# From Canonical Correlation Analysis to Self-supervised Graph Neural Networks

**Hengrui Zhang**[1]*, **Qitian Wu**[2], **Junchi Yan**[2], **David Wipf**[3], **Philip S. Yu**[1]†

[1] Department of Computer Science, University of Illinois at Chicago
[2] Department of Computer Science and Engineering, Shanghai Jiao Tong University
[3] AWS Shanghai AI Lab
hzhan55@uic.edu, {echo740, yanjunchi}@sjtu.edu.cn
daviwipf@amazon.com, psyu@uic.edu

## Abstract

We introduce a conceptually simple yet effective model for self-supervised representation learning with graph data. It follows the previous methods that generate two views of an input graph through data augmentation. However, unlike contrastive methods that focus on instance-level discrimination, we optimize an innovative feature-level objective inspired by classical Canonical Correlation Analysis. Compared with other works, our approach requires none of the parameterized mutual information estimator, additional projector, asymmetric structures, and most importantly, negative samples which can be costly. We show that the new objective essentially 1) aims at discarding augmentation-variant information by learning invariant representations, and 2) can prevent degenerated solutions by decorrelating features in different dimensions. Our theoretical analysis further provides an understanding for the new objective which can be equivalently seen as an instantiation of the Information Bottleneck Principle under the self-supervised setting. Despite its simplicity, our method performs competitively on seven public graph datasets. The code is available at: `https://github.com/hengruizhang98/CCA-SSG`.

## 1 Introduction

Self-supervised learning (SSL) has been a promising paradigm for learning useful representations without costly labels [7, 46, 5]. In general, it learns representations via a proxy objective between inputs and self-defined signals, among which contrastive methods [46, 40, 16, 5, 12] have achieved impressive performance on learning image representations by maximizing the mutual information of two views (or augmentations) of the same input. Such methods can be interpreted as a discrimination of a joint distribution (positive pairs) from the product of two marginal ones (negative pairs) [50].

Inspired by the success of contrastive learning in vision [17, 46, 40, 5, 16, 12, 6], similar methods have been adapted to learning graph neural networks [48, 15, 33, 57, 58]. Although these models have achieved impressive performance, they require complex designs and architectures. For example, DGI [48] and MVGRL [15] rely on a parameterized mutual information estimator to discriminate positive node-graph pairs from negative ones; GRACE [57] and GCA [58] harness an additional MLP-projector to guarantee sufficient capacity. Moreover, negative pairs sampled or constructed from data often play an indispensable role in providing effective contrastive signals and have a large impact on performance. Selecting proper negative samples is often nontrivial for graph-structured data, not to mention the extra storage cost for prohibitively large graphs. BGRL [39] is a recent endeavor on

---

*This work was done during the author's internship at AWS Shanghai AI Lab.
†Corresponding author.

Table 1: Technical comparison of self-supervised node representation learning methods. We provide a conceptual comparison with more self-supervised methods in Appendix G. *Target* denotes the comparison pair, N/G/F denotes node/graph/feature respectively. *MI-Estimator*: parameterized mutual information estimator. *Proj/Pred*: additional (MLP) projector or predictor. *Asymmetric*: asymmetric architectures such as EMA and Stop-Gradient, or two separate encoders for two branches. *Neg examples*: requiring negative examples to prevent trivial solutions. *Space* denotes space requirement for storing all the pairs. Our method is simple without any listed component and memory-efficient.

| | Methods | Target | MI-Estimator | Proj/Pred | Asymmetric | Neg examples | Space |
|---|---|---|---|---|---|---|---|
| Instance-level | DGI [48] | N-G | ✓ | - | - | ✓ | $O(N)$ |
| | MVGRL [15] | N-G | ✓ | - | ✓ | ✓ | $O(N)$ |
| | GRACE [57] | N-N | - | ✓ | - | ✓ | $O(N^2)$ |
| | GCA [58] | N-N | - | ✓ | - | ✓ | $O(N^2)$ |
| | BGRL [39] | N-N | - | ✓ | ✓ | - | $O(N)$ |
| | CCA-SSG (Ours) | F-F | - | - | - | - | $O(D^2)$ |

targeting a negative-sample-free approach for GNN learning through asymmetric architectures [12, 6]. However, it requires additional components, e.g., an exponential moving average (EMA) and Stop-Gradient, to empirically avoid degenerated solutions, leading to a more intricate architecture.

Deviating from the large body of previous works on contrastive learning, in this paper we take a new perspective to address SSL on graphs. We introduce Canonical Correlation Analysis inspired Self-Supervised Learning on Graphs (CCA-SSG), a simple yet effective approach that opens the way to a new SSL objective and frees the model from intricate designs. It follows the common practice of prior arts, generating two views of an input graph through random augmentation and acquiring node representations through a shared GNN encoder. Differently, we propose to harness a *non-contrastive* and *non-discriminative* feature-level objective, which is inspired by the well-studied Canonical Correlation Analysis (CCA) methods [18, 10, 11, 14, 2, 4]. More specifically, the new objective aims at maximizing the correlation between two augmented views of the same input and meanwhile decorrelating different (feature) dimensions of a single view's representation. We show that the objective 1) essentially pursuits discarding augmentation-variant information and preserving augmentation-invariant information, and 2) can prevent dimensional collapse [19] (i.e., different dimensions capture the same information) in nature. Furthermore, our theoretical analysis sheds more lights that under mild assumptions, our model is an instantiation of Information Bottleneck Principle [43, 44, 37] under SSL settings [53, 9, 45].

To sum up, as shown in Table 1, our new objective induces a simple and light model without reliance on negative pairs [48, 15, 57, 58], a parameterized mutual information estimator [48, 15], an additional projector or predictor [57, 58, 39] or asymmetric architectures [39, 15]. We provide a thorough evaluation for the model on seven node classification benchmarks. The empirical results demonstrate that despite its simplicity, CCA-SSG can achieve very competitive performance in general and even superior test accuracy in five datasets. It is worth noting that our approach is agnostic to the input data format, which means that it can potentially be applied to other scenarios beyond graph-structured data (such as vision, language, etc.). We leave such a technical extension for future works.

**Our contributions are as follows:**

**1)** We introduce a non-contrastive and non-discriminative objective for self-supervised learning, which is inspired by Canonical Correlation Analysis methods. It does not rely on negative samples, and can naturally remove the complicated components. Based on it we propose CCA-SSG, a simple yet effective framework for learning node representations without supervision (see Section 3).

**2)** We theoretically prove that the proposed objective aims at keeping augmentation-invariant information while discarding augmentation-variant one, and possesses an inherent relationship to an embodiment of Information Bottleneck Principle under self-supervised settings (see Section 4).

**3)** Experimental results show that without complex designs, our method outperforms state-of-the-art self-supervised methods MVGRL [15] and GCA [58] on 5 out of 7 benchmarks. We also provide thorough ablation studies on the effectiveness of the key components of CCA-SSG (see Section 5).

## 2   Related Works and Background

**Contrastive Learning on Graphs.**   Contrastive methods [46, 40, 17, 16, 5, 12] have been shown to be effective for unsupervised learning in vision, which have also been adapted to graphs. Inspired by the local-global mutual information maximization viewpoints [17], DGI [48] and InfoGraph [38] put forward unsupervised schemes for node and graph representation learning, respectively. MVGRL [15] generalizes CMC [40] to graph-structured data by introducing graph diffusion [23] to create another view for a graph. GCC [33] adopts InfoNCE loss [46] and MoCo-based negative pool [16] for large-scale GNN pretraining. GRACE [57], GCA [58] and GraphCL [52] follow the spirit of SimCLR [5] and learn node/graph representations by directly treating other nodes/graphs as negative samples. BGRL [39] targets a negative-sample-free model, inspired by BYOL [12], on node representation learning. But it still requires complex asymmetric architectures.

**Feature-level Self-supervised Objectives.**   The above-mentioned methods all focus on instance-level contrastive learning. To address their drawbacks, some recent works have been turning to feature-level objectives. For example, Contrastive Clustering [25] regards different feature dimensions as different clusters, thus combining the cluster-level discrimination with instance-level discrimination. W-MSE [8] performs a differentiable whitening operation on learned embeddings, which implicitly scatters data points in embedding space. Barlow Twins [53] borrows the idea of redundancy reduction and adopts a soft decorrelation term that makes the cross-correlation matrix of two views' representations close to an identity matrix. By contrast, our method is based on the classical Canonical Correlation Analysis, working by correlating the representations of two views from data augmentation and meanwhile decorrelating different feature dimensions of each view's representation.

**Canonical Correlation Analysis.**   CCA is a classical multivariate analysis method, which is first introduced in [18]. For two random variables $X \in \mathbb{R}^m$ and $Y \in \mathbb{R}^n$, their covariance matrix is $\Sigma_{XY} = Cov(X, Y)$. CCA aims at seeking two vectors $a \in \mathbb{R}^m$ and $b \in \mathbb{R}^n$ such that the correlation $\rho = \mathrm{corr}(a^\top X, b^\top Y) = \frac{a^\top \Sigma_{XY} b}{\sqrt{a^\top \Sigma_{XX} a}\sqrt{b^\top \Sigma_{YY} b}}$ is maximized. Formally, the objective is

$$\max_{a,b} a^\top \Sigma_{XY} b, \text{ s.t. } a^\top \Sigma_{XX} a = b^\top \Sigma_{YY} b = 1. \tag{1}$$

For multi-dimensional cases, CCA seeks two sets of vectors maximizing their correlation and subjected to the constraint that they are uncorrelated with each other [10]. Later studies apply CCA to multi-view learning with deep models [2, 11, 14], by replacing the linear transformation with neural networks. Concretely, assuming $X_1, X_2$ as two views of an input data, it optimizes

$$\max_{\theta_1,\theta_2} \mathrm{Tr}\left( P_{\theta_1}^\top(X_1) P_{\theta_2}(X_2) \right) \text{ s.t. } P_{\theta_1}^\top(X_1) P_{\theta_1}(X_1) = P_{\theta_2}^\top(X_2) P_{\theta_2}(X_2) = I. \tag{2}$$

where $P_{\theta_1}$ and $P_{\theta_2}$ are two feedforward neural networks and $I$ is an identity matrix. Despite its preciseness, such computation is really expensive [4]. Fortunately, soft CCA [4] removes the hard decorrelation constraint by adopting the following Lagrangian relaxation:

$$\min_{\theta_1,\theta_2} \mathcal{L}_{dist}\left( P_{\theta_1}(X_1), P_{\theta_2}(X_2) \right) + \lambda \left( \mathcal{L}_{SDL}(P_{\theta_1}(X_1)) + \mathcal{L}_{SDL}(P_{\theta_2}(X_2)) \right), \tag{3}$$

where $\mathcal{L}_{dist}$ measures correlation between two views' representations and $\mathcal{L}_{SDL}$ (called stochastic decorrelation loss) computes an $L_1$ distance between $P_{\theta_i}(X_i)$ and an identity matrix, for $i = 1, 2$.

## 3   Approach

### 3.1   Model Framework

In this paper we focus on self-supervised node representation learning, where we consider a single graph $\mathbf{G} = (\mathbf{X}, \mathbf{A})$. $\mathbf{X} \in \mathbb{R}^{N \times F}$ and $\mathbf{A} \in \mathbb{R}^{N \times N}$ denote node features and adjacency matrix respectively. Here $N$ is the number of nodes within the graph and $F$ denotes feature dimension.

Our model simply consists of three parts: 1) a random graph augmentation generator $\mathcal{T}$. 2) a GNN-based graph encoder $f_\theta$ where $\theta$ denotes its parameters. 3) a novel feature-level objective function based on Canonical Correlation Analysis. Fig. 1 is an illustration of the proposed model.
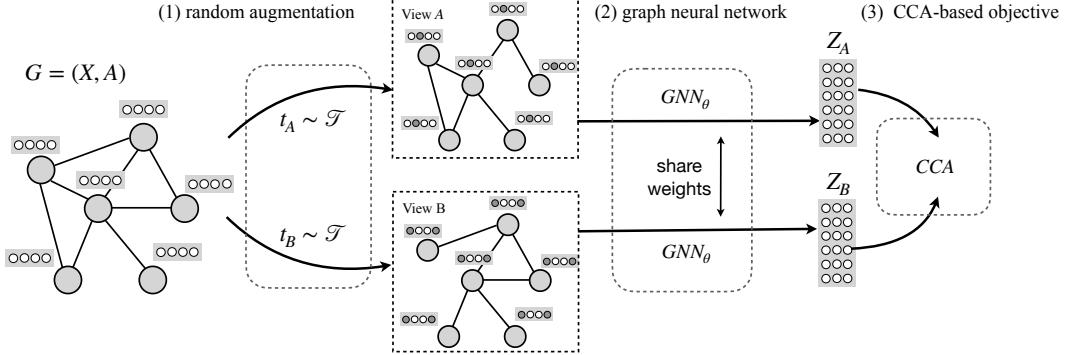
3

Figure 1: Illustration of the proposed model: given an input graph, we first generate two views through random augmentations: edge dropping and node feature masking. The two views are subsequently put into a shared GNN encoder to generate representations. The loss function is applied on the column-normalized embedding matrix of the two views. Note that this simple yet effective pipeline can also be conceptually applied for other data like vision and texts, which we leave for future works.

