire graph to more than two by data augmentation does not improve performance. Unlike the multi-view graph contrastive learning summarized in MVGRL which focuses on node attributes at graph-level, our multiple node-centered subgraphs in this paper focus more on the differences in structure at node-level.

## 3 METHODOLOGY

**Problem definition.** Given a graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ with $N$ nodes, where $\mathcal{V} = \{v_1, v_2, ..., v_N\}$ and $\mathcal{E}$ represent the node set and the edge set respectively. $\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, ..., \mathbf{x}_N\} \in \mathbb{R}^{N \times F}$ is the node features matrix, where $\mathbf{x}_i \in \mathbb{R}^F$ denotes the features of dimension $F$ for node $v_i$. We use the adjacency matrix $\mathbf{A} \in \mathbb{R}^{N \times N}$ to represent the connectivity of the graph, where $\mathbf{A}(i, j) = 1$ if nodes $v_i$ and $v_j$ are linked, and $\mathbf{A}(i, j) = 0$ otherwise. In this way, a graph can also be represented as $\mathcal{G} = (\mathbf{X}, \mathbf{A})$. If $\mathcal{V}' \subset \mathcal{V}$ is a subset of vertices of $\mathcal{G}$ and $\mathcal{E}'$ consists of all of the edges in $\mathcal{E}$ that have both endpoints in $\mathcal{V}$, then the subgraph $\mathcal{S} = (\mathcal{V}', \mathcal{E}')$ of graph $\mathcal{G}$ is an induced subgraph. The subgraphs mentioned in this paper are all induced subgraphs.

The goal of self-supervised graph representation learning is to learn a encoder $\mathcal{F} : \mathbb{R}^{N \times F} \times \mathbb{R}^{N \times N} \rightarrow \mathbb{R}^{N \times F'}$, which takes the features matrix $\mathbf{X}$ and the adjacency matrix $\mathbf{A}$ as input to get the node representation $\mathbf{H} = \{\mathbf{h}_1, \mathbf{h}_2, ..., \mathbf{h}_N\} \in \mathbb{R}^{N \times F'}$ without label information, formulated as $\mathbf{H} = \mathcal{F}(\mathbf{X}, \mathbf{A})$. The learned node representation $\mathbf{H}$ can be used directly for downstream tasks such

**Table 1: Summary of notations. The first five notations are described in detail in the following subsections.**

| Notation | Meaning |
|---|---|
| $\mathcal{F}$ | Shared encoder |
| $C$ | Corruption function |
| $\mathcal{T}$ | Subgraph generator |
| $\mathcal{R}$ | Readout function |
| $\mathcal{D}$ | Discriminator |
| $N$ | Number of nodes in the input graph |
| $F$ | Feature dimension of each node |
| $F'$ | Feature dimension of each node representation |
| $K$ | Number of node-centered subgraphs sampled by $\mathcal{T}$ |
| $N'$ | Number of nodes in the corresponding node-centered subgraph |
| $d$ | Range of neighbors in the **neighboring subgraph** |
| $l$ | Number of nodes that are most similar to the central node for **intimate subgraph** |
| $C$ | number of clusters for **communal subgraph** |
| $\eta$ | self-weighted factor for **full subgraph** |
| $\mathcal{G}, \tilde{\mathcal{G}}$ | Input graph and its negative example obtained by $C$ |
| $\mathcal{G}_i^k, \tilde{\mathcal{G}}_i^k$ | The $k$-th node-centered subgraph of node $v_i$ and its negative example |
| $\mathbf{H}_i^k, \tilde{\mathbf{H}}_i^k$ | Node representation matrix of the $k$-th node-centered subgraph of node $v_i$ and its negative example |
| $\mathbf{A}_i^k, \tilde{\mathbf{A}}_i^k$ | Adjacency matrix of the $k$-th node-centered subgraph of node $v_i$ and its negative example |
| $\mathbf{v}_i^k, \mathbf{u}_i^k$ | Representation of the $k$-th node-centered subgraph of node $v_i$ and its negative example, obtained by $\mathcal{R}$ |

as node classification and link prediction. For the sake of clarity, we list all important notations in Table 1.

**Overall framework.** Inspired by recent graph representation learning work based on contrastive learning, we propose MNCSCL algorithm for graph representation learning by maximizing MI of multiple node-centered subgraphs of nodes. As illustrated in Figure 2, if there is only a single graph provided as input, the summarized steps of MNCSCL are as follows:

- Utilize a corruption function $C$ to perturb the structure and attributes of the input graph $\mathcal{G}$ to obtain a negative example $\tilde{\mathcal{G}} = (\tilde{\mathbf{X}}, \tilde{\mathbf{A}}) \sim C(\mathbf{X}, \mathbf{A})$.
- Pass input graph $\mathcal{G}$ and negative example $\tilde{\mathcal{G}}$ into a shared encoder $\mathcal{F}$ to get node representation $\mathbf{H}$ and $\tilde{\mathbf{H}}$.
- Use a subgraph generator $\mathcal{T}$ to sample a series of node-centered subgraphs $\mathbf{G}_i = \{\mathcal{G}_i^1, \mathcal{G}_i^2, ..., \mathcal{G}_i^K\}$ for node $v_i$ from the input graph $\mathcal{G}$, where $K$ is the number of different subgraphs and $\mathcal{G}_i^k = (\mathbf{H}_i^k, \mathbf{A}_i^k)$, $k = 1, 2, ..., K$. The corresponding node-centered subgraphs set $\tilde{\mathbf{G}}_i = \{\tilde{\mathcal{G}}_i^1, \tilde{\mathcal{G}}_i^2, ..., \tilde{\mathcal{G}}_i^K\}$ from the negative example $\tilde{\mathcal{G}}$ are further obtained according to $\mathbf{G}_i$.
- Summarize all subgraphs in $\mathbf{G}_i$ and $\tilde{\mathbf{G}}_i$ through a readout function $\mathcal{R}$ to get their representations $\mathbf{V}_i = \{\mathbf{v}_i^1, \mathbf{v}_i^2, ..., \mathbf{v}_i^K\}$
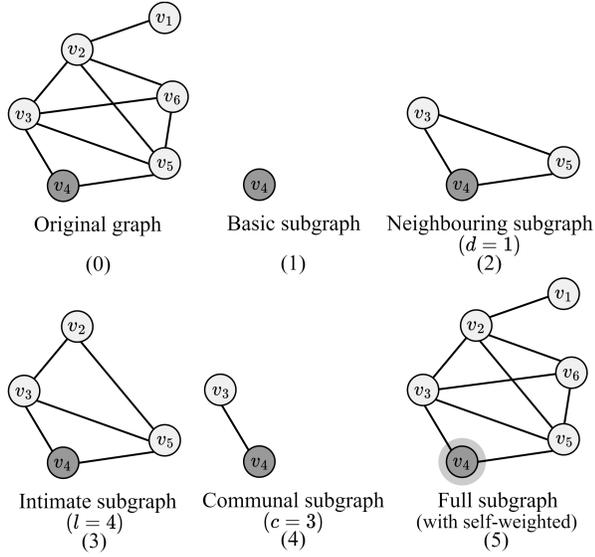
**Figure 3: Five types of node-centered subgraphs for node $v_4$ from original graph. Note that nodes are divided into 3 clusters in (4), which are $\{v_1, v_2\}$, $\{v_3, v_4\}$ and $\{v_5, v_6\}$.**

