

---

# Graph Contrastive Learning with Cross-view Reconstruction

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Although different graph self-supervised learning strategies have been proposed  
2 to tackle the supervision shortage issue in graph learning tasks, Graph contrastive  
3 learning (GCL) has been the most prevalent approach to this problem. Despite the  
4 remarkable performances those GCL methods have achieved, existing GCL meth-  
5 ods that heavily depend on various manually designed augmentation techniques  
6 still struggle to improve model robustness without risking losing task-relevant  
7 information. Consequently, the learned representation is either brittle or unillumi-  
8 nating. In light of this, we introduce the GraphCV, which follows the information  
9 bottleneck principle to learn minimal yet sufficient representations from graph  
10 data. Specifically, our proposed model elicits the predictive (useful for downstream  
11 instance discrimination) and other non-predictive features separately. Except for  
12 the conventional contrastive loss which guarantees the consistency and sufficiency  
13 of the representations across different augmentation views, we introduce a cross-  
14 view reconstruction mechanism to pursue the disentanglement of the two learned  
15 representations. Besides, an adversarial global view is added as the third view  
16 of contrastive loss to avoid the learned representation from being drafted too far  
17 away from the original distribution. We empirically demonstrate that our proposed  
18 model outperforms the state-of-the-art on graph classification task over multiple  
19 benchmark datasets.

## 20 1 Introduction

21 Graph representation learning (GRL) has attracted significant attention due to its widespread ap-  
22 plications in the real-world interaction systems, such as social, molecules, biological and citation  
23 networks [11]. The current state-of-the-art supervised GRL methods are mostly based on Graph  
24 Neural Networks (GNNs) [17, 35, 10, 41], which require a large amount of task-specific supervised  
25 information. Despite the remarkable performances, they are usually limited by the deficiency of label  
26 supervision in real-world graph data due to the fact that it is usually easy to collect unlabeled graph but  
27 could be very costly to obtain enough annotated label, especially in certain fields, like biochemistry.  
28 Therefore, many recent works [26, 12, 30] have studied how to fully utilize the unlabeled information  
29 on graph and thus stimulating the application of self-supervised learning (SSL) for GRL where only  
30 limited or even no label is needed.

31 As a prevalent and effective strategy of SSL, contrastive learning follows the mutual information  
32 maximization principle (InfoMax) [36] to maximize the agreements of the positive pairs while  
33 minimizing that of the negative pairs in embedding space. However, the graph contrastive learning  
34 paradigm guided with this InfoMax principle has been theoretically and empirically proved to be  
35 insufficient to learn the discriminative and robust representation. To be more specific, state-of-the-art  
36 GCL methods [26, 12, 44] usually implement specific perturbation  $A(\cdot)$  (e.g., Subgraph Sampling,  
37 Node Dropping, Edge Removing and Attributes Masking) on the anchor graph  $G$  to generate its

38 positive pair  $A(G)$  and then train the graph feature encoder to ensure the representation consistency  
 39 within the positive pair, i.e.,  $f(G) = f(A(G))$ . Consequently, such learning strategy is heavily  
 40 dependent on the design of graph augmentation techniques to guarantee the robustness of the learned  
 41 representation. Mild graph augmentation could push encoders to capture redundant and biased  
 42 information [33], which could inadvertently suppress the important predictive features and negatively  
 43 affect the downstream task via the so-called "shortcut" solution [27]. On the other hand, too aggressive  
 44 augmentation may easily lead to another extreme where many predictive features are randomly dropped  
 45 and the learned representations are not informative enough to identify different graphs. Although  
 46 recent works [31, 43, 42, 19] try to implement learnable augmentations on the anchor graph to extract  
 47 the most salient features (those are most easy to learn and resistant to augmentation) to improve  
 48 representation robustness, they still suffer from the difficulty in controlling the augmentation extent  
 49 to balance the representation sufficiency and robustness without the guidance of explicit domain  
 50 knowledge, thus usually lead to unsatisfying performance. To mitigate this issue, a method which  
 51 can disentangle the predictive and the non-predictive latent factors without sacrificing the sufficiency  
 52 of the original predictive graph features is in urgent need.