---

**Algorithm 1:** PyTorch-style code for CCA-SSG

```python
# f: encoder network
# lambda: trade-off
# D: embedding dimension
# g: input graph
# feat: node features

# generate two views through random augmentation
g1, feat1 = augment(g, feat)
g2, feat2 = augment(g, feat)
z1 = f(g1, feat1) # embedding of the 1st view
z2 = f(g2, feat2) # embedding of the 2st view

# batch normalization
z1_norm = ((z1 - z1.mean(0)) / z1.std(0))/ sqrt(D)
z2_norm = ((z2 - z2.mean(0)) / z2.std(0))/ sqrt(D)

# covariance matrix of each view
c1 = torch.mm(z1_norm.T(), z1_norm)
c2 = torch.mm(z2_norm.T(), z2_norm)

iden = torch.eye(D)
loss_inv = (z1_norm - z2_norm).pow(2).sum()
loss_dec_1 = (c1 - iden).pow(2).sum()
loss_dec_2 = (c2 - iden).pow(2).sum()
loss_dec = loss_dec_1 + loss_dec_2
loss = loss_inv + lambda * loss_dec
```

**Graph augmentations.** We consider the standard pipeline for random graph augmentation that has been commonly used in previous works [57, 39]. To be specific, we harness two ways for augmentation: **edge dropping** and **node feature masking**. Edge dropping randomly drops a fraction of edges from the original graph, while node feature masking randomly masks a fraction of features for all the nodes. In this way, $\mathcal{T}$ is composed of all the possible graph transformation operations and each $t \sim \mathcal{T}$ denotes a specific graph transformation for graph $G$. Note that we use commonly adopted augmentation methods to stay our focus on the design of objective function and conduct fair comparison with existing approaches. More complicated random augmentations [52, 58] can also be readily plugged into our model. Details for the used augmentation functions are in Appendix E.

**Training.** In each training iteration, we first randomly sample two graph transformations $t_A$ and $t_B$ from $\mathcal{T}$, and then generate two views $\tilde{\mathbf{G}}_A = (\tilde{\mathbf{X}}_A, \tilde{\mathbf{A}}_A)$ and $\tilde{\mathbf{G}}_B = (\tilde{\mathbf{X}}_B, \tilde{\mathbf{A}}_B)$ according to the transformations. The two views are subsequently fed into a shared GNN encoder to generate the node embeddings of the two views: $\mathbf{Z}_A = f_\theta(\tilde{\mathbf{X}}_A, \tilde{\mathbf{A}}_A)$, $\mathbf{Z}_B = f_\theta(\tilde{\mathbf{X}}_B, \tilde{\mathbf{A}}_B)$, where $\mathbf{Z}_A, \mathbf{Z}_B \in \mathbb{R}^{N \times D}$ and $D$ denotes embedding dimension. We further normalize the node embeddings along instance dimension so that each feature dimension has a 0-mean and $1/\sqrt{N}$-standard deviation distribution:

$$\tilde{\mathbf{Z}} = \frac{\mathbf{Z} - \mu(\mathbf{Z})}{\sigma(\mathbf{Z}) * \sqrt{N}} \tag{4}$$

The normalized $\tilde{\mathbf{Z}}_A$, $\tilde{\mathbf{Z}}_B$ will be used to compute a feature-level objective in Section 3.2. To help better understand the proposed framework, we provide the PyTorch-style pseudocode for training CCA-SSG in Algorithm 1.

**Inference.** To generate node embeddings for downstream tasks, we put the original graph $\mathbf{G} = (\mathbf{X}, \mathbf{A})$ into the trained graph neural network $f_\theta$ and obtain node embeddings $\mathbf{Z} = f_\theta(\mathbf{X}, \mathbf{A})$.

## 3.2 Learning Objective

Canonical Correlation Analysis has shown its great power in multi-view learning like instance recognition [4]. However, it still remains unexplored to leverage CCA for self-supervised learning. Note that in SSL, one generates two sets of data from the same input through transformation or random data augmentation, which could be regraded as two views of the input data. This inspires us to introduce the following objective for self-supervised representation learning:

$$\mathcal{L} = \underbrace{\left\| \tilde{\mathbf{Z}}_A - \tilde{\mathbf{Z}}_B \right\|_F^2}_{\text{invariance term}} + \lambda \underbrace{\left( \left\| \tilde{\mathbf{Z}}_A^\top \tilde{\mathbf{Z}}_A - \mathbf{I} \right\|_F^2 + \left\| \tilde{\mathbf{Z}}_B^\top \tilde{\mathbf{Z}}_B - \mathbf{I} \right\|_F^2 \right)}_{\text{decorrelation term}} \tag{5}$$

where $\lambda$ is a non-negative hyperparameter trading off two terms. Note that minimizing the invariance term is essentially maximizing the correlation between two views as their representations are already normalized. In SSL, as the two augmented views come randomly from the same distribution, we can adopt one encoder $f_\theta$ that is shared across two branches and seek for a regularization that encourages different feature dimensions to capture distinct semantics via the decorrelation term.

We next provide a variance-covariance perspective to the new objective, following similar lines of reasoning in [41, 42]. Assume that input data come from a distribution $x \sim p(x)$ and $s$ is a view of $x$ through random augmentation $s \sim p_{aug}(\cdot|x)$. Denote $z_s$ as the representation of $s$, then minimizing the invariance term, by expectation, is to minimize the variance of the normalized representation $\tilde{z}_s$, conditioned on $x$. Also, minimizing the decorrelation term is to push the off-diagonal elements of the covariance matrix (given by two $\tilde{z}_s$'s) close to $0$. Formally, we have

$$\mathcal{L}_{inv} = \left\| \tilde{\mathbf{Z}}_A - \tilde{\mathbf{Z}}_B \right\|_F^2 = \sum_{i=1}^{N} \sum_{k=1}^{D} (\tilde{z}_{i,j}^A - \tilde{z}_{i,j}^B)^2 \cong \mathbb{E}_x \left[ \sum_{k=1}^{D} \mathbb{V}_{s|x}[\tilde{z}_{s,k}] \right] * 2N, \tag{6}$$

$$\mathcal{L}_{dec} = \left\| \tilde{\mathbf{Z}}_S^\top \tilde{\mathbf{Z}}_S - \mathbf{I} \right\|_F^2 = \| \mathbf{Cov}_s[\tilde{z}] - I \|_F^2 \cong \sum_{i \neq j} \left( \rho_{i,j}^{z_s} \right)^2, \text{ for } \tilde{\mathbf{Z}}_S \in \{ \tilde{\mathbf{Z}}_A, \tilde{\mathbf{Z}}_B \}, \tag{7}$$

where $\rho$ is the Pearson correlation coefficient.

## 3.3 Advantages over Contrastive Methods

In this subsection we provide a systematic comparison with previous self-supervised methods for node representation learning, including DGI [48], MVGRL [15], GRACE [57], GCA [58] and BGRL [39], and highlight the merits of CCA-SSG. A quick overview is presented in Table 1.

**No reliance on negative samples**. Most of previous works highly rely on negative pairs to avoid collapse or interchangeable, trivial/degenerated solutions [48, 15, 57, 58]. E.g., DGI and MVGRL generate negative examples by corrupting the graph structure severely, and GRACE/GCA treats all the other nodes within a graph as negative examples. However, for self-supervised learning on graphs, it is non-trivial to construct informative negative examples since nodes are structurally connected, and selecting negative examples in an arbitrary manner may lead to large variance for stochastic gradients and slow training convergence [51]. The recently proposed BGRL model adopts asymmetric encoder architectures for SSL on graphs without the use of negative samples. However, though BGRL could avoid collapse empirically, it still remains as an open problem concerning its theoretical guarantee for preventing trivial solutions [41]. Compared with these methods, our model does not rely on negative pairs and asymmetric encoders. The feature decorrelation term can naturally prevent trivial solutions caused by the invariance term. We discuss the collapse issue detailedly in Appendix B.

**No MI estimator, projector network nor asymmetric architectures**. Most previous works rely on additional components besides the GNN encoder to estimate some score functions in final objectives. DGI and MVGRL require a parameterized estimator to approximate mutual information between two views, and GRACE leverages a MLP projector followed by an InfoNCE estimator. BGRL harnesses asymmetric encoder architecture which consists of EMA (Exponential Moving Average), Stop-Gradient and an additional projector. MVGRL also induces asymmetric architectures as it adopts two different GNNs for the input graph and the diffusion graph respectively. In contrast, our approach requires no additional components except a single GNN encoder.

**Better efficiency and scalability to large graphs**. Consider a graph with $N$ nodes. DGI and MVGRL contrast node embeddings with graph embedding, which would require $O(N)$ space cost.

GRACE treats two views of the same node as positive pairs and treat views of different nodes as negative pairs, which would take $O(N^2)$ space. BGRL focuses only on positive pairs, which will also take $O(N)$ space. By contrast, our method works on feature dimension. If we embed each node into a $D$-dimensional vector, the computation of the loss function would require $O(D^2)$ space. This indicates that the memory cost does not grow consistently as the size of graph increases. As a result, our method is promising for handling large-scale graphs without prohibitively large space costs.

## 4    Theoretical Insights with Connection to Information Theory

In this section we provide some analysis of the proposed objective function: 1) Interpretation of the loss function with entropy and mutual information. 2) The connection between the proposed objective and the Information Bottleneck principle. 3) Why the learned representations would be informative to downstream tasks. The proofs of propositions, theorems and corollaries are in Appendix D.

**Notations.** Denote the random variable of input data as $X$ and the downstream task as $T$ (it could be the label $Y$ if the downstream task is classification). Note that in SSL, we have no access to $T$ in training and here we introduce the notation for our analysis. Define $S$ as the self-supervised signal (i.e., an augmented view of $X$), and $S$ shares the same space as $X$. Our model learns a representation for the input, denoted by $Z_X$ and its views, denoted by $Z_S$. $Z_X = f_\theta(X), Z_S = f_\theta(S)$, $f_\theta(\cdot)$ is a encoder shared by the original data and its views, which is parameterized by $\theta$. The target of representation learning is to learn a optimal encoder parameter $\theta$. Furthermore, for random variable $A, B, C$, we use $I(A, B)$ to denote the mutual information between $A$ and $B$, $I(A, B|C)$ to denote conditional mutual information of $A$ and $B$ on a given $C$, $H(A)$ for the entropy, and $H(A|B)$ for conditional entropy. The proofs of propositions, theorems and corollaries are in Appendix D.

### 4.1    An Entropy and Mutual Information Interpretation of the Objective

We first introduce an assumption about the distributions of $P(Z_S)$ and $P(Z_S|X)$.

**Assumption 1.** *(Gaussian assumption of $P(Z_S|X)$ and $P(Z_S)$):*

$$P(Z_S|X) = \mathcal{N}(\mu_X, \Sigma_X), P(Z_S) = \mathcal{N}(\mu, \Sigma). \tag{8}$$

With Assumption 1, we can arrive at the following propositions:

**Proposition 1.** *In expectation, minimizing Eq.* (6) *is equivalent to minimizing the entropy of $Z_S$ conditioned on input $X$, i.e.,*

$$\min_\theta \mathcal{L}_{inv} \cong \min_\theta H(Z_S|X). \tag{9}$$

**Proposition 2.** *Minimizing Eq.* (7) *is equivalent to maximizing the entropy of $Z_S$, i.e.,*

$$\min_\theta \mathcal{L}_{dec} \cong \max_\theta H(Z_S). \tag{10}$$

The two propositions unveil the effects of two terms in our objective. Combining two propositions, we can further interpret Eq. (5) from an information-theoretic perspective.

**Theorem 1.** *By optimizing Eq* (5)*, we maximize the mutual information between the augmented view's embedding $Z_S$ and the input data $X$, and minimize the mutual information between $Z_S$ and the view itself $S$, conditioned on the input data $X$. Formally we have*

$$\min_\theta \mathcal{L} \Rightarrow \max_\theta I(Z_S, X) \text{ and } \min_\theta I(Z_S, S|X). \tag{11}$$

The proof is based on the facts $I(Z_S, X) = H(Z_S) - H(Z_S|X)$ and $I(Z_S, S|X) = H(Z_S|X) + H(Z_S|S) = H(Z_S|X)$. Theorem 1 indicates that our objective Eq. (5) learns representations that maximize the information of the input data, i.e., $I(Z_S, X)$, and meanwhile minimize the lost information during augmentation, i.e., $I(Z_S, S|X)$.