of node $v_i$ and the corresponding negative samples $\mathbf{U}_i = \{\mathbf{u}_i^1, \mathbf{u}_i^2, ..., \mathbf{u}_i^K\}$. As an example, $\mathbf{v}_i^k = \mathcal{R}(\mathbf{H}_i^k)$.
- Update parameters of $\mathcal{F}$, $\mathcal{R}$ and $\mathcal{D}$ (mentioned later) by applying gradient descent to maximize Eq. (13) or Eq. (14).

In the following sections, we will elaborate on the crucial components mentioned above.

## 3.1 Subgraph Generator

We can deal with a whole graph from different views [6], and the same is true for any node in the graph. For a single node, there are many subgraphs centered on it (e.g., ego network [42]). If these node-centered regional subgraphs have certain semantic information (e.g., ego network represents a specific individual and other persons who have a social relationship with him), then we can treat them as different perspectives of the central node.

The main role of the subgraph generator $\mathcal{T}$ is to sample these node-centered regional subgraphs from $\mathcal{G}$ and generate the corresponding negative samples from $\tilde{\mathcal{G}}$. Specifically, for a specific node $v_i$ and a prepared node-centered subgraph type $k$, the subgraph generator $\mathcal{T}$ first gets $idx$, a set which represents the index of chosen nodes for subgraph $\mathcal{G}_i^k$ to be obtained. Then, the node representation matrix $\mathbf{H}_i^k \in \mathbb{R}^{N' \times F'}$ and adjacency matrix $\mathbf{A}_i^k \in \mathbb{R}^{N' \times N'}$ of $\mathcal{G}_i^k$ are denoted respectively as

$$\mathbf{H}_i^k = \mathbf{H}_{idx,:}, \ \mathbf{A}_i^k = \mathbf{A}_{idx,idx}, \tag{1}$$

where $\cdot_{idx}$ is an indexing operation and $N'$ is the length of $idx$.

In this way, we can obtain the $k$-th node-centered subgraph $\mathcal{G}_i^k = (\mathbf{H}_i^k, \mathbf{A}_i^k) \sim \mathcal{T}(\mathbf{H}, \mathbf{A})$ for any specific node $v_i$. Likewise, the corresponding negative example can be obtained by $\tilde{\mathcal{G}}_i^k = (\tilde{\mathbf{H}}_i^k, \tilde{\mathbf{A}}_i^k) = (\tilde{\mathbf{H}}_{idx,:}, \tilde{\mathbf{A}}_{idx,idx})$.

## 3.2 Node-centered Subgraphs Design

To learn a more comprehensive representation, we carefully design five different node-centered subgraphs, as illustrated in Figure 3. The details of them are as follows:

**Subgraph 1: Basic subgraph.** The basic subgraph only contains the central node itself (i.e., $N' = 1$), that means for each node $v_i$:

$$idx = \{i\}. \tag{2}$$

Further, we can get the basic subgraph representation and its corresponding negative example by $\mathbf{v}_i^1 = \mathbf{h}_i$ and $\mathbf{u}_i^1 = \tilde{\mathbf{h}}_i$. For any specific node, basic subgraph is the purest "subgraph" as well as the main subgraph, which contains the most concentrated features of the node itself.

**Subgraph 2: Neighboring subgraph.** The neighbors of the central node are often closely related to the node in structure, and study with them can better capture the structural features of node [4]. The neighboring subgraph contains all nodes with a distance less than or equal to $d$ from the central node $v_i$, denoted as

$$idx = \{j | dis(v_i, v_j) \le d\}, \tag{3}$$

where $dis(\cdot, \cdot)$ is a function used to get the distance between two nodes. After obtaining the neighboring subgraph $\mathcal{G}_i^2$ by Eq. (1), a readout function $\mathcal{R} : \mathbb{R}^{N \times F'} \to \mathbb{R}^{F'}$ is used to obtain the neighboring subgraph representation and its corresponding negative example:

$$\mathbf{v}_i^2 = \mathcal{R}(\mathbf{H}_i^2), \ \mathbf{u}_i^2 = \mathcal{R}(\tilde{\mathbf{H}}_i^2). \tag{4}$$

It is beneficial to learn a more comprehensive representation when we take a larger range of neighbors. But at the same time, the features of the central node will be weakened, which will cause the model to be more inclined to learn the representation of a region or even the whole graph. It follows that it's important to choose the value of $d$. As shown in the Figure 5, the model reaches the best performance when $d = 1$ for the neighboring subgraph.

**Subgraph 3: Intimate subgraph.** The intimate subgraph takes into account the similarity that actually exists between two nodes in the input graph $\mathcal{G}$. It contains the first $l$ nodes that are most similar to the central node. This is equivalent to identifying the nodes that are structurally close to the central node from another perspective, and thus learning the structural features of the nodes more comprehensively.

The similarity between nodes is usually measured by a similarity scores matrix $\mathbf{S} \in \mathbb{R}^{N \times N}$, where $\mathbf{S}(i, j)$ measures the similarity between nodes $v_i$ and $v_j$. Here we follow the personalized pagerank (PPR) algorithm [9] as introduced in [36]. The similarity scores matrix $\mathbf{S}$ based on the PPR algorithm can be denoted as

$$\mathbf{S} = \alpha \cdot (\mathbf{I} - (1 - \alpha) \cdot \bar{\mathbf{A}})^{-1}, \tag{5}$$

where $\mathbf{I}$ is the identity matrix and $\bar{\mathbf{A}} = \mathbf{A}\mathbf{D}^{-1}$ denotes the column-normalized adjacency matrix. $\mathbf{D}$ is the diagonal matrix corresponding to $\mathbf{A}$ with $\mathbf{D}(i, i) = \sum_j \mathbf{A}(i, j)$ on its diagonal. $\alpha \in [0, 1]$ is a parameter which is set as 0.15 in [10]. For a specific node $v_i$, the subgraph generator $\mathcal{T}$ chooses its top-$l$ similar nodes (i.e., $N' = l$) to generate intimate subgraph with $\mathbf{S}(i, :)$, which can be denoted as

$$idx = top\_rank(\mathbf{S}(i, :), l), \tag{6}$$

where $top\_rank(\cdot)$ is a function that selects the top-$l$ values from a vector and return the corresponding indices. Same as *Subgraph*

*2* subsection, we can obtain the intimate subgraph representation and its corresponding negative example with $\mathbf{v}_i^3 = \mathcal{R}(\mathbf{H}_i^3)$ and $\mathbf{u}_i^3 = \mathcal{R}(\tilde{\mathbf{H}}_i^3)$.