53 In this paper, we address this challenge by proposing a novel graph contrastive learning model with  
 54 cross-view reconstruction, named GraphCV. Specifically, GraphCV consists of a graph encoder  
 55 followed with two decoders that are trained to extract information particular to the predictive and  
 56 trivial latent factors, respectively. To achieve this optimization objective, we propose a reconstruction-  
 57 based representation learning, including intra-view and inter-view reconstructions, to reconstruct  
 58 the original learned representation with corresponding learned predictive relevant and irrelevant  
 59 representations. Meanwhile, a adversarial graph perturbed from original view is added as the third  
 60 view of the contrastive loss besides the predictive relevant representations of the two augmented  
 61 graph views to ensure the predictive relevant representation maintaining the global semantics instead  
 62 of the partial or even trivial features. We provide theoretical analysis to show that GraphCV is capable  
 63 to learn a minimal sufficient representation with the designs above. Finally, we conduct experiments  
 64 to validate the effectiveness of GraphCV, on the commonly-used graph benchmark datasets. The  
 65 experiment results show that GraphCV achieves significant performance gains over different datasets  
 66 and settings compared with state-of-the-art baselines.

67 To sum up, our main contributions of this work include three aspects: (i) We propose GraphCV to  
 68 learn the disentangled and debiased representation with a cross-view reconstruction mechanism;  
 69 (ii) We provide solid theoretical analysis to analysis the rationality of our designs (iii) We conduct  
 70 thorough experiments to demonstrate that GraphCV can mitigate the "shortcut" solution in contrastive  
 71 learning and significantly outperforms the state-of-the-art baselines over multiple graph classification  
 72 benchmark datasets.

## 73 2 Preliminaries

### 74 2.1 Graph Representation Learning

75 In this work, we focus on the graph-level task, let  $\mathcal{G} = \{G_i = (V_i, E_i)\}_{i=1}^N$  denote a graph dataset  
 76 with  $N$  graphs, where  $V_i$  and  $E_i$  are the node set and edge set of graph  $G_i$ , respectively. We use  
 77  $x_v \in \mathbb{R}^d$  and  $x_e \in \mathbb{R}^d$  to denote the attribute vector of each node  $v \in V_i$  and edge  $e \in E_i$ . Each graph  
 78 is associated with a label, denoted as  $y_i$ , the goal the graph representation learning is to learn an  
 79 encoder  $f : G_i \rightarrow \mathbb{R}^d$  so that the learned representation  $\mathbf{z}_i = f(G_i)$  is sufficient to identify  $y_i$  in the  
 80 downstream task. We clarify sufficiency as  $\mathbf{z}_i$  containing the same amount of information as  $G_i$  for  
 81 label identification [1], and it is formulated as:

$$I(G_i; y_i | \mathbf{z}_i) = 0, \quad (1)$$

82 where  $I(\cdot)$  denotes the mutual information between two variables. We demonstrate the general  
 83 optimization result of classical representation learning in Figure 1(a).

### 84 2.2 Contrastive Learning

85 Contrastive Learning (CL) is a self-supervised representation learning method which leverages  
 86 instance-level identity as supervision. During the training phase, each instance  $x$  firstly goes through  
 87 proper data transformation to generate two data augmentation views  $t_1(x)$  and  $t_2(x)$ , where  $t_1(\cdot)$  and

88  $t_2(\cdot)$  are transformation functions. Then, the CL method learns an encoder  $f$  (a backbone network  
 89 plus a projection layer) which maps  $t_1(x)$  and  $t_2(x)$  closer in the hidden space so that the learned  
 90 representations  $\mathbf{z}_1$  and  $\mathbf{z}_2$  maintain all the information shared by  $t_1(x)$  and  $t_2(x)$ . The encoder is  
 91 usually optimized by contrastive loss, such as NCE loss [38], InfoNCE loss [34] and NT-Xent loss  
 92 [5]. We hereby provide the general optimization result of CL in Figure 1(b). In Graph Contrastive  
 93 Learning (GCL), we usually use GNNs, like GCN [17] and GIN [41], as the backbone networks and  
 94 the commonly-used graph data augmentation operators [44], such as node dropping, edge perturbation,  
 95 subgraph sampling, and attribute masking.