### 4.2    Connection with the Information Bottleneck Principle

The analysis in Section 4.1 enables us to further build a connection between our objective Eq. (5) and the well-studied Information Bottleneck Principle [43, 44, 37, 1] under SSL settings. Recall that the supervised Information Bottleneck (IB) is defined as follows:

**Definition 1.** *The supervised IB aims at maximizing an Information Bottleneck Lagrangian:*

$$\mathcal{IB}_{sup} = I(Y, Z_X) - \beta I(X, Z_X), \text{ where } \beta > 0. \tag{12}$$

As we can see, $\mathcal{IB}_{sup}$ attempts to maximize the information between the data representation $Z_X$ and its corresponding label $Y$, and concurrently minimize the information between $Z_X$ and the input data $X$ (i.e., exploiting compression of $Z_X$ from $X$). The intuition of IB principle is that $Z_X$ is expected to contain only the information that is useful for predicting $Y$.

Several recent works [9, 45, 53] propose various forms of IB under self-supervised settings. The most relevant one names Self-supervised Information Bottleneck:

**Definition 2.** *(Self-supervised Information Bottleneck [53]). The Self-supervised IB aims at maximizing the following Lagrangian:*

$$\mathcal{IB}_{ssl} = I(X, Z_S) - \beta I(S, Z_S), \text{ where } \beta > 0. \tag{13}$$

Intuitively, $\mathcal{IB}_{ssl}$ posits that a desirable representation is expected to be informative to augmentation invariant features, and to be a maximally compressed representation of the input.

Our objective Eq. (5) is essentially an embodiment of $\mathcal{IB}_{ssl}$:

**Theorem 2.** *Assume $0 < \beta \leq 1$, then by minimizing Eq. (5), the self-supervised Information Bottleneck objective is maximized, formally:*

$$\min_{\theta} \mathcal{L} \Rightarrow \max_{\theta} \mathcal{IB}_{ssl} \tag{14}$$

Theorem 2 also shows that Eq. (5) implicitly follows the same spirit of IB principle under self-supervised settings. As further enlightenment, we can relate Eq. (5) with the *multi-view Information Bottleneck* [9] and the *minimal and sufficient representations for self-supervision* [45]:

**Corollary 1.** *Let $X_1 = S$, $X_2 = X$ and assume $0 < \beta \leq 1$, then minimizing Eq. (5) is equivalent to minimizing the Multi-view Information Bottleneck Loss in [9]:*

$$\mathcal{L}_{MIB} = I(Z_1, X_1 | X_2) - \beta I(X_2, Z_1), \text{ where } 0 < \beta \leq 1. \tag{15}$$

**Corollary 2.** *When the data augmentation process is reversible, minimizing Eq. (5) is equivalent to learning the Minimal and Sufficient Representations for Self-supervision in [45]:*

$$Z_X^{ssl} = \arg\max_{Z_X} I(Z_X, S), Z_X^{ssl_{min}} = \arg\min_{Z_X} H(Z_X | S) \text{ s.t. } I(Z_X, S) \text{ is maximized.} \tag{16}$$

### 4.3 Influence on Downstream Tasks

We have provided a principled understanding for our new objective. Next, we discuss its effect on downstream tasks $T$. The rationality of data augmentations in SSL is rooted in a conjecture that an ideal data augmentation approach would not change the information related to its label. We formulate this hypothesis as a building block for analysis on downstream tasks [36, 9].

**Assumption 2.** *(Task-relevant information and data augmentation). All the task-relevant information is shared across the input data $X$ and its augmentations $S$, i.e., $I(X, T) = I(S, T) = I(X, S, T)$, or equivalently, $I(X, T | S) = I(S, T | X) = 0$.*

This indicates that all the task-relevant information is contained in augmentation invariant features. We proceed to derive the following theorem which reveals the efficacy of the learned representations by our objective with respect to downstream tasks.

**Theorem 3.** *(Task-relevant/irrelevant information). By optimizing Eq. (5), the task-relevant information $I(Z_S, T)$ is maximized, and the task-irrelevant information $H(Z_S | T)$ is minimized. Formally,*

$$\min_{\theta} \mathcal{L} \Rightarrow \max_{\theta} I(Z_S, T) \text{ and } \min_{\theta} H(Z_S | T). \tag{17}$$

Therefore, the learned representation $Z_S$ is expected to contain minimal and sufficient information about downstream tasks [45, 9], which further illuminates the reason why the embeddings given by SSL approaches have superior performance on various downstream tasks.

Table 2: Test accuracy on citation networks. The *input* column highlights the data used for training. ($\mathbf{X}$ for node features, $\mathbf{A}$ for adjacency matrix, $\mathbf{S}$ for diffusion matrix, and $\mathbf{Y}$ for node labels).

|  | Methods | Input | Cora | Citeseer | Pubmed |
|---|---|---|---|---|---|
| Supervised | MLP [47] | $\mathbf{X}, \mathbf{Y}$ | 55.1 | 46.5 | 71.4 |
|  | LP [56] | $\mathbf{A}, \mathbf{Y}$ | 68.0 | 45.3 | 63.0 |
|  | GCN [22] | $\mathbf{X}, \mathbf{A}, \mathbf{Y}$ | 81.5 | 70.3 | 79.0 |
|  | GAT [47] | $\mathbf{X}, \mathbf{A}, \mathbf{Y}$ | $83.0 \pm 0.7$ | $72.5 \pm 0.7$ | $79.0 \pm 0.3$ |
| Unsupervised | Raw Features [48] | $\mathbf{X}$ | $47.9 \pm 0.4$ | $49.3 \pm 0.2$ | $69.1 \pm 0.3$ |
|  | Linear CCA [18] | $\mathbf{X}$ | $58.9 \pm 1.5$ | $27.5 \pm 1.3$ | $75.8 \pm 0.4$ |
|  | DeepWalk [32] | $\mathbf{A}$ | $70.7 \pm 0.6$ | $51.4 \pm 0.5$ | $74.3 \pm 0.9$ |
|  | GAE [21] | $\mathbf{X}, \mathbf{A}$ | $71.5 \pm 0.4$ | $65.8 \pm 0.4$ | $72.1 \pm 0.5$ |
|  | DGI [48] | $\mathbf{X}, \mathbf{A}$ | $82.3 \pm 0.6$ | $71.8 \pm 0.7$ | $76.8 \pm 0.6$ |
|  | MVGRL[1] [15] | $\mathbf{X}, \mathbf{S}, \mathbf{A}$ | $83.5 \pm 0.4$ | $\mathbf{73.3 \pm 0.5}$ | $80.1 \pm 0.7$ |
|  | GRACE[2] [57] | $\mathbf{X}, \mathbf{A}$ | $81.9 \pm 0.4$ | $71.2 \pm 0.5$ | $80.6 \pm 0.4$ |
|  | CCA-SSG (Ours) | $\mathbf{X}, \mathbf{A}$ | $\mathbf{84.2 \pm 0.4}$ | $73.1 \pm 0.3$ | $\mathbf{81.6 \pm 0.4}$ |

[1] Results on Cora with authors' code is inconsistent with [15]. We adopt the results with authors' code.
[2] Results are from our reproducing with authors' code, as [57] did not use the public splits.

## 5 Experiments

We assess the quality of representations after self-supervised pretraining on seven node classification benchmarks: *Cora, Citeseer, Pubmed, Coauthor CS, Coauthor Physics* and *Amazon Computer, Amazon-Photo*. We adopt the public splits for *Cora, Citeseer, Pubmed*, and a 1:1:9 training/validation/testing splits for the other 4 datasets. Details of the datasets are in Appendix E.

**Evaluation protocol**. We follow the linear evaluation scheme as introduced in [48]: **i)** We first train the model on all the nodes in a graph without supervision, by optimizing the objective in Eq. (5). **ii)** After that, we freeze the parameters of the encoder and obtain all the nodes' embeddings, which are subsequently fed into a linear classifier (i.e., a logistic regression model) to generate a predicted label for each node. In the second stage, only nodes in training set are used for training the classifier, and we report the classification accuracy on testing nodes.

We implement the model with PyTorch. All experiments are conducted on a NVIDIA V100 GPU with 16 GB memory. We use the Adam optimizer [20] for both stages. The graph encoder $f_\theta$ is specified as a standard two-layer GCN model [22] for all the datasets except *citeseer* (where we empirically find that a one-layer GCN is better). We report the mean accuracy with a standard deviation through 20 random initialization (on *Coauthor CS, Coauthor Physics* and *Amazon Computer, Amazon-Photo*, the split is also randomly generated). Detailed hyperparameter settings are in Appendix E.

### 5.1 Comparison with Peer Methods

We compare CCA-SSG with classical unsupervised models, Deepwalk [32] and GAE [21], and self-supervised models, DGI [48], MVGRL [15], GRACE [57] and GCA [58]. We also compare with supervised learning models, including MLP, Label Propagation (LP) [56], and supervised baselines GCN [22] and GAT [47][3]. The results of baselines are quoted from [15, 57, 58] if not specified.

We report the node classification results of citation networks and other datasets in Table 2 and Table 3 respectively. As we can see, CCA-SSG outperforms both the unsupervised competitors and the fully supervised baselines on *Cora* and *Pubmed*, despite its simple architecture. On *Citeseer*, CCA-SSG achieves competitive results as of the most powerful baseline MVGRL. On four larger benchmarks, CCA-SSG also achieves the best performance in four datasets except *Coauther-Physics*. It is worth mentioning that we empirically find that on *Coauthor-CS* a pure 2-layer-MLP encoder is better than GNN models. This might because the graph-structured information is much less informative than the node features, presumably providing harmful signals for classification (in fact, on *Coauthor-CS*, linear models using merely node features can greatly outperform DeepWalk/DeepWalk+features).

---

[3]The BGRL [39] is not compared as its source code has not been released.

Table 3: Test accuracy on co-author and co-purchase networks. We report both mean accuracy and standard deviation. Results of baseline models are from [58].

| | Methods | Input | Computer | Photo | CS | Physics |
|---|---|---|---|---|---|---|
| | Supervised GCN [22] | $\mathbf{X}, \mathbf{A}, \mathbf{Y}$ | $86.51 \pm 0.54$ | $92.42 \pm 0.22$ | $93.03 \pm 0.31$ | $95.65 \pm 0.16$ |
| | Supervised GAT [47] | $\mathbf{X}, \mathbf{A}, \mathbf{Y}$ | $86.93 \pm 0.29$ | $92.56 \pm 0.35$ | $92.31 \pm 0.24$ | $95.47 \pm 0.15$ |
| Unsupervised | Raw Features [48] | $\mathbf{X}$ | $73.81 \pm 0.00$ | $78.53 \pm 0.00$ | $90.37 \pm 0.00$ | $93.58 \pm 0.00$ |
| | Linear CCA [18] | $\mathbf{X}$ | $79.84 \pm 0.53$ | $86.92 \pm 0.72$ | $93.13 \pm 0.18$ | $95.04 \pm 0.17$ |
| | DeepWalk [32] | $\mathbf{A}$ | $85.68 \pm 0.06$ | $89.44 \pm 0.11$ | $84.61 \pm 0.22$ | $91.77 \pm 0.15$ |
| | DeepWalk + features | $\mathbf{X}, \mathbf{A}$ | $86.28 \pm 0.07$ | $90.05 \pm 0.08$ | $87.70 \pm 0.04$ | $94.90 \pm 0.09$ |
| | GAE [21] | $\mathbf{X}, \mathbf{A}$ | $85.27 \pm 0.19$ | $91.62 \pm 0.13$ | $90.01 \pm 0.71$ | $94.92 \pm 0.07$ |
| | DGI [48] | $\mathbf{X}, \mathbf{A}$ | $83.95 \pm 0.47$ | $91.61 \pm 0.22$ | $92.15 \pm 0.63$ | $94.51 \pm 0.52$ |
| | MVGRL [15] | $\mathbf{X}, \mathbf{S}, \mathbf{A}$ | $87.52 \pm 0.11$ | $91.74 \pm 0.07$ | $92.11 \pm 0.12$ | $95.33 \pm 0.03$ |
| | GRACE[1] [57] | $\mathbf{X}, \mathbf{A}$ | $86.25 \pm 0.25$ | $92.15 \pm 0.24$ | $92.93 \pm 0.01$ | $95.26 \pm 0.02$ |
| | GCA[1] [58] | $\mathbf{X}, \mathbf{A}$ | $87.85 \pm 0.31$ | $92.49 \pm 0.09$ | $93.10 \pm 0.01$ | $\mathbf{95.68 \pm 0.05}$ |
| | CCA-SSG (Ours) | $\mathbf{X}, \mathbf{A}$ | $\mathbf{88.74 \pm 0.28}$ | $\mathbf{93.14 \pm 0.14}$ | $\mathbf{93.31 \pm 0.22}$ | $95.38 \pm 0.06$ |

[1] GCA is essentially an enhanced version of GRACE by adopting adaptive augmentations. Both GRACE and GCA would suffer from *out of memory* on *Coauthor-Physics* using a GPU wth 16GB memory. The reported results are from authors' papers using a 32GB GPU.