**Subgraph 4: Communal subgraph.** In graph clustering, nodes in a uniform cluster tend to have similarity in attributes. Therefore, we can select all nodes in the cluster to which the central node belongs to get the communal subgraph. These nodes with similar attributes to the central node can help the model better learn the attribute features of the central node.

It is particularly noteworthy that the other subgraphs are obtained independently of the node attributes (i.e., $\mathbf{H}$), but the communal subgraph is attribute-related. Since $\mathbf{H}$ will change during training, a fixed *idx* before the model training as other subgraphs may lead to undesirable results.

There are already many methods to perform graph clustering [26]. We tried three different clustering strategies based on *K*-means clustering [18] due to the stable and excellent performance of it.

- **Strategy 1: Precomputed K-means.** Same as other subgraphs, the node-centered subgraphs are sampled before model training starts. Specifically, the traditional *K*-means clustering algorithm is used to cluster the nodes in the input graph $\mathcal{G}$. For any specific node $v_i$, take all the indices of nodes in its cluster as *idx* and further compute $\mathbf{v}_i^4$ and $\mathbf{u}_i^4$.

- **Strategy 2: A differentiable version of K-means.** To update the communal subgraph during training, we need an end-to-end clustering algorithm. Here we follow a differentiable version of K-means as introduced in [31]. For each node $v_i$, let $\mu_c$ denote the center of cluster $c$ and $\hat{\gamma}_{ic}$ (s.t., $\sum_c \hat{\gamma}_{ic} = 1, \forall i$) denotes the degree to which node $v_i$ is assigned to cluster $c$. Suppose that the number of clusters to be obtained is $C$, ClusterNet updates $\mu_c$ via an iterative process by alternately setting.

$$\mu_c = \frac{\sum_i \hat{\gamma}_{ic} \mathbf{h}_i}{\sum_i \hat{\gamma}_{ic}} \quad c = 1, ..., C, \tag{7}$$

and
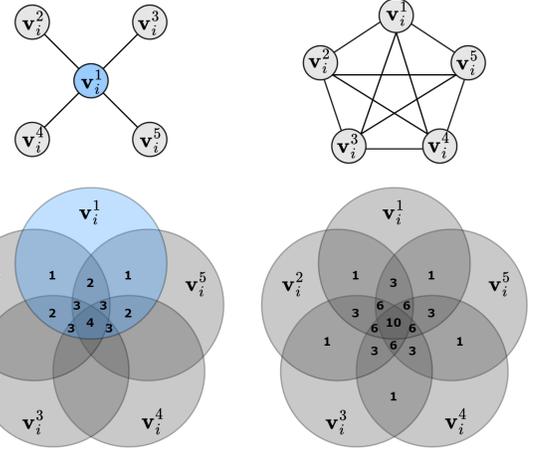
$$\hat{\gamma}_{ic} = \frac{exp(-\beta \cdot sim(\mathbf{h}_i, \mu_c))}{\sum_c exp(-\beta \cdot sim(\mathbf{h}_i, \mu_c))} \quad c = 1, ..., C, \tag{8}$$

where $\beta$ is an inverse-temperature hyperparameter, the standard K-means assignment is recovered when $\beta \to \infty$. $sim(\cdot, \cdot)$ denotes a similarity function between two instances. Eventually, we can get the communal subgraph representation by
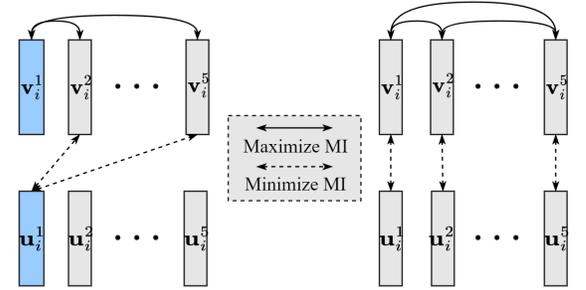
$$\mathbf{v}_i^4 = \sigma \left( \sum_{c=1}^{C} \hat{\gamma}_{ic} \mu_c \right), \tag{9}$$

where $\sigma(\cdot)$ is the logistic sigmoid nonlinearity. Since this method does not get *idx*, we simply get negative example $\{\mathbf{u}_1^4, \mathbf{u}_2^4, ..., \mathbf{u}_N^4\}$ of all nodes by row-wise shuffling of $\{\mathbf{v}_1^4, \mathbf{v}_2^4, ..., \mathbf{v}_N^4\}$.

- **Strategy 3: An end-to-end version of K-means with an estimation network.** In this way, we replace the iterative process in *Strategy 2* with an estimation network, which utilizes a multi-layer neural network to directly predict the degree to which each node belongs to each cluster $\hat{\gamma} \in \mathbb{R}^{N \times C}$,



(a) The "core view" and "full graph" paradigms



(b) Contrastive lossess under core view and full graph cases

**Figure 4: Take the five node-centered subgraphs of node $v_i$ as an example. (a) The "core view" (left) and "full graph" (right) paradigms. The numbers within the regions represent the number of MI in this region. For example, if we select all 5 node-centered subgraphs under the full graph case, MI will be calculated once between every two subgraphs, and hence is marked with the number 10. (b) Contrastive lossess under core view and full graph cases. Refer to Section 3.3 for more details.**

denoted as

$$\hat{\gamma} = softmax(MLP(\mathbf{H}; \theta)), \tag{10}$$

where $softmax(\cdot)$ is a softmax nonlinearity and $MLP(\cdot)$ is a multi-layer neural network with trainable parameters $\theta$. Then we get the communal subgraph representation by Eq. (8) and Eq. (9) as same as *Strategy 2*.

The comparison of the three strategies is shown in Figure 5. After weighing both accuracy and efficiency, we chose *Strategy 2* to generate the communal subgraph.