96 All the CL-based methods are built on an assumption that augmentations do not change the infor-  
 97 mation regarding to the label. Here, we follow [7] to clear up the definition of redundancy.  $t_1(G)$   
 98 is redundant to  $t_2(G)$  for  $y$  iff  $t_1(G)$  and  $t_2(G)$  share the same label-relevant information. In CL,  
 99  $t_1(G)$  and  $t_2(G)$  are supposed to be mutually redundant, we define the mutual redundancy as:

$$I(t_1(G); y | t_2(G)) = I(t_2(G); y | t_1(G)) = 0. \quad (2)$$

### 100 2.3 Adversarial Training

101 Deep neural networks have been demonstrated to be vulnerable to adversarial attacks [8]. Among  
 102 all approaches proposed against adversarial attacks, Adversarial Training (AT) achieves remarkable  
 103 robustness. Specifically, AT introduces a perturbation variable  $\delta$  to improve the model robustness  
 104 by training the network on the adversarial samples  $x + \delta$ . During the training phase, the model is  
 105 optimized to minimize the training loss, and the  $\delta$  is optimized within the radius  $\epsilon$  to maximize the  
 106 loss. The supervised setting of AT is defined as:

$$\min_{\theta} \mathbb{E}_{(x,y) \in D} \max_{\|\delta\|_{\infty} \leq \epsilon} \mathcal{L}(x + \delta, y; \theta), \quad (3)$$

107 where  $(x, y)$  are the data feature and label sampled from training set  $D$  respectively, and  $\mathcal{L}$  denotes the  
 108 supervised training objective, such as the cross-entropy loss. Except improving model robustness to  
 109 adversarial attacks, AT is also capable of reducing overfitting and further increasing the generalization  
 110 performance [39, 47]. One possible reason behind this phenomena is that AT follows the Information  
 111 Bottleneck principle [32, 24], in which the optimal representations only contain minimal yet sufficient  
 112 information. Depending on whether relevant to label or not, the mutual information between  
 113 representation  $\mathbf{z}$  and  $x$  can be decomposed into two parts:

$$I(x; \mathbf{z}) = \underbrace{I(y; \mathbf{z})}_{\text{predictive information}} + \underbrace{I(x; \mathbf{z} | y)}_{\text{non-predictive information}}. \quad (4)$$

114 Ideally, the learned representation  $\mathbf{z}$  through AT is expected to keep all the predictive information  
 115 from  $x$  intact while ignoring other non-predictive (trivial) information, we plot this ideal situation in  
 116 Figure 1(c).

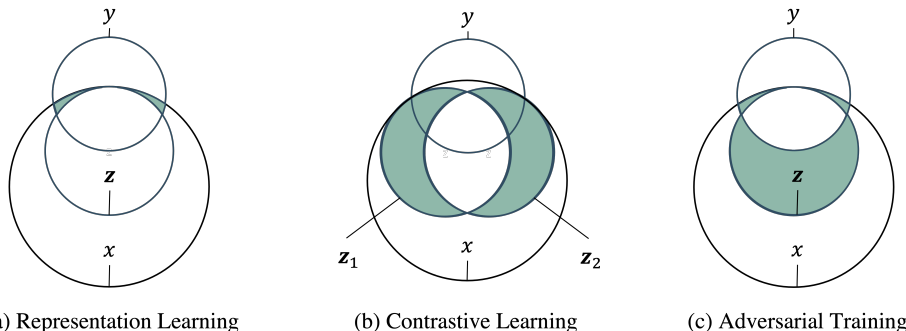


Figure 1: Illustration of the relation between feature  $x$ , label  $y$  and representation  $\mathbf{z}$  in terms of information entropy under the three scenarios above. (a) The ideal optimization result of classical graph representation learning, the green area becomes null when Equation 1 is satisfied. (b) The green area is optimized to null in CL, and  $\mathbf{z}_1$  and  $\mathbf{z}_2$  share the same intersection with  $y$  when the assumption in Equation 2 holds. (c) The green area becomes null when AT achieves its optimal optimization results, in which the second term in Equation 4 is minimized to 0.