## 5.2 Ablation Study and Scalability Comparison

**Effectiveness of invariance/decorrelation terms**. We alter our loss by removing the invariance/decorrelation term respectively to study the effects of each component, with results reported in Table 4. We find that only using the invariance term will lead to merely performance drop instead of completely collapsed solutions. This is because node embeddings are normalized along the instance dimension to have a zero-mean and fixed-standard deviation, and the worst solution is no worse than dimensional collapse (i.e., all the embeddings lie in an line, and our decorrelation term can help to prevent it) instead of complete collapse (i.e., all the embeddings degenerate into a single point). As expected, only optimizing the decorrelation term will lead to poor result, as the model learns nothing meaningful but disentangled representation. In Appendix B we discuss the relationship between complete/dimensional collapse, when the two cases happen and how to avoid them.

**Effect of decorrelation intensity**. We study how the intensity of feature decorrelation improves/degrades the performance by increasing the trade-off hyper-parameter $\lambda$. Fig. 2 shows test accuracy w.r.t. different $\lambda$'s on *Cora, Citeseer* and *Pubmed*. The performance benefits from a proper selection of $\lambda$ (from 0.0005 to 0.001 in our experiments). When $\lambda$ is too small, the decorrelation term does not work; if it is too large, the invariance term would be neglected, leading to serious performance degrade. An interesting finding is that even when $\lambda$ is very small or even equals to 0 (w/o $\mathcal{L}_{dec}$ in Table 4), the test accuracy on *Citeseer* does not degrade as much as that on *Cora* and *Citeseer*. The reason is that node embeddings of *Citeseer* is already highly uncorrelated even without the decorrelation term. Appendix F visualizes the correlation matrices without/with decorrelations.

**Effect of embedding dimension**. Fig. 3 shows the effect of the embedding dimension. Similar to contrastive methods [48, 15, 57, 58], CCA-SSG benefits from a large embedding dimension (compared with supervised learning), while the optimal embedding dimension of CCA-SSG (512 on most benchmarks) is a bit larger than other methods (usually 128 or 256). Yet, we notice a performance drop as the embedding dimension increases. We conjecture that the CCA is essentially a dimension-reduction method, the ideal embedding dimension ought to be smaller than the dimension of input. Hence we do not apply it on well-compressed datasets (e.g. ogbn-arXiv and ogbn-product).

**Scalability Comparison.** Table 5 compares model size, training time (till the epoch that gives the highest evaluation accuracy) and memory cost of CCA-SSG with other methods, on *Cora, Pubmed* and *Amazon-Computers*. Overall, our method has fewer parameters, shorter training time, and fewer memory cost than MVGRL, GRACE and GCA in most cases. DGI is another simple and efficient model, but it yields much poorer performance. The results show that despite its simplicity and efficiency, our method achieves even better (or competitive) performance.

Table 4: Ablation study of node classification accuracy (%) on the key components of CCA-SSG.

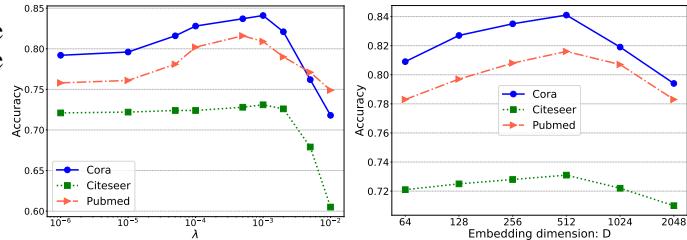| Variants | Cora | Citeseer | Pubmed |
|---|---|---|---|
| Baseline | 84.2 | 73.1 | 81.6 |
| w/o $\mathcal{L}_{dec}$ | 79.1 | 72.2 | 75.3 |
| w/o $\mathcal{L}_{inv}$ | 40.1 | 28.9 | 46.5 |



Figure 2: Effect of $\lambda$.


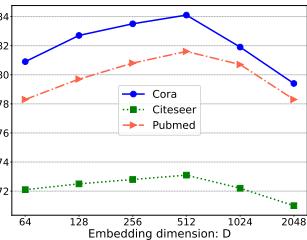
Figure 3: Effect of $D$.

Table 5: Comparison of the number of parameters, training time for achieving the best performance, and the memory cost of different methods on *Cora, Pubmed* and *Amazon-Computer*. MVGRL on Pubmed and Computer requires subgraph sampling with graph size $4000$. Others are full-graph.

| Methods | Cora ($N$: 2,708) | | | Pubmed ($N$: 19,717) | | | Computer ($N$: 13,752) | | |
|---|---|---|---|---|---|---|---|---|---|
| | #Paras | Time | Mem | #Paras | Time | Mem | #Paras | Time | Mem |
| DGI | 1260K | 6.4s | 1.4G | 782K | 5.9s | 1.9G | 919K | 14.1s | 1.9G |
| MVGRL | 1731K | 26.9s | 4.6G | 775K | 29s | 5.4G | 1049K | 31.5s | 5.5G |
| GRACE/GCA | 997K | 8.3s | 1.7G | 520K | 756s | 12.6G | 273K | 314s | 7.6G |
| CCA-SSG(Ours) | 997K | 3.8s | 1.6G | 519K | 9.6s | 2.7G | 656K | 14.8s | 2.5G |

# 6 Conclusion and Discussions

In this paper, we have introduced CCA-SSG, a conceptually simple, efficient yet effective method for self-supervised representation learning on graphs, based on the idea of Canonical Correlation Analysis. Compared with contrastive methods, our model does not require additional components except random augmentations and a GNN encoder, whose effectiveness is justified in experiments.

**Limitations of the work.** Despite the theoretical grounds and the promising experimental justifications, our method would suffer from several limitations. 1) The objective Eq. (5) is essentially performing dimension reduction, while SSL approach usually requires a large embedding dimension. As a result, our method might not work well on datasets where input data does not have a large feature dimension. 2) Like other augmentation based methods, CCA-SSG highly relies on a high-quality, informative and especially, label-invariant augmentations. However, the augmentations used in our model might not perfectly meet these requirements, and it remains an open problem how to generate informative graph augmentations that have non-negative impacts on the downstream tasks.

**Potential negative societal impacts.** This work explores a simple pipeline for representation learning without large amount of labeled data. However, in industry there are many career workers whose responsibility is to label or annotate data. The proposed method might reduce the need for labeling data manually, and thus makes a few individuals unemployed (especially for developing countries and remote areas). Furthermore, our model might be biased, as it tends to pay more attention to the majority and dominant features (shared information across most of the data). The minority group whose features are scare are likely to be downplayed by the algorithm.

# Acknowledgments and Disclosure of Funding

# References

[1] Rana Ali Amjad and Bernhard C. Geiger. Learning representations for neural network-based classification using the information bottleneck principle. *IEEE Trans. Pattern Anal. Mach.*

*Intell.*, 42(9):2225–2239, 2020.

[2] Galen Andrew, Raman Arora, Jeff A. Bilmes, and Karen Livescu. Deep canonical correlation analysis. In *ICML*, volume 28, pages 1247–1255, 2013.

[3] Mohamed Ishmael Belghazi, Aristide Baratin, Sai Rajeshwar, Sherjil Ozair, Yoshua Bengio, Aaron Courville, and Devon Hjelm. Mutual information neural estimation. In *International Conference on Machine Learning*, pages 531–540, 2018.

[4] Xiaobin Chang, Tao Xiang, and Timothy M. Hospedales. Scalable and effective deep CCA via soft decorrelation. In *CVPR*, pages 1488–1497, 2018.

[5] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey E. Hinton. A simple framework for contrastive learning of visual representations. In *ICML*, Proceedings of Machine Learning Research, pages 1597–1607, 2020.

[6] Xinlei Chen and Kaiming He. Exploring simple siamese representation learning. *arXiv preprint arXiv:2011.10566*, 2020.

[7] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: pre-training of deep bidirectional transformers for language understanding. In *NAACL*, pages 4171–4186, 2019.

[8] Aleksandr Ermolov, Aliaksandr Siarohin, Enver Sangineto, and Nicu Sebe. Whitening for self-supervised representation learning. *arXiv preprint arXiv:2007.06346*, 2020.

[9] Marco Federici, Anjan Dutta, Patrick Forré, Nate Kushman, and Zeynep Akata. Learning robust representations via multi-view information bottleneck. In *ICLR*, 2020.

[10] Gene H Golub and Hongyuan Zha. The canonical correlations of matrix pairs and their numerical computation. *Linear algebra for signal processing*, pages 27–29, 1995.

[11] Yunchao Gong, Qifa Ke, Michael Isard, and Svetlana Lazebnik. A multi-view embedding space for modeling internet images, tags, and their semantics. *Int. J. Comput. Vis.*, 106(2):210–233, 2014.

[12] Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre H. Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Ávila Pires, Zhaohan Guo, Mohammad Gheshlaghi Azar, Bilal Piot, Koray Kavukcuoglu, Rémi Munos, and Michal Valko. Bootstrap your own latent - A new approach to self-supervised learning. In *NeurIPS*, 2020.

[13] Michael Gutmann and Aapo Hyvärinen. Noise-contrastive estimation: A new estimation principle for unnormalized statistical models. In *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*, pages 297–304, 2010.

[14] David R. Hardoon, Sándor Szedmák, and John Shawe-Taylor. Canonical correlation analysis: An overview with application to learning methods. *Neural Comput.*, 16(12):2639–2664, 2004.

[15] Kaveh Hassani and Amir Hosein Khas Ahmadi. Contrastive multi-view representation learning on graphs. In *ICML*, volume 119 of *Proceedings of Machine Learning Research*, pages 4116–4126. PMLR, 2020.

[16] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross B. Girshick. Momentum contrast for unsupervised visual representation learning. In *CVPR*, pages 9726–9735, 2020.

[17] R. Devon Hjelm, Alex Fedorov, Samuel Lavoie-Marchildon, Karan Grewal, Philip Bachman, Adam Trischler, and Yoshua Bengio. Learning deep representations by mutual information estimation and maximization. In *ICLR*, 2019.

[18] Harold Hotelling. Relations between two sets of variates. *Biometrika*, 28:322–377, 1936.

[19] Tianyu Hua, Wenxiao Wang, Zihui Xue, Yue Wang, Sucheng Ren, and Hang Zhao. On feature decorrelation in self-supervised learning. *arXiv preprint arXiv:2105.00470*, 2021.

[20] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015.

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

[22] Thomas N. Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. In *ICLR*, 2017.

[23] Johannes Klicpera, Stefan Weißenberger, and Stephan Günnemann. Diffusion improves graph learning. In *NeurIPS*, pages 13333–13345, 2019.

[24] Qimai Li, Zhichao Han, and Xiao-Ming Wu. Deeper insights into graph convolutional networks for semi-supervised learning. In *AAAI*, volume 32, 2018.

[25] Yunfan Li, Peng Hu, Jerry Zitao Liu, Dezhong Peng, Joey Tianyi Zhou, and Xi Peng. Contrastive clustering. *arXiv preprint arXiv:2009.09687*, 2020.

[26] Julian J. McAuley, Christopher Targett, Qinfeng Shi, and Anton van den Hengel. Image-based recommendations on styles and substitutes. In *SIGIR*, pages 43–52, 2015.

[27] Jin Ming, Zheng Yizhen, Li Yuan-Fang, Gong Chen, Zhou Chuan, and Pan Shirui. Multi-scalecontrastive siamese networks for self-supervised graph representation learning. In *IJCAI*, 2021.

[28] Xu Minghao, Wang Hang, Ni Bingbing, Guo Hongyu, and Tang Jian. Self-supervised graph-levelrepresentation learning with local and global structure. In *ICML*, 2021.

[29] Galileo Namata, Ben London, Lise Getoor, Bert Huang, and UMD EDU. Query-driven active surveying for collective classification. In *10th International Workshop on Mining and Learning with Graphs*, volume 8, 2012.