**Subgraph 5: Full subgraph.** To learn a comprehensive representation of a node, it is essential to observe it from a global perspective. The full subgraph contains all the nodes (i.e. $N' = N$) in the input graph $\mathcal{G}$, e.g., for any specific node $v_i$,

$$idx = \{j | j = 1, 2, ..., N\}. \tag{11}$$

**Table 2: The statistics of all the five datasets. *Note that the node classification on PPI dataset is a multilabel classification problem.**

| Task | Dataset | Type | #Nodes | #Edges | #Features | #Classes | Train / Val / Test |
|---|---|---|---|---|---|---|---|
| Node classification | Cora | Citation network | 2,708 | 5,429 | 1,433 | 7 | 0.05 / 0.18 / 0.37 |
| & Link prediction | Citeseer | Citation network | 3,327 | 4,732 | 3,703 | 6 | 0.04 / 0.15 / 0.30 |
| (Transductive) | Pubmed | Citation network | 19,717 | 44,338 | 500 | 3 | 0.003 / 0.03 / 0.05 |
| Node classification | Reddit | Social network | 232,965 | 11,606,919 | 602 | 41 | 0.66 / 0.10 / 0.24 |
| (Inductive) | PPI | Protein network | 56,944 | 818,716 | 50 | 121* | 0.79 / 0.11 / 0.10 |

Compared with the previous subgraphs, the full subgraph contains far more nodes than they do, which extremely weakens the specificity of the central node. At the same time, it is not conducive to learning a specialized node representation as all nodes have the same full subgraph. Based on these ideas, we propose full subgraph for specific node $v_i$ with self-weighted, denoted as

$$\mathbf{v}_i^5 = (1-\eta)\mathcal{R}(\mathbf{H}_i^5) + \eta\mathbf{h}_i, \tag{12}$$

where $\eta \in [0, 1]$ is a self-weighted factor.

So far, we have introduced five carefully designed node-centered subgraphs. It is noted that subgraphs other than *Subgraph 4* can be precomputed before model training starts, which allows us to quickly obtain these subgraphs during training by performing only one calculation before training.

### 3.3 Contrastive Loss

The key idea of self-supervised comparative learning is to define a proxy task to generate positive and negative samples. The encoder $\mathcal{F}$ that generates the node representation is trained by contrast between positive and negative samples.

To handle multiple subgraphs we obtained before, we take the "core view" (CV) and "full graph" (FG) paradigms as introduced in CMC [28], as shown in Figure 4.(a). Further, the contrastive lossess under core view and full graph cases are illustrated in Figure 4.(b). Specifically, in the core view case, we regard *Subgraph 1* as the most critical node-centered subgraph. It and its corresponding negative sample constitute positive and negative pairs with *Subgraph 2~5*, respectively. As for the full graph case, any two between *Subgraph 1~5* constitute positive pairs, and *Subgraph 1~5* respectively form negative pairs with their corresponding negative examples.

After defining the proxy task, we follow the intuitions from DGI [30] and use a noise-contrastive type objective with a standard binary cross-entropy (BCE) loss between positive examples and negative examples. MNCSCL's objective under core view case is

$$\mathcal{L}_{CV} = \sum_{i=1}^{N} \sum_{j=2}^{K} \mathbb{E}_{(\mathbf{X},\mathbf{A})} \left[ log\mathcal{D}(\mathbf{v}_i^1, \mathbf{v}_i^j) \right]$$
$$+ \sum_{i=1}^{N} \sum_{j=2}^{K} \mathbb{E}_{(\tilde{\mathbf{X}},\tilde{\mathbf{A}})} \left[ log(1 - \mathcal{D}(\mathbf{u}_i^1, \mathbf{v}_i^j)) \right], \tag{13}$$

where $K$ is the number of selected node-centered subgraphs and $\mathcal{D} : \mathbb{R}^{F'} \times \mathbb{R}^{F'} \to \mathbb{R}$ is a discriminator which is used for estimating the MI by assigning higher scores to positive examples than negatives.

When using the full graph case, the objective becomes

$$\mathcal{L}_{FG} = \sum_{i=1}^{N} \sum_{j=1}^{K-1} \sum_{k=j+1}^{K} \mathbb{E}_{(\mathbf{X},\mathbf{A})} \left[ log\mathcal{D}(\mathbf{v}_i^j, \mathbf{v}_i^k) \right]$$
$$+ \sum_{i=1}^{N} \sum_{j=1}^{K} \mathbb{E}_{(\tilde{\mathbf{X}},\tilde{\mathbf{A}})} \left[ log(1 - \mathcal{D}(\mathbf{v}_i^j, \mathbf{u}_i^j)) \right]. \tag{14}$$

The Eq. (13) and Eq. (14) are used as the contrastive loss in the experiments respectively.

## 4 EXPERIMENTS

### 4.1 Experimental Settings

**Datasets.** We use 5 commonly used benchmark datasets in the previous work [4, 19] for node classification and link prediction downstream tasks, including 3 transductive citation networks (i.e., Cora, Citeseer, and Pubmed), a inductive large social network (i.e., Reddit) and a inductive protein-protein interaction dataset that contains multiple graphs with multiple labels (i.e., PPI).

- **Cora, Citeseer, and PubMed** are all citation networks, with Cora and Citeseer focusing on papers in computer science and information science, and Pubmed containing a large amount of literature information in medical and life science fields. They represent citation relationships between papers through graph data structure, where each node represents a paper and the edges represent the citation relationships between papers.
- **Reddit** is a collection of information from Reddit, the world's largest social news aggregation, discussion and community site. the Reddit dataset provides a large amount of user-generated content, including posts, comments, polls and more. In Reddit, posts are represented as nodes, and the connections between them correspond to user comments.
- **PPI** is a protein-protein interaction dataset that contains multiple graphs with multiple labels. PPI contains the interaction relationships between proteins that can form a network or graph structure. Each node represents a protein, while edges indicate interactions between proteins.

We use all five datasets in the node classification task and follow the settings in the division of the training set and the test set as same as [30]. In the link prediction task, we use Cora, Citeseer, and Pubmed datasets and follow the setup described in [15]. The statistics of all the datasets are shown in Table 2.

**Table 3: The classification accuracy (in %) on the transductive datasets and the micro-averaged F1 ($\times 100$) on the inductive datasets of the node classification task. Some results are directly taken from their original papers (DGI, GMI#inductive, GIC and GRACE#inductive), and other compared results are taken from [10, 20]. The second column is the data used in the training process (X: features matrix, A: adjacency matrix, Y: labels). The best result for each dataset is indicated by bolded.**