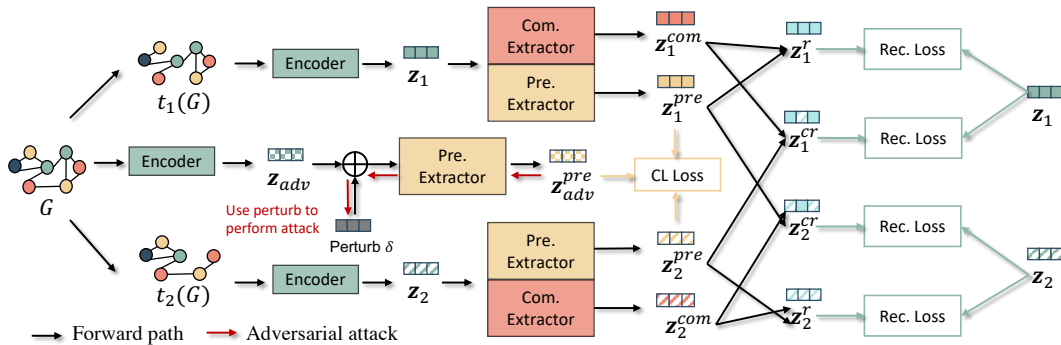


Figure 2: The illustration of proposed GraphCV. (1) Graph augmentations are applied to the input graph  $G$  to produce two augmented graphs, which are then fed into the shared graph encoder  $f(\cdot)$  to generate two graph representations  $\mathbf{z}_1$  and  $\mathbf{z}_2$ . (2)  $\mathbf{z}_1$  and  $\mathbf{z}_2$  are used as the inputs of the two decoder to generate two pairs of graph representations,  $\mathbf{z}^{pre}$  captures the predictive factors and  $\mathbf{z}^{com}$  keep other complementary non-predictive features. Then we use the two pairs of representations to reconstruct  $\mathbf{z}_1$  and  $\mathbf{z}_2$  in both of the intra-view and inter-view. (3) An adversarial sample generated by  $G$  will go through the same procedure to generate  $\mathbf{z}_{adv}^{pre}$ . We take it as the third view besides  $\mathbf{z}_1^{pre}$  and  $\mathbf{z}_2^{pre}$  in CL guarantee the  $\mathbf{z}^{pre}$  can keep the global semantics.

### 117 3 Proposed Model

118 In this section, we introduce the proposed GraphCV, and its framework is shown in Figure 2. Before  
 119 we dive into the details of GraphCV, we will briefly analyze the sufficiency and robustness in GCL,  
 120 and provide illustration of the disentanglement hypothesis.

#### 121 3.1 Motivation of GraphCV

122 To minimize the redundant information in graph representation, the principle of information bottleneck  
 123 (IB) has been introduced in graph representation learning [45], and the learned model is empirically  
 124 proved to be more robust to adversarial attacks. In the circumstances of graph contrastive learning  
 125 (GCL), labels are not accessible to guide the optimization process, thus it is more challenging to  
 126 discern the predictive information and redundant information. For each graph  $G \in \mathcal{G}$ , GCL methods  
 127 generally pick two augmentation operators  $t_1(\cdot)$  and  $t_2(\cdot)$  IID sampled from the same family of  
 128 augmentation  $\mathcal{T}$  to generate two augmented graph  $t_1(G)$  and  $t_2(G)$  for contradistinction. However,  
 129 as we stated in Section 1, strong augmentation will randomly drop a large part of graph features  
 130 and cause overwhelming distribution shift between the two augmented views. As a result, the two  
 131 augmented graph views may fail to hold the mutual redundant assumption in Section 2.2 and thereby  
 132 he learned representation will not satisfy the sufficiency requirement in Section 2.1.

133 To address this dilemma, we propose the GraphCV with reconstruction mechanism to improve its  
 134 robustness without rely on carefully balancing the representation sufficiency and robustness or training  
 135 an extra augmentation module. Given an augmented graph  $t(G)$ , we aim to learn a air of disentangled  
 136 representation  $\mathbf{z} = (\mathbf{z}^{pre}, \mathbf{z}^{com})$ , where  $\mathbf{z}^{pre} \in \mathbb{R}^d$  is expected to be specific to the predictive  
 137 information, while  $\mathbf{z}^{com} \in \mathbb{R}^d$  is optimized to elicit the complementary non-predictive factors from  
 138  $t(G)$ , i.e.,  $\mathbf{z}^{pre} \sim G$ . We provide the illustration of the optimal representation disentanglement in  
 139 Figure 3(a) and (b): (1)  $\mathbf{z}^{pre}$  is sufficient for its corresponding graph view  $t(G)$  regarding to  $y$ , and  
 140 the union of the  $\mathbf{z}^{com}$  and  $\mathbf{z}^{pre}$  cover all the information in  $t(G)$ ; (2)  $\mathbf{z}^{pre}$  and  $\mathbf{z}^{com}$  are optimized  
 141 toward mutually exclusion (disentangled), and  $\mathbf{z}^{pre}$  is trained to capture all the shared information  
 142 between the two augmentation views, thus maintains all the information for label identification  
 143 (under the assumption that implemented mild augmentation does not cause predictive features loss).  
 144 To approximate to this optimal point, we need to propose corresponding designs to guarantee the  
 145 sufficiency and disentanglement of the representation. Next, we will introduce our framework details  
 146 to explain how do we achieve the two restrictions, respectively.