[30] Sebastian Nowozin, Botond Cseke, and Ryota Tomioka. f-gan: Training generative neural samplers using variational divergence minimization. In *NIPS*, pages 271–279, 2016.

[31] Zhen Peng, Wenbing Huang, Minnan Luo, Qinghua Zheng, Yu Rong, Tingyang Xu, and Junzhou Huang. Graph representation learning via graphical mutual information maximization. In *WWW*, pages 259–270, 2020.

[32] Bryan Perozzi, Rami Al-Rfou, and Steven Skiena. Deepwalk: online learning of social representations. In *KDD*, pages 701–710. ACM, 2014.

[33] Jiezhong Qiu, Qibin Chen, Yuxiao Dong, Jing Zhang, Hongxia Yang, Ming Ding, Kuansan Wang, and Jie Tang. GCC: graph contrastive coding for graph neural network pre-training. In *KDD*, pages 1150–1160, 2020.

[34] Prithviraj Sen, Galileo Namata, Mustafa Bilgic, Lise Getoor, Brian Galligher, and Tina Eliassi-Rad. Collective classification in network data. *AI magazine*, 29(3):93–93, 2008.

[35] Arnab Sinha, Zhihong Shen, Yang Song, Hao Ma, Darrin Eide, Bo-June Paul Hsu, and Kuansan Wang. An overview of microsoft academic service (MAS) and applications. In *WWW*, pages 243–246, 2015.

[36] Karthik Sridharan and Sham M. Kakade. An information theoretic framework for multi-view learning. In *COLR*, pages 403–414. Omnipress, 2008.

[37] DJ Strouse and David J. Schwab. The deterministic information bottleneck. In *UAI*, 2016.

[38] Fan-Yun Sun, Jordan Hoffmann, Vikas Verma, and Jian Tang. Infograph: Unsupervised and semi-supervised graph-level representation learning via mutual information maximization. In *ICLR*, 2020.

[39] Shantanu Thakoor, Corentin Tallec, Mohammad Gheshlaghi Azar, Rémi Munos, Petar Velickovic, and Michal Valko. Bootstrapped representation learning on graphs. *arXiv preprint arXiv:2102.06514*, 2021.

[40] Yonglong Tian, Dilip Krishnan, and Phillip Isola. Contrastive multiview coding. In *ECCV*, pages 776–794, 2020.

[41] Yuandong Tian, Xinlei Chen, and Surya Ganguli. Understanding self-supervised learning dynamics without contrastive pairs. *arXiv preprint arXiv:2102.06810*, 2021.

[42] Yuandong Tian, Lantao Yu, Xinlei Chen, and Surya Ganguli. Understanding self-supervised learning with dual deep networks. *arXiv preprint arXiv:2010.00578*, 2020.

[43] Naftali Tishby, Fernando C. N. Pereira, and William Bialek. The information bottleneck method. *arXiv preprint physics/0004057*, 2000.

[44] Naftali Tishby and Noga Zaslavsky. Deep learning and the information bottleneck principle. In *ITW*, pages 1–5, 2015.

[45] Yao-Hung Hubert Tsai, Yue Wu, Ruslan Salakhutdinov, and Louis-Philippe Morency. Self-supervised learning from a multi-view perspective. In *ICLR*, 2021.

[46] Aäron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018.

[47] Petar Velickovic, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio. Graph attention networks. In *ICLR*, 2018.

[48] Petar Velickovic, William Fedus, William L. Hamilton, Pietro Liò, Yoshua Bengio, and R. Devon Hjelm. Deep graph infomax. In *ICLR*, 2019.

[49] Minjie Wang, Lingfan Yu, Da Zheng, Quan Gan, Yu Gai, Zihao Ye, Mufei Li, Jinjing Zhou, Qi Huang, Chao Ma, Ziyue Huang, Qipeng Guo, Hao Zhang, Haibin Lin, Junbo Zhao, Jinyang Li, Alexander J. Smola, and Zheng Zhang. Deep graph library: Towards efficient and scalable deep learning on graphs. *arXiv*, 1909.01315, 2019.

[50] Tongzhou Wang and Phillip Isola. Understanding contrastive representation learning through alignment and uniformity on the hypersphere. In *ICML*, pages 9929–9939, 2020.

[51] Lee Xiong, Chenyan Xiong, Ye Li, Kwok-Fung Tang, Jialin Liu, Paul Bennett, Junaid Ahmed, and Arnold Overwijk. Approximate nearest neighbor negative contrastive learning for dense text retrieval. In *ICLR*, 2021.

[52] Yuning You, Tianlong Chen, Yongduo Sui, Ting Chen, Zhangyang Wang, and Yang Shen. Graph contrastive learning with augmentations. In *NeurIPS*, 2020.

[53] Jure Zbontar, Li Jing, Ishan Misra, Yann LeCun, and Stéphane Deny. Barlow twins: Self-supervised learning via redundancy reduction. In *ICML*, 2021.

[54] Jiaqi Zeng and Pengtao Xie. Contrastive self-supervised learning for graph classification. *arXiv preprint arXiv:2009.05923*, 2020.

[55] Hanlin Zhang, Shuai Lin, Weiyang Liu, Pan Zhou, Jian Tang, Xiaodan Liang, and Eric P Xing. Iterative graph self-distillation. *arXiv preprint arXiv:2010.12609*, 2020.

[56] Xiaojin Zhu, Zoubin Ghahramani, and John D. Lafferty. Semi-supervised learning using gaussian fields and harmonic functions. In *ICML*, pages 912–919, 2003.

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

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

# A Algorithm

We provide the pseudo code for our method in Algorithm 2, the detailed description of which is in Section 3.1.

---

**Algorithm 2:** Algorithm for CCA-SSG

---

**Input:** A graph $\mathbf{G} = (\mathbf{X}, \mathbf{A})$ with $N$ nodes, where $\mathbf{X}$ is node feature matrix, and $\mathbf{A}$ is the adjacency matrix. Augmentations $\mathcal{T}$, encoder $f_\theta$ by random initialization, trade-off hyper-parameter $\lambda$, maximum training steps $T$.
**Output:** Learned encoder $f_\theta$.

1 **while** *not reaching $T$* **do**
2      Sample two augmentation functions $t_A, t_B \sim \mathcal{T}$;
3      Generate transformed graphs: $\tilde{\mathbf{G}}_A = (\tilde{\mathbf{X}}_A, \tilde{\mathbf{A}}_A), \tilde{\mathbf{G}}_B = (\tilde{\mathbf{X}}_B, \tilde{\mathbf{A}}_B)$ ;
4      Get node embeddings through the graph neural network as encoder:
        $\mathbf{Z}_A = f_\theta(\tilde{\mathbf{X}}_A, \tilde{\mathbf{A}}_A), \mathbf{Z}_B = f_\theta(\tilde{\mathbf{X}}_B, \tilde{\mathbf{A}}_B)$ ;
5      Normalize embeddings along instance dimension: $\tilde{\mathbf{Z}}_A = \frac{\mathbf{Z}_A - \mu(\mathbf{Z}_A)}{\sigma(\mathbf{Z}_A) * \sqrt{N}}, \tilde{\mathbf{Z}}_B = \frac{\mathbf{Z}_B - \mu(\mathbf{Z}_B)}{\sigma(\mathbf{Z}_B) * \sqrt{N}}$ ;
6      Calculate the loss function $\mathcal{L}$ according to Eq. (5) ;
7      Update $\theta$ by gradient descent;
8 **Inference**: $\mathbf{Z} = f_\theta(\mathbf{X}, \mathbf{A})$, where $\theta$ is the frozen parameters of the encoder.

---

# B Discussions on Degenerated Solutions in SSL

In this section we provide an illustration and some discussions for degenerated (collapsed) solutions, or namely trivial solutions, in self-supervised representation learning. The discussion is inspired by the separation of complete collapse and dimensional collapse proposed in [19]. We show that our method naturally avoids complete collapse through feature-wise normalization, and could prevent/alleviate dimensional collapse through the decorrelation term Eq. (7).

In most contrastive learning methods especially the augmentation-based ones [46, 16, 5, 40], both positive pairs and negative pairs are required for learning a model. For instance, the widely adopted InfoNCE [46] loss has the following formulation:

$$\mathcal{L}_{\text{InfoNCE}} = -\log \frac{\exp\left(z_i^A \cdot z_i^B / \tau\right)}{\sum_j \exp\left(z_i^A \cdot z_j^B / \tau\right)}, \tag{18}$$

where $z_i^A$ and $z_i^B$ are the (normalized) embeddings of two views of the same instance $i$, and $\tau$ is the temperature hyperparameter. The numerator enforces similarity between positive pairs (two views of the same instance), while the denominator promotes dis-similarity between negative pairs (two views of different instances). Therefore, minimizing Eq. (18) is equivalent to maximizing the cosine similarity of positive pairs and meanwhile minimizing the cosine similarity of negative pairs. Note that the normalization is applied for each instance (projecting the embedding onto a hypersphere), so we are essentially minimizing the distances between positive pairs and maximizing the distance between negative pairs. The previous work [50] provides a thorough analysis on the behaviors of the objective by decomposing it into two terms: 1) alignment term (for positive pairs) and 2) uniformity term (for negative pairs).

The alignment loss is defined as the expected distance between positive pairs:

$$\mathcal{L}_{\text{align}} \triangleq \mathbb{E}_{(x,y) \sim p_{\text{pos}}} \|f(x) - f(y)\|_2^\alpha, \text{ with } \alpha > 0. \tag{19}$$

The uniformity loss is the logarithm of the average pairwise Gaussian potential:

$$\mathcal{L}_{\text{uniformity}} \triangleq \log \mathbb{E}_{x,y \sim p_{\text{data}}} e^{-t\|f(x) - f(y)\|_2^2}, \text{ with } t > 0 \tag{20}$$

Intuitively, the alignment term makes the positive pairs close to each other on the hypersphere, while the uniformity term makes different data points distribute on the hypersphere uniformly.

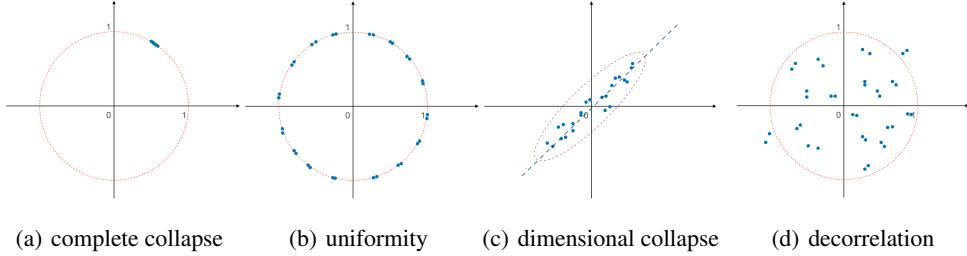|  |  |  |  |
|---|---|---|---|
| (a) complete collapse | (b) uniformity | (c) dimensional collapse | (d) decorrelation |

Figure 4: An illustration of the two types of collapse and how to deal with them with a 2-d case. Blue circles are data points. Fig. 4(a): complete collapse, when all data samples degenerate to a same point on the hypersphere. Fig. 4(b): the uniformity loss keeps positive pairs close, but forces all data points to distribute on the hypersphere uniformly. Fig. 4(c): in dimensional collapse, the data points are not projected onto the hypersphere, but they distribute nearly as a line in the space, making them hard to discriminate. Fig. 4(d): the decorrelation term prevents dimensional collapse by directly decorrelating each dimensional representations, which implicitly scatters the data points.

In particular, only considering the alignment term in Eq. (19) will lead to trivial solutions: all the embeddings would degenerate to a fixed point on the hypersphere. This phenomenon is called **complete collapse** [19]. Denote $\mathbf{Z}_A$ and $\mathbf{Z}_B$ as two embedding matrix of two views ($\mathbf{Z} \in \mathbb{R}^{N \times D}$ and is row normalized), then in this case $\mathbf{Z}_A \mathbf{Z}_B^\top \cong \mathbf{1}$ is an all-one matrix (so as $\mathbf{Z}_A \mathbf{Z}_A^\top$ and $\mathbf{Z}_B \mathbf{Z}_B^\top$).

The uniformity term in Eq. (20) prevents complete collapse by separating the embeddings of arbitrary two data points, so that the data points would be embedded uniformly on the hypersphere. Fig. 4(a) and 4(b) provide an illustration for complete collapse and how the uniformity term prevents it.