| Method | Input | | | Transductive | | | Inductive | |
|---|---|---|---|---|---|---|---|---|
| | X | A | Y | Cora | Citeseer | Pubmed | Reddit | PPI |
| Raw features | ✓ | | | 56.6 ± 0.4 | 57.8 ± 0.2 | 69.1 ± 0.2 | 58.5 ± 0.1 | 42.5 ± 0.3 |
| Deep Walk | | ✓ | | 67.2 | 43.2 | 65.3 | 32.4 | 52.9 |
| GCN | ✓ | ✓ | ✓ | 81.4 ± 0.6 | 70.3 ± 0.7 | 76.8 ± 0.6 | 93.3 ± 0.1 | 51.5 ± 0.6 |
| FastGCN | ✓ | ✓ | ✓ | 78.0 ± 2.1 | 63.5 ± 1.8 | 74.4 ± 0.8 | 89.5 ± 1.2 | 63.7 ± 0.6 |
| DGI | ✓ | ✓ | | 82.3 ± 0.6 | 71.8 ± 0.7 | 76.8 ± 0.6 | 94.0 ± 0.1 | 63.8 ± 0.2 |
| GMI | ✓ | ✓ | | 83.0 ± 0.2 | 72.4 ± 0.2 | 79.9 ± 0.4 | 95.0 ± 0.02 | 65.0 ± 0.02 |
| GIC | ✓ | ✓ | | 81.7 ± 0.8 | 71.9 ± 0.9 | 77.4 ± 0.5 | - | - |
| GRACE | ✓ | ✓ | | 83.1 ± 0.2 | 72.1 ± 0.1 | 79.6 ± 0.5 | 94.2±0.0 | 66.2±0.1 |
| MVGRL | ✓ | ✓ | | 82.9 ± 0.3 | 72.6 ± 0.4 | 80.1 ± 0.7 | - | - |
| MNCSCL-FG | ✓ | ✓ | | 84.3 ± 0.5 | 73.2 ± 0.6 | 80.0 ± 0.4 | 95.2 ± 0.1 | **67.3 ± 0.2** |
| MNCSCL-CV | ✓ | ✓ | | **84.7 ± 0.3** | **73.8 ± 0.5** | **81.5 ± 0.4** | **95.8 ± 0.1** | 67.1 ± 0.2 |

**Baseline methods.** In node classification task, the compared methods include direct use of row features, 1 traditional unsupervised algorithm (i.e., Deep Walk [24]), two supervised graph neural networks (i.e., GCN [14] and FastGCN [1]) and 5 state-of-the-art self-supervised methods(i.e., DGI, GMI, GIC, GRACE [41] and MVGRL [6]).

- **DGI** is one of the classical methods of graph representation learning based on contrastive learning. It aims to maximize the MI between the local perspective and the global perspective of the input graph, as well as the corresponding corrupted graph.
- **GMI** draws on the ideas of DGI, but rather than employing a corruption function, this approach focuses on maximizing the MI between the hidden representations of nodes and their original local structure.
- **GIC** is also inspired by DGI, its objective is to maximize the MI between the node's representation and the representation of the cluster to which it is assigned.
- **GRACE** propose a novel framework for unsupervised graph representation learning by leveraging a contrastive objective at the node level. In order to enhance the contrast effect, they created two sets of negative pairs, one within the same view and the other across different views.
- **MVGRL** performs self-supervised learning by contrasting structural views of graphs, where they contrast first-order neighbors of nodes as well as a graph diffusion, with good results.

In link prediction, we directly follow the effective link prediction methods used in GIC (i.e., Deep Walk, Spectral Clustering (SC) [27], VGAE [15], ARGVA [21], DGI and GIC).

- **Spectral Clustering** is initially introduced to solve the node partitioning problem in graph analysis. It has demonstrated satisfactory performance across diverse domains, such as graphs, text, images, and microarray data. Its effectiveness has been widely acknowledged in these areas.

- **VGAE** is an unsupervised learning framework designed for graph-structured data, utilizing the variational auto-encoder (VAE) methodology. In VGAE, a GNN-based encoder is employed to generate node embeddings, while a straightforward decoder is used to reconstruct the adjacency matrix.
- **ARGVA** is a graph embedding framework specifically designed for graph data, incorporating adversarial learning techniques. Similar to VGAE, ARGVA adopts a similar structure, but it learns the underlying data distribution through an adversarial approach.

**Evaluation metrics.** For the node classification task, we classify the test set by a logistic regression classifier, and then evaluate the performance using classification accuracy for transductive datasets (i.e., Cora, Citeseer and Pubmed) and micro-averaged F1 score for inductive datasets (i.e., Reddit and PPI). Suppose that TP, FN, FP and TN represent the number of true positives, false negatives, false positives, and true negatives, respectively. Then classification accuracy can be calculated by $accuracy = (TP + TN)/(TP + FP + TN + FN)$. Also micro-averaged F1 score can be calculated by $F1 - Score = 2 * precision * recall/(precision + recall)$, where $precision = TP/(TP + FP)$ and $recall = TP/(TP + FN)$. For the link prediction task, we use the AUC score (the area under ROC curve) and the AP score (the area under Precision-Recall curve) for evaluation. The closer the AUC score and the AP score approaches 1, the better the performance of the algorithm is.

**Training strategy.** We implement MNCSCL using PyTorch [12] on 4 NVIDIA GeForce RTX 3090 GPUs and use Adam optimizer [13] with an initial learning rate of 0.001 (specially, 0.0001 for Reddit) during training. We follow settings in DGI that using an early stopping strategy with a patience of 20 epochs for transductive datasets and a fixed number of epochs (150 on Reddit, 20 on PPI) for inductive datasets. For large graphs, we adopt the sampling strategy used by GraphSAGE [4].

**Table 4: The AUC scores (in %) of the link prediction task. The results of the compared methods are replicated from [19]. The best result for each dataset is indicated by bolded.**

| Method | Cora | | Citeseer | | Pubmed | |
|---|---|---|---|---|---|---|
| | AUC | AP | AUC | AP | AUC | AP |
| DeepWalk | 83.1 ± 0.01 | 85.0 ± 0.00 | 80.5 ± 0.02 | 83.6 ± 0.01 | 84.4 ± 0.00 | 84.1 ± 0.00 |
| Spectral Clustering | 84.6 ± 0.01 | 88.5 ± 0.00 | 80.5 ± 0.01 | 85.0 ± 0.01 | 84.2 ± 0.02 | 87.8 ± 0.01 |
| VGAE | 91.4 ± 0.01 | 92.6 ± 0.01 | 90.8 ± 0.02 | 92.0 ± 0.02 | 96.4 ± 0.00 | 96.5 ± 0.00 |
| ARGVA | 92.4 ± 0.004 | 93.2 ± 0.003 | 92.4 ± 0.003 | 93.0 ± 0.003 | **96.8 ± 0.001** | **97.1 ± 0.001** |
| DGI | 89.8 ± 0.8 | 89.7 ± 1.0 | 95.5 ± 1.0 | 95.7 ± 1.0 | 91.2 ± 0.6 | 92.2 ± 0.5 |
| GIC | 93.5 ± 0.6 | 93.3 ± 0.7 | 97.0 ± 0.5 | 96.8 ± 0.5 | 93.7 ± 0.3 | 93.5 ± 0.3 |
| MNCSCL-CV | **94.8 ± 0.4** | **94.2 ± 0.6** | **97.7 ± 0.4** | **97.2 ± 0.5** | 94.8 ± 0.2 | 95.4 ± 0.4 |