#### 147 3.2 Disentanglement by Cross-View Reconstruction

148 In GCL, we usually leverage a graph encoder to aggregate the feature of graph data as its repre-  
 149 sentation. There are multiple choices of graph encoders in GCL, including GCN [17] and GIN  
 150 [41], etc. In this work, we adopt GIN as the backbone network  $f$  for simplicity. Note that any  
 151 other commonly-used graph encoders can also be applied to our model. Given two augmentation

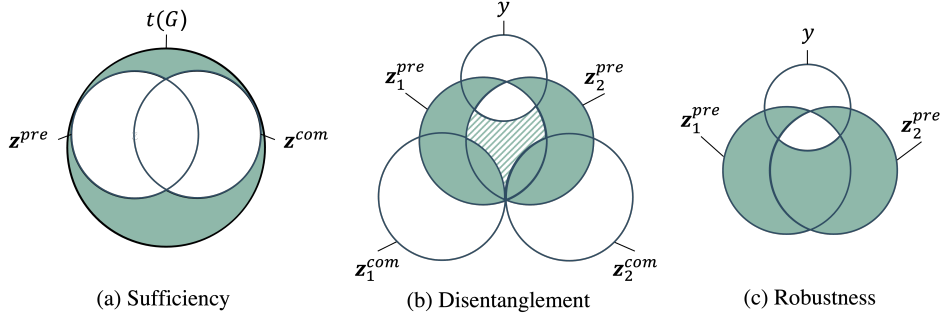


Figure 3: Illustration of the relation between augmented graph  $t(G)$ , label  $y$ , predictive representation  $\mathbf{z}^{pre}$  and non-predictive representation  $\mathbf{z}^{com}$  in terms of information entropy under the optimal situation. The green areas in the three figure become null when the optimal situation is achieved. (a) The union of  $\mathbf{z}^{pre}$  and  $\mathbf{z}^{com}$  covers all the information of its corresponding augmented view  $t(G)$ . (b)  $\mathbf{z}^{pre}$  and  $\mathbf{z}^{com}$  are optimized to approximate mutually exclusion.  $\mathbf{z}_1^{pre}$  and  $\mathbf{z}_2^{pre}$  are optimized to be mutually redundant, including the predictive (white) and part of redundant (shadow) information.  $\mathbf{z}_1^{com}$  and  $\mathbf{z}_2^{com}$  is responsible to capture all the complementary non-predictive information of the two augmented views, where non-overlapping can exist in the two representations since the two randomly sampled augmentation operators can cause the distribution shift between the two views while  $\mathbf{z}^{pre}$  only extract the shared part of the two views. (c)  $\mathbf{z}_1^{pre}$  and  $\mathbf{z}_2^{pre}$  are optimized to more robust, and the redundant information left within them is further minimized.

152 views  $t_1(G)$  and  $t_2(G)$ , we firstly use the encoder  $f(\cdot)$  to map the them into a lower dimension  
 153 hidden space to generate two embeddings  $\mathbf{z}_1$  and  $\mathbf{z}_2$ . Instead of directly maximizing the agreement  
 154 between the two entangled representations  $\mathbf{z}_1$  and  $\mathbf{z}_2$ , we further feed them into a pair decoders  
 155  $g = (g_{pre}, g_{com})$  (both of them are MLP-based networks) to extract the disentangled predictive and  
 156 non-predictive embeddings:

$$[\mathbf{z}^{pre} = g_{pre}(f(t(G))), \mathbf{z}^{com} = g_{com}(f(t(G)))], \quad (5)$$

157 where we can generate a pair of disentangled embeddings for both  $t_1(G)$  and  $t_2(G)$  through the  
 158 procedure above. Ideally, the mutual redundancy assumption between  $(\mathbf{z}_1^{pre}, \mathbf{z}_2^{pre})$  can thus be  
 159 guaranteed because  $t_1(G)$  and  $t_2(G)$  are augmented from the same original graph, and they naturally  
 160 share the same predictive factors, unless aggressive augmentations are implemented. Here, we clarify  
 161 the lower bound of the mutual information between one augmentation view and the learned predictive  
 162 representation of another augmentation view in Theorem 1.