Another kind of collapse that has been neglected by most existing works is **dimensional collapse** [19]. Different from complete collapse where all the data points degenerate into **a single point**, dimensional collapse means data points are distributed on **a line**, and each dimension captures exactly the same features (or different dimensions are highly correlated can capture the same information). Note that if the data representations are normalized along feature dimensions, all the data points would be projected onto a hypersphere. Under this circumstance there will not be dimensional collapse. However, if we normalize the output along the instance dimension so that each column has zero-mean and $1/\sqrt{N}$-standard deviation, as is done in this paper in Eq. (6), merely optimizing Eq. (6) would not prevent dimensional collapse, i.e. $\tilde{\mathbf{Z}}^\top \tilde{\mathbf{Z}} \cong \mathbf{1}$ ($\tilde{\mathbf{Z}} \in \mathbb{R}^{N \times D}$ is normalized by column).

In our model, the feature decorrelation term in Eq. (7) exactly prevents dimensional collapse by minimizing $\left\| \tilde{\mathbf{Z}}^\top \tilde{\mathbf{Z}} - \mathbf{I} \right\|_F^2$. Note that the diagonal term is always equal to 1, so we are pushing each dimension to capture orthogonal features. Also, the feature decorrelation term implicitly scatters the data points in the space, making them distinguishable for downstream tasks [19, 8]. An illustration of the dimensional collapse and the effect of feature decorrelation is provided in Fig. 4(c) and 4(d), respectively.

## C  Properties of Mutual Information and Entropy

In this section, we enumerate some useful properties of mutual information and entropy that will be used in Appendix D for proving the theorems. For any random variables $A, B, C, X, Z$, we have:

- **Property 1**. Non-negativity:

$$I(A, B) \geq 0, I(A, B|C) \geq 0. \tag{21}$$

- **Property 2**. Chain rule:

$$I(A, B, C) = I(A, B) - I(A, B|C). \tag{22}$$

- **Property 3**. Data Processing Inequality (DPI). $Z = f_\theta(X)$, then:

$$I(Z, A) = I(f_\theta(X), A) \leq I(X, A) \tag{23}$$

- **Property 4**. Non-negativity of discrete entropy. For discrete random variable:
$$H(A) \geq 0, H(A|B) \geq 0. \tag{24}$$

- **Property 5**. Relationship between entropy and mutual information:
$$H(A) = H(A|B) + I(A, B). \tag{25}$$

- **Property 6**. Entropy of deterministic function. If $Z$ is deterministic given $X$:
$$H(Z|X) = 0 \tag{26}$$

- **Property 7**. Entropy of Gaussian distribution. Assume $X$ obeys a $k$-dimensional Gaussian distribution, $X \sim \mathcal{N}(\mu, \Sigma)$, and we have
$$H(X) = \frac{k}{2}(\ln 2\pi + 1) + \frac{1}{2}\ln|\Sigma|. \tag{27}$$

# D   Proofs in Section 4

As already introduced in Section 4, we use $X$ and $S$ to denote the data and its augmentations respectively. We use $Z_X$ and $Z_S$ to denote their embeddings through the encoder $f_\theta$: $Z_X = f_\theta(X)$, $Z_S = f_\theta(S)$. We aim to learn the optimal encoder parameters $\theta$.

## D.1   Proof of Proposition 1

Restate Proposition 1:

**Proposition 1.** *In expectation, minimizing Eq.* (6) *is equivalent to minimizing the entropy of $Z_S$ conditioned on the input $X$, i.e.,:*
$$\min_\theta \mathcal{L}_{inv} \cong \min_\theta H(Z_S|X) \tag{28}$$

*Proof.* Assume input data come from a distribution $\boldsymbol{x} \sim p(\boldsymbol{x})$ and $\boldsymbol{s}$ is a view of $\boldsymbol{x}$ through random augmentation $\boldsymbol{s} \sim p_{aug}(\cdot|\boldsymbol{x})$. Denote $\boldsymbol{z_s}$ as the representation of $\boldsymbol{s}$. Note that $\boldsymbol{s}_1$ and and $\boldsymbol{s}_2$ both come from $p_{aug}(\cdot|\boldsymbol{x})$.

Recall the invariance term: $\mathcal{L}_{inv} = \left\|\tilde{\mathbf{Z}}_A - \tilde{\mathbf{Z}}_B\right\|_F^2 = \sum_{i=1}^{N}\sum_{k=1}^{D}\left(\tilde{z}_{i,k}^A - \tilde{z}_{i,k}^B\right)^2$. If we ignore the normalization and use $\boldsymbol{s}_1$ and $\boldsymbol{s}_2$ to represent view $A$ and view $B$. We have:

$$
\begin{aligned}
\mathcal{L}_{inv}/N &\cong \mathbb{E}_{\boldsymbol{x}}\left(\sum_{k=1}^{D}\mathbb{E}_{\boldsymbol{s}_1,\boldsymbol{s}_2 \sim p(\cdot|\boldsymbol{x})}(z_k^{\boldsymbol{s}_1} - z_k^{\boldsymbol{s}_2})^2\right) \\
&= \mathbb{E}_{\boldsymbol{x}}\left(\sum_{k=1}^{D}\mathbb{E}_{\boldsymbol{s}_1,\boldsymbol{s}_2 \sim p(\cdot|\boldsymbol{x})}(z_k^{\boldsymbol{s}_1\,2} + z_k^{\boldsymbol{s}_2\,2} - 2*z_k^{\boldsymbol{s}_1}z_k^{\boldsymbol{s}_2})\right) \\
&= 2*\mathbb{E}_{\boldsymbol{x}}\left(\sum_{k=1}^{D}\mathbb{V}_{\boldsymbol{s} \sim p(\cdot|\boldsymbol{x})}z_k^{\boldsymbol{s}}\right) \\
&= 2*\sum_{k=1}^{D}\mathbb{E}_{\boldsymbol{x}}\left(\mathbb{V}_{\boldsymbol{s} \sim p(\cdot|\boldsymbol{x})}z_k^{\boldsymbol{s}}\right)
\end{aligned}
\tag{29}
$$

This indicates that minimizing $\mathcal{L}_{inv}$ is to minimize the variance of augmentation's representations conditioned on the input data.

Note the decorrelation term Eq. (7) aims to learn orthogonal representations at each dimension. If the representations are perfectly decorrelated, then $H(Z_S|X) = \sum_k H(Z_{S,k}|X)$. With Assumption 1, each dimensional representation also obeys a 1-dimensional Gaussian distribution, whose entropy is $H(Z_{S,k}|X) = \frac{1}{2}\log 2\pi e \sigma_k^{s\,2}$. This indicates by minimizing the variance of features at each dimension, its entropy is also minimized. Hence we have Proposition 1. $\square$

**Remark 1.** $I(Z_S, S|X) = H(Z_S|X) - H(Z_S|S, X) = H(Z_S|X)$ *(Property 6 in Appendix C). So $I(Z_S, S|X)$ is also minimized.*

### D.2 Proof of Proposition 2

Restate Proposition 2:

**Proposition 2.** *Minimizing Eq. (7) is equivalent to maximizing the entropy of $Z_S$, i.e.,*

$$\min_{\theta} \mathcal{L}_{dec} \cong \max_{\theta} H(Z_S). \tag{30}$$

*Proof.* With the assumption that $Z_S$ obeys a Gaussian distribution, we have:

$$\max_{\theta} H(Z_S) \cong \max_{\theta} \log |\Sigma_{Z_S}|, \tag{31}$$

where $|\Sigma_{Z_S}|$ is the determinant of the covariance matrix of the embeddings of the augmented data. Note that in our implementation we normalize the embedding matrix along the instance dimension: $\Sigma_{Z_S} \cong \tilde{\mathbf{Z}}_S^\top \tilde{\mathbf{Z}}_S$, so the diagonal entries of $\Sigma_{Z_S}$ are all 1's. And $\Sigma_{Z_S} \in \mathbb{R}^{D \times D}$ is a symmetric matrix.

If $\lambda_1, \lambda_2, \cdots, \lambda_D$ are the $D$ eigenvalues of $\Sigma_{Z_S}$, then $\sum_{i=1}^{D} \lambda_i = \text{trace}(\Sigma_{Z_S}) = D$. We have

$$\log |\Sigma_{Z_{\theta,X'}}| = \log \prod_{i=1}^{D} \lambda_i = \underbrace{\sum_{i=1}^{D} \log \lambda_i \leq D \log \frac{\sum_{i=1}^{D} \lambda_i}{D}}_{\text{Jensen Inequality}} = 0. \tag{32}$$

This means that the upper bound of $|\Sigma_{Z_S}|$ is 1, and the upper bound is achieved if and only if $\lambda_i = 1$ for $\forall i$, which indicates $\Sigma_{Z_S}$ is an identity matrix. This global optimum is exactly the same as that of the feature decorrelation term $\mathcal{L}_{dec}$ in Eq. (7). Therefore we conclude the proof. $\qquad\square$

### D.3 Proof of Theorem 1

Restate Theorem 1:

**Theorem 1.** *By optimizing Eqn (5), we maximize the mutual information between the view's embedding $Z_S$ and the input data $X$, and minimize the mutual information between the view's embedding $Z_S$ and the view it self $S$, conditioned on the input data $X$. Formally,*

$$\min_{\theta} \mathcal{L} \Rightarrow \max_{\theta} I(Z_S, X) \quad and \quad \min_{\theta} I(Z_S, S|X). \tag{33}$$

*Proof.* According to Remark 1, we have:

$$I(Z_S, S|X) = H(Z_S|X). \tag{34}$$

According to Property 5 in Appendix C, we have:

$$I(Z_S, X) = H(Z_S) - H(Z_S|X). \tag{35}$$

Then combining Proposition 1 and Proposition 2, we conclude the proof. $\qquad\square$

### D.4 Proof of Theorem 2

Restate Theorem 2:

**Theorem 2.** *Assume $0 < \beta \leq 1$, then by minimizing Eq. (5), the self-supervised Information Bottleneck objective is maximized, formally:*

$$\min_{\theta} \mathcal{L} \Rightarrow \max_{\theta} \mathcal{IB}_{ssl}. \tag{36}$$

*Proof.* According to Property 5 in Appendix C, we can rewrite the IB principle in SSL setting as:
$$\mathcal{IB}_{ssl} = [H(Z_S) - H(Z_S|X)] - \beta \left[ H(Z_S) - H(Z_S|S) \right]. \tag{37}$$

Notice that $Z_S$ is deterministic given $S$: $Z_S = f_\theta(S)$. According to Property 5 in Appendix C, we have $H(Z_S|S) = 0$. Hence, we further have the following relationship
$$\mathcal{IB}_{ssl} = (1 - \beta)H(Z_S) - H(Z_S|X). \tag{38}$$

Let $\lambda = 1 - \beta \geq 0$. Now we can decompose the objective $\mathcal{IB}_{ssl}$ into two terms: 1) maximizing $H(Z_S)$, which increases the information entropy of the embeddings of augmented data. 2) minimizing $H(Z_S|X)$, which decreases the entropy of the embeddings of augmented data, conditioned on the original data.

With Proposition 1 and Proposition 2, we complete the proof. $\qquad\square$

### D.5 Proof of Corollary 1

Restate Corollary 1:

**Corollary 1.** *Let $X_1 = S$, $X_2 = X$ and assume $0 < \beta \leq 1$, then minimizing Eq. (5) is equivalent to minimizing the Multi-view Information Bottleneck Loss in [9]:*
$$\mathcal{L}_{MIB} = I(Z_1, X_1|X_2) - \beta I(X_2, Z_1), 0 < \beta \leq 1 \tag{39}$$

By maximizing $I(X_2, Z_1)$, the model could obtain sufficient information for downstream tasks by ensuring the representation $Z_1$ of $X_1$ is sufficient for $X_2$, and decreasing $I(Z_1, X_1|X_2)$ will increase the robustness of the representation by discarding irrelevant information.

*Proof.* Let $X_1 = S, X_2 = X$ be two views of the input data. We have:
$$\begin{aligned} \mathcal{L}_{MIB} &= I(Z_S, S|X) - \beta I(X, Z_S) \\ &= [H(Z_S|X) - H(Z_S|S, X)] - \beta[H(Z_S) - H(Z_S|X)] \\ &= (1 - \beta)H(Z_S|X) - \beta H(Z_S) - H(Z_S|S, X). \end{aligned} \tag{40}$$
As $Z_S$ is deterministic given $S$, we can obtain $H(Z_S|S, X) = 0$. Based on this, we can further simplify $\mathcal{L}_{MIB}$ as
$$\mathcal{L}_{MIB} = H(Z_S|X) - \lambda H(Z_S), \text{ with } \lambda > 0. \tag{41}$$
With Proposition 1 and Proposition 2, we complete the proof. $\qquad\square$

### D.6 Proof of Corollary 2

Restate Corollary 2:

**Corollary 2.** *When the data augmentation process is reversible, minimizing Eq. (5) is equivalent to learning the Minimal and Sufficient Representations for Self-supervision in [45]:*
$$Z_X^{ssl} = \underset{Z_X}{\arg\max}\, I(Z_X, S), Z_X^{ssl_{min}} = \underset{Z_X}{\arg\min}\, H(Z_X|S) \text{ s.t. } I(Z_X, S) \text{ is maximized.} \tag{42}$$

$Z_X^{ssl}$ is the sufficient self-supervised representation by maximizing $I(Z_X, S)$, and $Z_X^{ssl_{min}}$ is the minimal and sufficient representation by minimizing $H(Z_X|S)$.