## 4.2 Implementation Details

**Encoder design.** For transductive datasets, we adopt a one-layer Graph Convolutional Network (GCN) as our encoder, with the following propagation rule:

$$\mathcal{F}(\mathbf{X}, \mathbf{A}) = \sigma(\hat{\mathbf{D}}^{-\frac{1}{2}} \hat{\mathbf{A}} \hat{\mathbf{D}}^{-\frac{1}{2}} \mathbf{X} \mathbf{W}), \tag{15}$$

where $\hat{\mathbf{A}} = \mathbf{A} + \mathbf{I}_N$ is the adjacency matrix with self-loops and $\hat{\mathbf{D}}(i,i) = \sum_j \hat{\mathbf{A}}(i,j)$ is its corresponding degree matrix. $\sigma(\cdot)$ is the PReLU nonlinearity [7] and $\mathbf{W}$ is a learnable parameter matrix with $F' = 512$ (specially, $F' = 256$ on Pubmed). As for inductive datasets, we adopt a one-layer GCN with skip connections [35] as our encoder, with the following propagation rule:

$$\mathcal{F}(\mathbf{X}, \mathbf{A}) = \sigma(\hat{\mathbf{D}}^{-\frac{1}{2}} \hat{\mathbf{A}} \hat{\mathbf{D}}^{-\frac{1}{2}} \mathbf{X} \mathbf{W} + \hat{\mathbf{A}} \mathbf{W}_{skip}), \tag{16}$$

where $\mathbf{W}_{skip}$ is a learnable parameter matrix with $F' = 512$ for skip connections.

**Corruption function.** For transductive datasets, we transform adjacency matrix $\mathbf{A}$ to a diffusion matrix $\mathbf{U}$. Specifically, we compute diffusion using fast approximation and sparsification methods [2] with heat kernel [16]:

$$\mathbf{U} = \exp\left(t \mathbf{A} \mathbf{D}^{-1} - t\right), \tag{17}$$

where $\mathbf{D}$ is a diagonal degree matrix as in Section 3.1 and $t$ is diffusion time [2]. For Reddit dataset, we implement the corruption function $C$ by keeping the adjacency matrix $\mathbf{A}$ unchanged (i.e. $\tilde{\mathbf{A}} = \mathbf{A}$) and perturbing the feature matrix $\mathbf{X}$ by row-wise shuffling. And for PPI dataset, we simply samples a different graph from the training set due to it's a multiple-graph dataset.

**Readout function.** We use identical readout function $\mathcal{R}$ for all datasets, which performs a simple average of all node features for a given subgraph with $N'$ nodes:

$$\mathcal{R}(\mathbf{H}) = \sigma\left(\frac{1}{N'} \sum_{i=1}^{N'} \mathbf{h}_i\right), \tag{18}$$

where $\sigma(\cdot)$ is the logistic sigmoid nonlinearity.

**Discriminator.** We use a simple bilinear scoring function as discriminator $\mathcal{D}$:

$$\mathcal{D}(\mathbf{h}_i, \mathbf{h}_j) = \sigma(\mathbf{h}_i^T \mathbf{W}_d \mathbf{h}_j), \tag{19}$$

where $\mathbf{W}_d$ is a learnable scoring matrix and $\sigma(\cdot)$ is the logistic sigmoid nonlinearity.

**Table 5: The classification accuracy (in %) of different node-centered subgraphs combinations (*Subgraph 1*: Basic subgraph, *Subgraph 2*: Neighboring subgraph, *Subgraph 3*: Intimate subgraph, *Subgraph 4*: Communal subgraph, *Subgraph 5*: Full subgraph) on Cora, Citeseer and Pubmed datasets with same hyperparameter setting. Note that due to the use of the "core view" paradigm, there must be at least the basic subgraph and another subgraph. The best result for each dataset is indicated by bolded.**

| Subgraphs | | | | | Dataset | | |
|---|---|---|---|---|---|---|---|
| 1 | 2 | 3 | 4 | 5 | Cora | Citeseer | Pubmed |
| ✓ | ✓ | | | | 82.5 ± 0.7 | 72.2 ± 0.7 | 77.5 ± 0.1 |
| ✓ | | ✓ | | | 82.8 ± 0.3 | 72.5 ± 0.7 | 78.6 ± 0.9 |
| ✓ | | | ✓ | | 82.7 ± 0.6 | 72.1 ± 0.6 | 78.2 ± 0.4 |
| ✓ | | | | ✓ | 82.7 ± 0.5 | 71.8 ± 0.6 | 78.1 ± 0.6 |
| AVG of 2 Subgraphs | | | | | 82.7 ± 0.6 | 72.2 ± 0.7 | 78.2 ± 0.6 |
| ✓ | ✓ | ✓ | | | 83.1 ± 0.5 | 72.7 ± 0.9 | 78.3 ± 0.8 |
| ✓ | ✓ | | ✓ | | 82.9 ± 0.9 | 72.3 ± 0.5 | 78.0 ± 0.9 |
| ✓ | ✓ | | | ✓ | 83.2 ± 0.4 | 72.1 ± 0.4 | 77.6 ± 0.8 |
| ✓ | | ✓ | ✓ | | 82.9 ± 0.5 | 72.0 ± 0.7 | 78.6 ± 0.4 |
| ✓ | | ✓ | | ✓ | 83.0 ± 0.5 | 72.9 ± 0.5 | 78.7 ± 0.8 |
| ✓ | | | ✓ | ✓ | 83.1 ± 0.8 | 72.9 ± 0.6 | 78.5 ± 0.6 |
| AVG of 3 Subgraphs | | | | | 83.0 ± 0.6 | 72.5 ± 0.7 | 78.3 ± 0.8 |
| ✓ | ✓ | ✓ | ✓ | | 83.5 ± 0.4 | 72.6 ± 0.7 | 79.2 ± 0.4 |
| ✓ | ✓ | ✓ | | ✓ | 83.5 ± 0.4 | 73.0 ± 0.6 | 78.8 ± 0.9 |
| ✓ | ✓ | | ✓ | ✓ | 83.6 ± 0.6 | 72.4 ± 0.6 | 78.5 ± 0.8 |
| ✓ | | ✓ | ✓ | ✓ | 83.3 ± 0.7 | 73.1 ± 1.0 | 78.6 ± 0.6 |
| AVG of 4 Subgraphs | | | | | 83.5 ± 0.5 | 72.7 ± 0.8 | 78.7 ± 0.7 |
| ✓ | ✓ | ✓ | ✓ | ✓ | **84.7 ± 0.3** | **73.8 ± 0.5** | **81.5 ± 0.4** |