163 **Theorem 1** Suppose  $f(\cdot)$  is a GNN encoder as powerful as 1-WL test. Let  $g_{pre}(\cdot)$  elicits only the  
 164 predictive information from  $\mathbf{z}$  meanwhile  $g_{com}(\cdot)$  extracts the non-predictive factors of  $G$  from  $\mathbf{z}_1$   
 165 and  $\mathbf{z}_2$ . Then we have:

$$I(t_1(G); \mathbf{z}_2^{pre}, \mathbf{z}_2^{com}) \geq I(\mathbf{z}_1^{pre}; \mathbf{z}_2^{pre}) \text{ where } G \in \mathcal{G} \text{ and } t_1(\cdot), t_2(\cdot) \in \mathcal{T}.$$

166 The detailed proof is provided in Section C of Appendix. Therefore, we can maximize the consistency  
 167 between the representations of the two views by maximizing the mutual information of between  $\mathbf{z}_1^{pre}$   
 168 and  $\mathbf{z}_2^{pre}$ . Therefore, we can derive our objective to ensure the view invariance as follow:

$$\mathcal{L}_{pred} = \frac{1}{N} \sum_{i=1}^N \mathcal{L}_{CL}(\mathbf{z}_{1,i}^{pre}, \mathbf{z}_{2,i}^{pre}), \quad (6)$$

169 where  $\mathcal{L}_{CL}(\cdot)$  denotes the contrastive objective and we adopt InfoNCE loss in this work [34]. Mean-  
 170 while, to pursue the feature sufficiency and disentanglement as stated above, we thus propose to use  
 171 the cross-view reconstruction mechanism to approximate these two objectives. To be specific, we  
 172 will use the representation pair  $(\mathbf{z}^{pre}, \mathbf{z}^{com})$  within and cross the augmentation views to recover the  
 173 original raw data so that the two objectives can be guaranteed simultaneously. Due to the reason that  
 174 graph data is a kind of non-Euclidean structured data which can not be represented in the euclidean  
 175 space like the raw data in computer vision domain, we turn to infer the output of  $\mathbf{z} = f(t(G))$  based  
 176 on  $(\mathbf{z}^{com}, \mathbf{z}^{pre})$ . Firstly, we do the reconstruction within the augmentation view, namely mapping  
 177  $(\mathbf{z}_w^{pre}, \mathbf{z}_w^{com})$  to  $\mathbf{z}_w$ , where  $w \in \{1, 2\}$  representing the augmentation view. The optimal result of the  
 178 this step is shown in Figure 3(a), where the joint of  $\mathbf{z}_w^{pre}$  and  $\mathbf{z}_w^{com}$  can cover all the information  
 179 in its corresponding augmentation graph view  $w$ . Since only  $\mathbf{z}^{pre}$  is involved in the contrastive

180 loss in Equation 6, the graph encoder  $f$  is thus not optimized to focus only on the easy-learned  
 181 salient features and less powerful than 1-WL test. Then, we define the  $(\mathbf{z}_w^{pre}, \mathbf{z}_{w'}^{com})$  as a cross-view  
 182 representation pair and the reconstruction procedure will be repeated on it to predict  $\mathbf{z}_{w'}$ , aiming to  
 183 ensure  $\mathbf{z}^{pre}$  and  $\mathbf{z}^{com}$  is optimized to approximate mutual disentanglement, where  $w = 1, w' = 2$   
 184 or  $w = 2, w' = 1$ . The disentanglement optimization process is illustrated in Figure 3(b), where  
 185  $I(\mathbf{z}^{pre}; \mathbf{z}^{com}) \rightarrow 0$  in the ideal situation. Here, we formulate the reconstruction procedures as:

$$\mathbf{z}_w^r = g_{rec}(\mathbf{z}_w^{pre} \odot \mathbf{z}_w^{com}), \quad \mathbf{z}_w^{cr} = g_{rec}(\mathbf{z}_w^{pre} \odot \mathbf{z}_{w'}^{com}), \quad (7)$$