*Proof.* Eq. (42) can be converted to minimizing the relaxed Lagrangian objective as below
$$\mathcal{L}_{ssl_{min}} = H(Z_X|S) - \beta I(Z_X, S), \text{ with } \beta > 0. \tag{43}$$
Then $\mathcal{L}_{ssl_{min}}$ could be decomposed into
$$\begin{aligned} \mathcal{L}_{ssl_{min}} &= H(Z_X|S) - \beta I(S, Z_X) \\ &= H(Z_X|S) - \beta[H(Z_X) - H(Z_X|S)] \\ &= (1 + \beta)H(Z_X|S) - \beta H(Z_X) \end{aligned} \tag{44}$$

With $\beta > 0$, $\mathcal{L}_{ssl_{min}}$ is essentially a symmetric formulation of Eq. (38), by exchanging $X$ with $S$, and $Z_X$ with $Z_S$. With the assumption that the data augmentation process is reversible and according to Proposition 1 and Proposition 2, we conclude the proof. $\qquad\square$

### D.7 Proof of Theorem 3

Restate Theorem 3:

**Theorem 3** *(task-relevant/irrelevant information). By optimizing Eq. (5), the task-relevant information $I(Z_S, T)$ is maximized, and the task-irrelevant information $H(Z_S|T)$ is minimized. Formally:*

$$\min_{\theta} \mathcal{L} \Rightarrow \max_{\theta} I(Z_S, T) \text{ and } \min_{\theta} H(Z_S|T). \tag{45}$$

*Proof.* Note that with Assumption 2, we have $I(X, T|S) = I(S, T|X) = 0$. Therefore we obtain $0 \leq I(Z_S, T|X) \leq I(S, T|X) = 0$, which induces $I(Z_S, T|X) = 0$. Then we can derive

$$
\begin{aligned}
I(Z_S, T) =& I(Z_S, T|X) + I(Z_S, X, T) \\
=& 0 + I(Z_S, X) - I(Z_S, X|T) \\
=& I(Z_S, X) - I(Z_S, X|T) \\
\geq& I(Z_S, X) - I(X, S|T)
\end{aligned}
\tag{46}
$$

and

$$
\begin{aligned}
H(Z_S|T) =& H(Z_S|X, T) + I(Z_S, X|T) \\
=& H(Z_S|X) - I(Z_S, T|X) + I(Z_S, X|T) \\
=& H(Z_S|X) - 0 + I(Z_S, X|T) \\
\leq& H(Z_S|X) + I(X, S|T)
\end{aligned}
\tag{47}
$$

Note that $I(X, S|T)$ is a fixed gap indicating the amount of task-irrelevant information shared between $X$ and $S$.

With Theorem 1, by optimizing the objective Eq. (5), we maximize the lower bound of the task-relevant information $I(Z_S, T)$, and minimize the upper bound of the task-irrelevant information $H(Z_S|T)$. Then the proof is completed. $\qquad\square$

## E  Implementation Details

### E.1  Loss function

In our implementation we did not directly use the original loss function as given in Eqn. (5). For simplicity, we use its equivalent form, which can be easily derived from the following equation:

$$
\begin{aligned}
\left\| \tilde{\mathbf{Z}}_A - \tilde{\mathbf{Z}}_B \right\|_F^2 =& \sum_{k=1}^{D} \sum_{i=1}^{N} (\tilde{z}_{i,k}^A - \tilde{z}_{i,k}^B)^2 \\
=& \sum_{k=1}^{D} \sum_{i=1}^{N} \left( (\tilde{z}_{i,k}^A)^2 + (\tilde{z}_{i,k}^B)^2 - 2 * \tilde{z}_{i,k}^A \tilde{z}_{i,k}^B \right) \\
=& \sum_{k=1}^{D} (2 - 2 * \tilde{Z}_{A,k}^\top \tilde{Z}_{B,k}) \\
=& 2D - 2 * \mathrm{trace}(\tilde{\mathbf{Z}}_A^\top \tilde{\mathbf{Z}}_B)
\end{aligned}
\tag{48}
$$

So we can rewrite the objective function Eqn. (5) as the following one:

$$\mathcal{L} = -\mathrm{trace}(\tilde{\mathbf{Z}}_A^\top \tilde{\mathbf{Z}}_B) + \lambda' \left( \left\| \tilde{\mathbf{Z}}_A^\top \tilde{\mathbf{Z}}_A - \mathbf{I} \right\|_F^2 + \left\| \tilde{\mathbf{Z}}_B^\top \tilde{\mathbf{Z}}_B - \mathbf{I} \right\|_F^2 \right) \tag{49}$$

where the $\lambda'$ here should be half of the $\lambda$ in Eqn. (5). For simplicity we do not discriminate between these two symbols. The values of the trade-off parameter $\lambda$ in Fig 2 as well as that in Appendix E.4 are actually denoted as $\lambda'$ in Eqn. (49).

Table 6: Statistics of benchmark datasets

| Dataset | #Nodes | #Edges | #Classes | #Features |
|---------|--------|--------|----------|-----------|
| Cora | 2,708 | 10,556 | 7 | 1,433 |
| Citeseer | 3,327 | 9,228 | 6 | 3,703 |
| Pubmed | 19,717 | 88,651 | 3 | 500 |
| Coauthor CS | 18,333 | 327,576 | 15 | 6,805 |
| Coauthor Physics | 34,493 | 991,848 | 5 | 8,451 |
| Amazon Computer | 13,752 | 574,418 | 10 | 767 |
| Amazon Photo | 7,650 | 287,326 | 8 | 745 |

## E.2 Graph augmentations

We adopt two random data augmentations strategies on graphs: 1) **Edge dropping**. 2) **Node feature masking**. The two strategies are widely used in node-level contrastive learning [57, 58, 39].

- **Edge dropping** works on the graph structure level, where we randomly remove a portion of edges in the original graph. Formally, given the edge dropping ratio $p_e$, for each edge we have $p_e$ probability to drop this edge from the graph. When calculating the degree for each node, the dropped edge will not be considered.

- **Node feature masking** works on the node feature level, where we randomly set a fraction of features of all nodes as $0$. Formally, given the node feature masking ratio $p_f$, for each input feature, we set it as $0$ with a probability of $p_f$. Note that the masking operation is applied to the selected feature columns of all the nodes.

Note that the previous works [57, 58, 39] use two separate sets of edge dropping ratio $p_e$ and node feature dropping ratio $p_f$ for generating two views, i.e. $p_{e_1}$ and $p_{f_1}$ for view $A$, $p_{e_2}$ and $p_{f_2}$ for view $B$. However, in our implementation, we let $p_{e_1} = p_{e_2}$ and $p_{f_1} = p_{f_2}$, so that the two transformations $t_A$ and $t_B$ come from the same distribution $\mathcal{T}$.

## E.3 Datasets

We evaluate our models on seven node classification benchmarks: *Cora, Citeseer, Pubmed, Coauthor CS, Coauthor Physics, Amazon Computer* and *Amazon Photo*. We provide dataset statistics in Table 6, and brief introduction and settings are as follows:

**Cora[4], Citeseer, Pubmed[5]** are three widely used node classification benchmarks [34, 29]. Each dataset contains one citation network, where nodes mean papers and edges mean citation relationships. We use the public split for linear evaluation, where each class has fixed 20 nodes for training, another fixed 500 nodes and 1000 nodes are for validation/test respectively.

**Coauther CS, Coauther Physics** are co-authorship graphs based on the Microsoft Academic Graph from the KDD Cup 2016 challenge [35]. Nodes are authors, that are connected by an edge if they co-authored a paper; node features represent paper keywords for each author's papers, and class labels indicate most active fields of study for each author. As there is no public split for these datasets, we randomly split the nodes into train/validation/test (10%/10%/80%) sets.

**Amazon Computer, Amazon Photo** are segments of the Amazon co-purchase graph [26], where nodes represent goods, edges indicate that two goods are frequently bought together; node features are bag-of-words encoded product reviews, and class labels are given by the product category. We also use a 10%/10%/80% split for these two datasets.

For all datasets, we use the processed version provided by Deep Graph Library [49][6]. All datasets are public available and do not have licenses.

---

[4]`https://relational.fit.cvut.cz/dataset/CORA`

[5]Citeseer and Pubmed: `https://linqs.soe.ucsc.edu/data`

[6]`https://docs.dgl.ai/en/0.6.x/api/python/dgl.data.html`, Apache License 2.0

Table 7: Details of hyper-parameters of the experimental results in Table 2 and Table 3.

| Dataset | CCA-SSG | | | | | | | | Logistic Regression | |
|---|---|---|---|---|---|---|---|---|---|---|
| | Steps | # layers | # hidden units | $\lambda$ | lr | wd | $p_f$ | $p_e$ | lr | wd |
| Cora | 50 | 2 | 512-512 | 1e-3 | 1e-3 | 0 | 0.1 | 0.4 | 1e-2 | 1e-4 |
| Citeseer | 20 | 1 | 512 | 5e-4 | 1e-3 | 0 | 0.0 | 0.4 | 1e-2 | 1e-2 |
| Pubmed | 100 | 2 | 512-512 | 1e-3 | 1e-3 | 0 | 0.3 | 0.5 | 1e-2 | 1e-4 |
| Computer | 50 | 2 | 512-512 | 5e-4 | 1e-3 | 0 | 0.1 | 0.3 | 1e-2 | 1e-4 |
| Photo | 50 | 2 | 512-512 | 1e-3 | 1e-3 | 0 | 0.2 | 0.3 | 1e-2 | 1e-4 |
| CS[1] | 50 | 2 | 512-512 | 1e-3 | 1e-3 | 0 | 0.2 | - | 5e-3 | 1e-4 |
| Physics | 100 | 2 | 512-512 | 1e-3 | 1e-3 | 0 | 0.5 | 0.5 | 5e-3 | 1e-4 |

[1] We use MLP (instead of GCN) as the encoder on Coauthor-CS, which is essentially equivalent to setting $p_e = 1.0$ (drop all the edges except the self-loops).

In Table 2, we have mentioned that for MVGRL [15] and GRACE [57], we reproduce the experiments with authors' codes, both of which are publicly available: MVGRL[7] and GRACE[8].

### E.4 Hyper-parameters

We provide all the detailed hyper-parameters on the seven benchmarks in Table 7. All hyper-parameters are selected through small grid search, and the search space is provided as follows:

- Training steps: {20, 50, 100, 200}
- Number of layers: {1, 2, 3}
- Number of hidden units: {128, 256, 512, 1024}
- $\lambda$: {1e-4, 5e-4, 1e-3, 5e-3, 1e-2}
- learning rate of CCA-SSG: {5e-4, 1e-3, 5e-3}
- weight decay of CCA-SSG: {0, 1e-5, 1e-4, 1e-3}
- edge dropping ratio: {0, 0.1, 0.2, 0.3, 0.4, 0.5}
- node feature masking ratio: {0, 0.1, 0.2, 0.3, 0.4, 0.5}
- learning rate of logistic regression: {1e-3, 5e-3, 1e-2}
- weight decay of logistic regression: {1e-4, 1e-3, 1e-2}

## F Additional Experiments

### F.1 Visualizations of Correlation Matrix

In Fig. 5, we provide visualizations of the absolute correlation matrix of the raw input features, the embeddings without decorrelation term, and embeddings with decorrelation term on three datasets: *Cora, Citeseer* and *Pubmed*.

As we can see, the raw input feature of the three datasets are all nearly fully uncorrelated (Fig. 5(a), 5(d) and 5(g)). Specifically, the on-diagonal term is close to 1 while the off-diagonal term is close to 0. When training without the decorrelation term Eq. (7), the off-diagonal elements of the correlation matrix of node embeddings increase dramatically as shown in Fig. 5(b) and 5(h), indicating that different dimensions fail to capture orthogonal information. Fig. 5(c) and 5(i) show that with the decorrelation term Eq. (7), our method could learn nearly highly disentangled representations. An interesting finding is that even without the decorrelation term, on Citeseer our method could still generate fairly uncorrelated representations (Fig. 5(e)). The possible reason is that: 1) on Citeseer, we use a one-layer GCN as the encoder, which is less expressive than a two-layer one and alleviate the trend of collapsing. 2) The number of training steps on Citeseer is much smaller than others, so that the impact of invariance term is weaken.