**Details of node-centered subgraphs.** We conduct experiments with all five node-centered subgraphs on two different contrastive losses (named MNCSCL-FG under full graph case and MNCSCL-CV under core view case, respectively). More specifically, we choose $d = 1$ for the neighboring subgraph and *Strategy 2* for the communal subgraph. For hyperparameters, we set the size of intimate subgraph $l$ as 20 (specially, 10 on citeseer). The number of clusters $C$ and
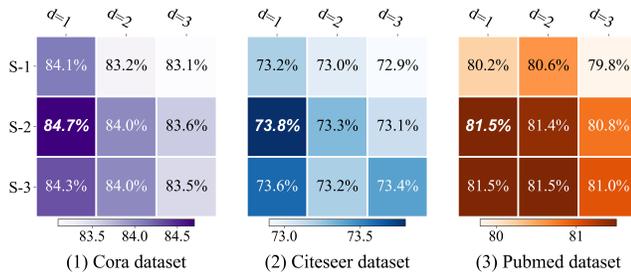
**Figure 5: Heat map of classification accuracy (in %) on three transductive datasets when choosing different range of neighbors in the neighboring subgraph ($d = 1$, $d = 2$ and $d = 3$) and different clustering strategies in the communal subgraph (S-1, S-2 and S-3, where S-1 denotes *Strategy 1*, S-2 denotes *Strategy 2* and S-3 denotes *Strategy 3*). The bolded value in each sub-plot represents the maximum classification accuracy for the current dataset.**

inverse-temperature hyperparameter $\beta$ for communal subgraph is set to 128 and 10 respectively. To build a proper full subgraph, we also set the self-weighted factor as 0.01.

## 4.3 Experimental Results and Analysis

**Node classification.** Experiments show that MNCSCL achieves the best performance on all five datasets compared to other competing self-supervised methods, as shown in Table 3. We believe this robust performance is due to our comparison of multiple node-centered subgraphs, resulting in learning a more comprehensive node representation. Although MNCSCL-FG and MNCSCL-CV both have shown excellent performance, MNCSCL-FG performs better on PPI dataset and MNCSCL-CV performs better on other datasets. We think this is due to the very sparse available features on the PPI (over 40% of the nodes have all-zero features). It is more effective to use more perspectives for comparison on such sparse graph datasets. For other datasets, the contrastive loss under core view case is enough for MNCSCL to learn comprehensive information about the nodes, too much comparison will instead lead to overfitting and waste of resources. Compared to supervised methods, MNCSCL also outperforms the two compared methods on all datasets. This shows that our method is also very competitive compared to traditional supervision methods.

**Link prediction.** To test the generalization capability of MNCSCL, we intend to further investigate the performance of MNCSCL in link prediction task, as shown in Table 4. We find that MNCSCL outperforms all competing methods on both the Cora and Citeseer datasets, suggesting that a multiple node-centered subgraphs based comparison can help the model learn node representations with good generalizability. The excellent performance on different downstream tasks further proves the feasibility of our method.

## 4.4 Ablation Study

**Node-centered subgraph combination.** To investigate how the number of node-centered subgraphs affects the performance of MNCSCL, we permuted and combined five previously mentioned

node-centered subgraphs under the core view case and observed the classification accuracy of different subgraph combinations on Cora, Citeseer and Pubmed datasets with same hyperparameter setting, as shown in Table 5. Obviously, as the number of node-centered subgraphs increases, the classification accuracy continues to increase, and the best results are achieved when all five subgraphs are used. This indicates that multiple node-centered subgraphs contrastive learning can indeed learn better node representation. We also notice that the classification accuracy improvement is more obvious with the increase in the number of node-centered subgraphs.

**Range of neighbors and clustering strategy.** We investigate the value of $d$ in the neighboring subgraph and the selection of different clustering strategies in the communal subgraph, and the results are shown in Fig 5. Here we use all five node-centered subgraphs with the contrastive loss under core view case. It can be seen that the optimal choice in all three datasets is $d = 1$ neighbors and *Strategy 2*. Our analyses are as follows.

- For the neighboring subgraph, the classification accuracy tends to decrease as the vale of $d$ increases. We believe that too many neighboring nodes will lead to cause overfitting of the model and performance degradation. This viewpoint can also be verified from the pubmed dataset, where it is not obvious whether the classification accuracy is better or worse when setting $d = 1$ and $d = 2$. Because the attributes of pubmed are relatively sparse, sometimes its neighboring subgraph needs to contain more neighbors to get better results.
- In the choice of clustering strategy, *Strategy 1* is significantly less effective than the other two end-to-end strategies. *Strategy 2* and *Strategy 3* show comparable performance, but considering that *Strategy 3* involves an estimation network, which means more resource consumption, we finally choose *Strategy 2* as the clustering strategy for the communal subgraph.

## 5 CONCLUSION

We propose Multiple Node-centered Subgraphs Contrastive Representation Learning (MNCSCL), a novel approach to self-supervised graph representation learning. MNCSCL obtains five different node-centered subgraphs carefully designed by us through a subgraph generator on each node, and maximizes the mutual information between them through two types of contrastive loss, thus allowing us to obtain comprehensive node representation that combines information from multiple node-centered subgraphs of nodes. Experiments show that MNCSCL reach the advanced level of self-supervised learning in both transductive and inductive node classification tasks as well as in link prediction task.

# REFERENCES

[1] Jie Chen, Tengfei Ma, and Cao Xiao. 2018. Fastgcn: fast learning with graph convolutional networks via importance sampling. *arXiv preprint arXiv:1801.10247* (2018).

[2] Johannes Gasteiger, Stefan Weißenberger, and Stephan Günnemann. 2019. Diffusion improves graph learning. *Advances in neural information processing systems* 32 (2019).

[3] Aditya Grover and Jure Leskovec. 2016. node2vec: Scalable feature learning for networks. In *Proceedings of the 22nd ACM SIGKDD international conference on Knowledge discovery and data mining*. 855–864.

[4] Will Hamilton, Zhitao Ying, and Jure Leskovec. 2017. Inductive representation learning on large graphs. *Advances in neural information processing systems* 30 (2017).

[5] William L Hamilton, Rex Ying, and Jure Leskovec. 2017. Representation learning on graphs: Methods and applications. *arXiv preprint arXiv:1709.05584* (2017).

[6] Kaveh Hassani and Amir Hosein Khasahmadi. 2020. Contrastive multi-view representation learning on graphs. In *International Conference on Machine Learning*. PMLR, 4116–4126.

[7] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2015. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *Proceedings of the IEEE international conference on computer vision*. 1026–1034.

[8] R Devon Hjelm, Alex Fedorov, Samuel Lavoie-Marchildon, Karan Grewal, Phil Bachman, Adam Trischler, and Yoshua Bengio. 2018. Learning deep representations by mutual information estimation and maximization. *arXiv preprint arXiv:1808.06670* (2018).