186 where  $g_{rec}$  is the parameterized reconstruction model and  $\odot$  is the pre-defined fusion opera-  
 187 tor, like element-wise product or concatenation. The reconstruction procedures are optimized  
 188 by minimizing the entropy  $H(\mathbf{z}_w | \mathbf{z}_w^{pre}, \mathbf{z}_{w'}^{com})$ , where  $w = w'$  or  $w \neq w'$ . Ideally, we  
 189 can reach the optimal situation demonstrated in Figure 3(a) and (b) iff  $H(\mathbf{z}_w | \mathbf{z}_w^{pre}, \mathbf{z}_{w'}^{com}) =$   
 190  $-\mathbb{E}_{p(\mathbf{z}_w, \mathbf{z}_w^{pre}, \mathbf{z}_{w'}^{com})} [\log p(\mathbf{z}_w | \mathbf{z}_w^{pre}, \mathbf{z}_{w'}^{com})] = 0$ , where  $\mathbf{z}_w$  is exactly recovered given its non-  
 191 predictive representation and the predictive representation of any view. Nevertheless, the condi-  
 192 tion probability  $p(\mathbf{z}_w | \mathbf{z}_w^{pre}, \mathbf{z}_{w'}^{com})$  is unknown for us, we hence use the variation distribution  
 193 approximated by  $g_{rec}$  instead, denoted as  $q(\mathbf{z}_w | \mathbf{z}_w^{pre}, \mathbf{z}_{w'}^{com})$ . We provide the upper bound of  
 194  $H(\mathbf{z}_w | \mathbf{z}_w^{pre}, \mathbf{z}_{w'}^{com})$  in Theorem 2.

195 **Theorem 2** Assume  $q$  is a Gaussian distribution,  $g_r$  is the parameterized reconstruction model which  
 196 infer  $\mathbf{z}_w$  from  $(\mathbf{z}_w^{pre}, \mathbf{z}_{w'}^{com})$ . Then we have:

$$H(\mathbf{z}_w | \mathbf{z}_w^{pre}, \mathbf{z}_{w'}^{com}) \leq \|\mathbf{z}_w - g_r(\mathbf{z}_w^{pre} \odot \mathbf{z}_{w'}^{com})\|_2^2 \text{ where } w = w' \text{ or } w \neq w'.$$

197 The detailed proof is demonstrated in Section C of Appendix. Since we adopt two augumentation  
 198 views, the objective function constraining representation sufficiency and disentanglement can be  
 199 formulated as:

$$\mathcal{L}_{recon} = \frac{1}{2N} \sum_{i=1}^N \sum_{w=1}^2 \left[ \|\mathbf{z}_{w,i} - \mathbf{z}_{w,i}^r\|_2^2 + \|\mathbf{z}_{w,i} - \mathbf{z}_{w,i}^{cr}\|_2^2 \right]. \quad (8)$$

### 200 3.3 Adversarial Contrastive View

201 With the cross-view reconstruction mechanism above, the two learned representations stated above  
 202 are optimized towards the disentangled manner. However, we still need to ensure  $\mathbf{z}^{pre}$  is more  
 203 relevant to predictive information than  $\mathbf{z}^{com}$ . Therefore, we extend the Equation 6 to three contrastive  
 204 views and add an extra global view without dropping any features as the third views to guarantee  
 205 the learned  $\mathbf{z}^{pre}$  main the global semantics instead of partial or even trivial features, i.e.,  $\mathbf{z}_1^{com} \sim G$   
 206 and  $\mathbf{z}_2^{com} \sim G$ . During the experiments, we find a adversarial graph sample perturbed from original  
 207 graph view can not only play the same role of original view in our expectation, but also achieve  
 208 better robustness. A possible explanation is that there is still redundant information left the shared  
 209 part of the two augumentation views as illustrated in Figure 3(b), especially when the implemented  
 210 augmentations are not very aggressive. Thus, an adversarial view can further eliminate the trivial  
 211 information from the learned  $\mathbf{z}^{pre}$  for better generalization ability. We define the adversarial objective  
 212 as follows:

$$\delta^* = \operatorname{argmax}_{\|\delta\|_\infty \leq \epsilon} \mathcal{L}_{adv}(t_1(G), t_2(G), G + \delta), \quad (9)$$

213 where the adversarial sample  $G + \delta$  together with the two augumentation views, i.e.,  $t_1(G)$  and  
 214  $t_2(G)$  are employed as the positive pair. Our implementation of crafting perturbation is spurred  
 215 by recent work [42] that add perturbation  $\delta$  on the output of first hidden layer  $\mathbf{h}^{(1)}$  because it is  
 216 empirically proved to generate more challenging view than adding perturbation on the initial node  
 217 feature. Therefore, the adversarial contrastive objective is defined as:

$$\mathcal{L}_{adv} = \frac{1}{N} \sum_{i=1}^N \max_{\delta^*} [\mathcal{L}_{CL}(\mathbf{z}_{1,i}^{com}, G + \delta^*) + \mathcal{L}_{CL}(\mathbf{z}_{2,i}^{com}, G + \delta^*)]. \quad (10)$$