---

[7] https://github.com/kavehhassani/mvgrl, no license.
[8] https://github.com/CRIPAC-DIG/GRACE, Apache License 2.0.

(a) Cora: raw feature　　　(b) Cora: w/o decorrelation　　　(c) Cora: with decorrelation

(d) Citeseer: raw feature　　　(e) Citeseer: w/o decorrelation　　　(f) Citeseer: with decorrelation

(g) Pubemd: raw feature　　　(h) Pubmed: w/o decorrelation　　　(i) Pubmed: with decorrelation
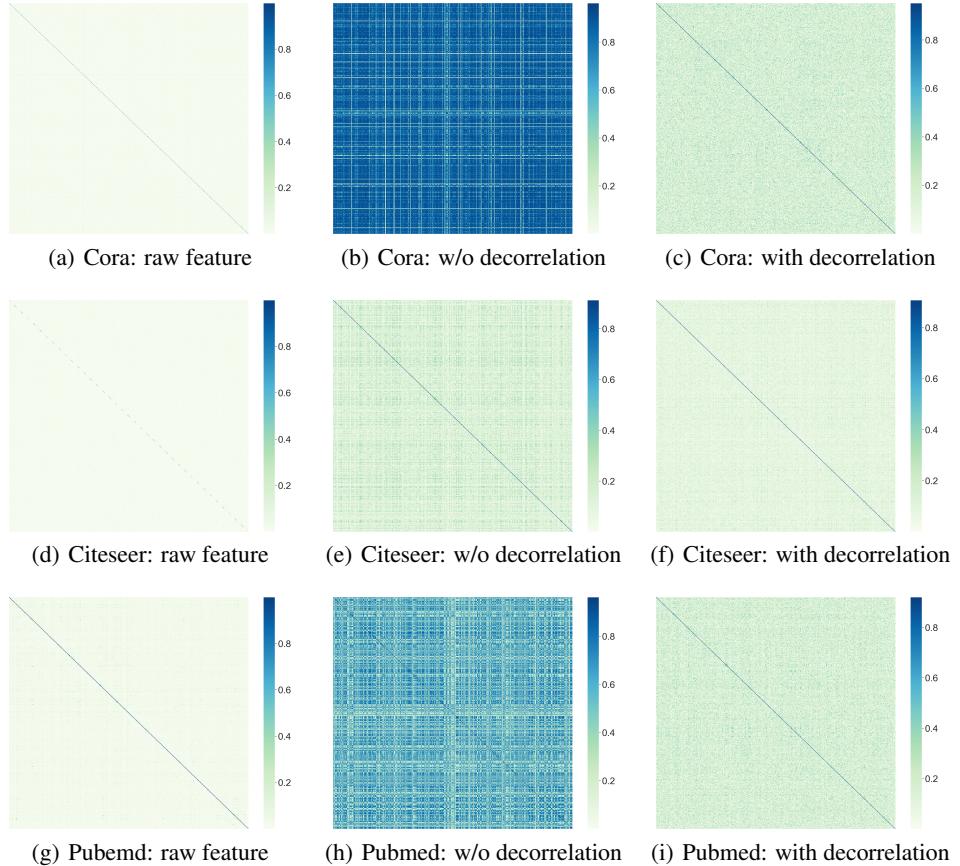
Figure 5: Visualizations of the correlation matrix (absolute value) of the raw input features, the embeddings without decorrelation term, and embeddings with decorrelation term on *Cora, Citeseer* and *Pubmed*. Light green: $\rightarrow 0$; Dark blue: $\rightarrow 1$.

These visualizations also echo the dimensional collapse issue as discussed in Appendix B: without the feature decorrelation term Eq. (7), there is a high probability that all the dimensions capture similar semantic information, thus leading to the dimensional collapse issue. The dimensional collapse can be fundamentally avoided by the decorrelation term Eq. (7).

## F.2　Effects of Augmentation Intensity

We further explore the effects of augmentation intensity on downstream node classification tasks. We try different combinations of the feature masking ratio $p_f$ and edge dropping ratio $p_e$, and report the corresponding performance on the 7 benchmarks mentioned in Appendix E.3. Other hyper-parameters are the same as reported in Table 7. As we can see in Fig. 6, for each dataset there exists an optimal $p_e/p_f$ combination, that could help the model reach the best performance. Also, we find that our method is not that sensitive to the augmentation intensity: as long as $p_e$ and $p_f$ are in a proper range, our method could still achieve impressive and competitive performance. However, it is still very important to select a proper augmentation intensity as well as augmentation method, in order for label-invariant data augmentations for learning informative representations.

## F.3　Performance under Low Label Rates

We further evaluate the node embeddings learned through CCA-SSG on downstream node classification tasks (still using linear, logistic regression), with respect to various label rates (ratio of training nodes). The experiments are conducted on three citation networks: *Cora, Citeseer* and *Pubmed*. In the linear evaluation step, we follow the setups in [24]: we train the linear classifier with 1%, 2%,
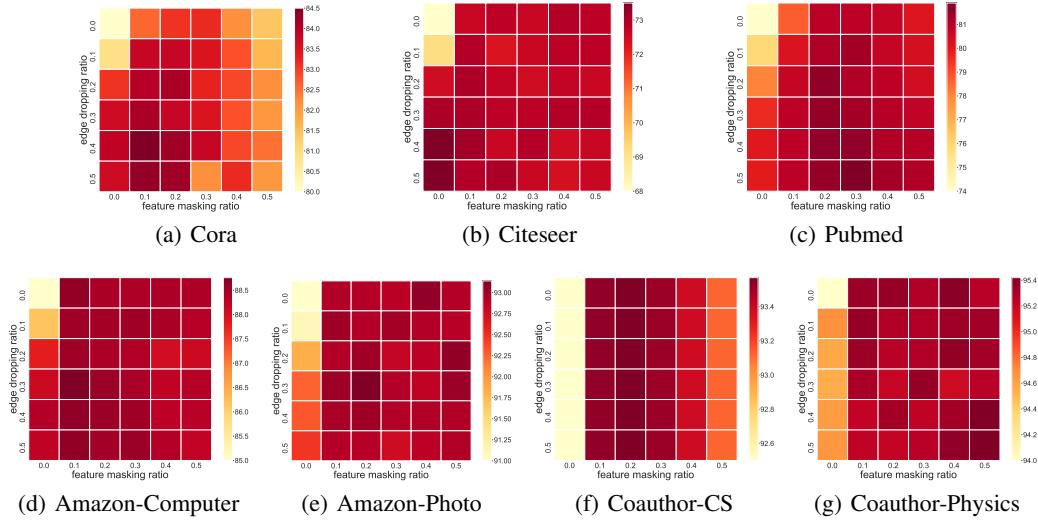
Figure 6: Visualizations of the effects of different augmentation intensity, by adopting different combinations of feature masking ratio $p_f$ and edge dropping ratio $p_e$, and we report test accuracy (%). Each row represents a specific setting for edge dropping ratio $p_e$ and each column represents a specific setting for feature masking ratio $p_f$. Note that when $p_e = p_f = 0$ (the upper left entry in each subfigure), the test accuracy for different datasets should be: 50.2 on *Cora*; 30.5 on *Citeseer*; 46.4 on *Pubmed*; 54.56 on *Computer*; 83.95 on *Photo*; 90.4 on *CS*; 87.85 on *Physics*. Since these results are much worse than the others, we raise their values in each subfigure for better visualization. On *CS*, the edge dropping ratio $p_e$ would make no difference to the performance as we use MLP as the encoder, which does not take graph structure as input.

Table 8: Node classification Accuracy under low label rates (%).

| Dataset | Cora | | | | | Citeseer | | | | | Pubmed | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Label Rate | 1% | 2% | 3% | 4% | 5% | 1% | 2% | 3% | 4% | 5% | 0.05% | 0.1% | 0.3% |
| LP | 62.3 | 65.4 | 67.5 | 69.0 | 70.2 | 40.2 | 43.6 | 45.3 | 46.4 | 47.3 | 66.4 | 65.4 | 66.8 |
| Cheby | 52.0 | 62.4 | 70.8 | 74.1 | 77.6 | 42.8 | 59.9 | 66.2 | 68.3 | 69.3 | 47.3 | 51.2 | 72.8 |
| GCN | 62.3 | 72.2 | 76.5 | 78.4 | 79.7 | 55.3 | 64.9 | 67.5 | 68.7 | 69.6 | 57.5 | 65.9 | 77.8 |
| CCA-SSG | **72.5** | **79.3** | **81.0** | **82.0** | **82.3** | **58.9** | **65.6** | **68.6** | **70.8** | **71.7** | **68.8** | **73.1** | **81.1** |

3%, 4%, 5% (resp. 0.05%, 0.1%, 0.3%) training nodes on *Cora* and *Citeseer* (resp. *Pubmed*), and then test the model with another 1000 nodes. Both training nodes and testing nodes are randomly selected for each trial, and we report the mean accuracy through 20 trials with random splits and random initialization in Table 8.

We compare our method with Label Propagation, GCN with Chebyshev filter(Cheby) and the vanilla GCN [22], whose results are taken from [24] as well. As we can see in Table 8, our method achieves very impressive performance under low label rates, especially when the labeled nodes are really scarce (i.e. 1% on *Cora* and *Citeseer*, 0.05% on *Pubmed*). This is because through self-supervised pretraining, our method could fully utilize the information of unlabeled nodes, and learn good representations for them, which make them easy to distinguish even with only a few number of labeled nodes for training.

## G  Further Comparisons with previous contrastive methods

In Table 1 we have made a thorough comparison with typical contrastive methods from the **technical details**. Here, we further compare our method with more existing contrastive self-supervised graph models (both node-level and graph-level) from the perspective of their general, conceptual designs: 1) How they generate views. 2) The pairs for contrasting. 3) The loss function. 4) Downstream tasks

Table 9: Further conceptual comparison with existing contrastive learning methods on graphs. *View generation in general* (how the method generate views): Cross-scale means this method treat elements in different scales of the graph data as different views (e.g. node and graph); Fix-Diff means using fixed graph diffusion [23] operation to create another view; Rand-Aug means using random graph augmentations (e.g. edge dropping, feature masking, etc.) to generate views. *Pairs* represents the contrasting components, where $N$ is node and $G$ is graph. *Loss* (i.e. the used loss function): NCE represents Noise-Contrastive Estimation [13]; JSD represents Jensen-Shannon Mutual Information Estimator [30]; InfoNCE represents InfoNCE Estimator [46]; MINE means Mutual Information Neural Estimator [3]; BYOL means the asymmetric objective proposed in the BYOL paper [12]. *Tasks* denotes the downstream tasks (node-level, graph-level or edge-level) to which the method has been applied.

| | Methods | View generation in general | Pairs | Loss | Tasks |
|---|---|---|---|---|---|
| Instance-level | DGI [48] | Cross-scale | N-G | NCE | Node |
| | InfoGraph [38] | Cross-scale | N-G | JSD | Graph |
| | MVGRL [15] | Fix-Diff + Cross-scale | N-G | NCE/JSD | Node/Graph |
| | GCC [33] | Rand-Aug | N-N | InfoNCE | Node/Graph |
| | GMI [31] | Hybrid[1] | Hybrid | MINE/JSD | Node/Edge |
| | GRACE [57] | Rand-Aug | N-N | InfoNCE | Node |
| | GraphCL [52] | Rand-Aug | G-G | InfoNCE | Graph |
| | GCA [58] | Rand-Aug | N-N | InfoNCE | Node |
| | CSSL [54] | Rand-Aug | G-G | InfoNCE | Graph |
| | IGSD [55] | Rand-Aug | G-G | BYOL+InfoNCE | Graph |
| | GraphLog [28] | Rand-Aug | N-P-G[2] | InfoNCE | Graph |
| | BGRL [39] | Rand-Aug | N-N | BYOL | Node |
| | MERIT [27] | Fix-Diff + Rand-Aug | N-N | BYOL+InfoNCE | Node |
| | CCA-SSG (Ours) | Rand-Aug | F-F | CCA | Node |

[1] The view generation and contrasting pairs in GMI [31] is unique and complex, and could not be classified into any category.
[2] P denotes hierarchical prototype and could be seen as clustering centroid.

(i.e., node-level, edge-level or graph-level). The comparison is shown in Table 9. Note that this is a high-level comparison with general taxonomy, and each method may have distinct implementation details and specific designs.

We highlight that all of the previous methods focus on contrastive learning at instance level. Our paper proposes a non-contrastive and non-discriminative objective as a new self-supervised representation learning framework, inspired by canonical correlation analysis.