[9] Glen Jeh and Jennifer Widom. 2003. Scaling personalized web search. In *Proceedings of the 12th international conference on World Wide Web*. 271–279.

[10] Yizhu Jiao, Yun Xiong, Jiawei Zhang, Yao Zhang, Tianqi Zhang, and Yangyong Zhu. 2020. Sub-graph contrast for scalable self-supervised graph representation learning. In *2020 IEEE international conference on data mining (ICDM)*. IEEE, 222–231.

[11] Longlong Jing and Yingli Tian. 2020. Self-supervised visual feature learning with deep neural networks: A survey. *IEEE transactions on pattern analysis and machine intelligence* 43, 11 (2020), 4037–4058.

[12] Nikhil Ketkar and Jojo Moolayil. 2021. Introduction to pytorch. In *Deep learning with python*. Springer, 27–91.

[13] Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014).

[14] Thomas N Kipf and Max Welling. 2016. Semi-supervised classification with graph convolutional networks. *arXiv preprint arXiv:1609.02907* (2016).

[15] Thomas N Kipf and Max Welling. 2016. Variational graph auto-encoders. *arXiv preprint arXiv:1611.07308* (2016).

[16] Risi Imre Kondor and John Lafferty. 2002. Diffusion kernels on graphs and other discrete structures. In *Proceedings of the 19th international conference on machine learning*, Vol. 2002. 315–322.

[17] Yingming Li, Ming Yang, and Zhongfei Zhang. 2018. A survey of multi-view representation learning. *IEEE transactions on knowledge and data engineering* 31, 10 (2018), 1863–1883.

[18] J MacQueen. 1967. Classification and analysis of multivariate observations. In *5th Berkeley Symp. Math. Statist. Probability*. 281–297.

[19] Costas Mavromatis and George Karypis. 2020. Graph infoclust: Leveraging cluster-level node information for unsupervised graph representation learning. *arXiv preprint arXiv:2009.06946* (2020).

[20] Yujie Mo, Liang Peng, Jie Xu, Xiaoshuang Shi, and Xiaofeng Zhu. 2022. Simple unsupervised graph representation learning. AAAI.

[21] Shirui Pan, Ruiqi Hu, Sai-fu Fung, Guodong Long, Jing Jiang, and Chengqi Zhang. 2019. Learning graph embedding with adversarial training methods. *IEEE transactions on cybernetics* 50, 6 (2019), 2475–2487.

[22] Qiyao Peng, Hongtao Liu, Yinghui Wang, Hongyan Xu, Pengfei Jiao, Minglai Shao, and Wenjun Wang. 2022. Towards a multi-view attentive matching for personalized expert finding. In *Proceedings of the ACM Web Conference 2022*. 2131–2140.

[23] Zhen Peng, Wenbing Huang, Minnan Luo, Qinghua Zheng, Yu Rong, Tingyang Xu, and Junzhou Huang. 2020. Graph representation learning via graphical mutual information maximization. In *Proceedings of The Web Conference 2020*. 259–270.

[24] Bryan Perozzi, Rami Al-Rfou, and Steven Skiena. 2014. Deepwalk: Online learning of social representations. In *Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining*. 701–710.

[25] Jiezhong Qiu, Qibin Chen, Yuxiao Dong, Jing Zhang, Hongxia Yang, Ming Ding, Kuansan Wang, and Jie Tang. 2020. Gcc: Graph contrastive coding for graph neural network pre-training. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 1150–1160.

[26] Satu Elisa Schaeffer. 2007. Graph clustering. *Computer science review* 1, 1 (2007), 27–64.

[27] Lei Tang and Huan Liu. 2011. Leveraging social media networks for classification. *Data Mining and Knowledge Discovery* 23, 3 (2011), 447–478.

[28] Yonglong Tian, Dilip Krishnan, and Phillip Isola. 2020. Contrastive multiview coding. In *European conference on computer vision*. Springer, 776–794.

[29] Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, and Yoshua Bengio. 2017. Graph attention networks. *arXiv preprint arXiv:1710.10903* (2017).

[30] Petar Velickovic, William Fedus, William L Hamilton, Pietro Liò, Yoshua Bengio, and R Devon Hjelm. 2019. Deep Graph Infomax. *ICLR (Poster)* 2, 3 (2019), 4.

[31] Bryan Wilder, Eric Ewing, Bistra Dilkina, and Milind Tambe. 2019. End to end learning and optimization on graphs. *Advances in Neural Information Processing Systems* 32 (2019).

[32] Felix Wu, Amauri Souza, Tianyi Zhang, Christopher Fifty, Tao Yu, and Kilian Weinberger. 2019. Simplifying graph convolutional networks. In *International conference on machine learning*. PMLR, 6861–6871.

[33] Lirong Wu, Haitao Lin, Cheng Tan, Zhangyang Gao, and Stan Z Li. 2021. Self-supervised learning on graphs: Contrastive, generative, or predictive. *IEEE Transactions on Knowledge and Data Engineering* (2021).

[34] Minghao Xu, Hang Wang, Bingbing Ni, Hongyu Guo, and Jian Tang. 2021. Self-supervised graph-level representation learning with local and global structure. In *International Conference on Machine Learning*. PMLR, 11548–11558.

[35] Jiawei Zhang and Lin Meng. 2019. Gresnet: Graph residual network for reviving deep gnns from suspended animation. *arXiv preprint arXiv:1909.05729* (2019).

[36] Jiawei Zhang, Haopeng Zhang, Congying Xia, and Li Sun. 2020. Graph-bert: Only attention is needed for learning graph representations. *arXiv preprint arXiv:2001.05140* (2020).

[37] Chen Zhao. 2021. *Fairness-Aware Multi-Task and Meta Learning*. Ph. D. Dissertation.

[38] Chen Zhao and Feng Chen. 2019. Rank-based multi-task learning for fair regression. In *2019 IEEE International Conference on Data Mining (ICDM)*. IEEE, 916–925.

[39] Chen Zhao, Feng Chen, and Bhavani Thuraisingham. 2021. Fairness-aware online meta-learning. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*. 2294–2304.

[40] Jun Zhao, Xudong Liu, Qiben Yan, Bo Li, Minglai Shao, and Hao Peng. 2020. Multi-attributed heterogeneous graph convolutional network for bot detection. *Information Sciences* 537 (2020), 380–393.

[41] Yanqiao Zhu, Yichen Xu, Feng Yu, Qiang Liu, Shu Wu, and Liang Wang. 2020. Deep graph contrastive representation learning. *arXiv preprint arXiv:2006.04131* (2020).

[42] Thomas Zimmermann and Nachiappan Nagappan. 2008. Predicting defects using network analysis on dependency graphs. In *Proceedings of the 30th international conference on Software engineering*. 531–540.