218 where the optimized perturbation  $\delta'$  is solved by projected gradient descent (PGD) [22]

### 219 3.4 The Joint Objective

220 We design the joint objective of GraphCV by combining all of objectives above together. Given the  
 221 graph  $G \in \mathbf{G}$ , the graph encoder  $f$  and all the decoders  $g$  can be optimized with the objective below:

$$\min_{f,g} \mathbb{E}_{G \in \mathbf{G}} \left[ \mathcal{L}_{pre} + \lambda_r \mathcal{L}_{recon} + \lambda_a \max_{\|\delta\|_\infty \leq \epsilon} \mathcal{L}_{adv} \right], \quad (11)$$

222 where  $\lambda_r$  and  $\lambda_a$  are the coefficients to balance the magnitude of each loss term. Our proposed model  
 223 is able to learn optimal representation illustrated in Figure 3(c) with the joint objective. We present  
 224 the training algorithm of our model in Section F of Appendix.

## 225 4 Experiments

226 In this section, we demonstrate the empirical evaluation results of GraphCV on public graph bench-  
 227 mark datasets. Ablation study and robustness analysis are conducted to evaluate the effectiveness of  
 228 the designs in GraphCV. We provide the dataset statistics, training details and more analysis about  
 229 hyper-parameters and representation disentanglement can be found in the Appendix.

### 230 4.1 Experimental Setups

231 **Datasets.** We evaluate our model on five graph benchmark datasets from the field of bioinformatics,  
 232 including MUTAG, PTC-MR, NCI1, DD, and PROTEINS, and other five from the field of social  
 233 network, which are COLLAB, IMDB-B, RDT-B, RDT-M5K, and IMDB-M, for the task of graph-  
 234 level property classification. Additionally, We use ogbg-molhiv from Open Graph Benchmark Dataset  
 235 [15] to demonstrate our model’s advantages over large-scale dataset. More details about dataset  
 236 statistics are included in Section A of Appendix.

237 **Baselines.** Under the unsupervised representation learning setting, we compare GraphCV with  
 238 the seven SOTA self-supervised learning methods GraphCL [44], InfoGraph[30], MVGRL [12],  
 239 AD-GCL[31], GASSL[42], InfoGCL[40] and DGCL[18], as well as four classical unsupervised rep-  
 240 resentation learning methods, including node2vec [9], sub2vec [2], graph2vec [25], and GVAE[16].

241 **Evaluation Protocol.** We follow the evaluation protocols in the previous works [30, 44, 18] to verify  
 242 the effectiveness of our model. The learned representation is fine-tuned by a linear SVM classifier for  
 243 task-specific prediction. We report the mean test accuracy evaluated by a 10-fold cross validation with  
 244 standard deviation of five random seeds as the final performance. In addition, we follow the setting of  
 245 semi-supervised representation learning from GraphCL on the ogbg-molhiv dataset, with the finetune  
 246 label rates as 1%, 10%, and 20%. The final performance is reported as the mean ROC-AUC of five  
 247 initialization random seeds

248 **Implementation Details.** We implement our framework with PyTorch and employ the data  
 249 augmentation function provided by PyGCL library [48]. We choose GIN [41] as the backbone encoder  
 250 and the model is optimized by Adam optimizer. There are two specific hyper-parameters in our model,  
 251 namely  $\lambda_r$  and  $\lambda_a$ , the search space of them are  $\{0.0, 1.0, 3.0, 5.0, 10.0\}$  and  $\{0.0, 0.5, 1.0, 1.5, 2.0\}$ ,  
 252 respectively. More details about implementation details is provided in the Section B of Appendix.  
 253 All of the experiments are conducted on Nvidia GeForce RTX 2080ti GPU.

### 254 4.2 Overall Performance Comparison

255 **Unsupervised representation learning.** The overall performance comparison is shown in Table 1.  
 256 From the results, we can have three observations: (1) The GCL-based methods generally yield higher  
 257 performances than classical unsupervised learning methods, indicating the effectiveness of utilizing  
 258 instance-level supervision; (2) InfoGCL and GASSL achieve better performances than GraphCL,  
 259 which empirically proves the conclusion that InfoMax object could suffer from the overwhelmed  
 260 information and thus more challenging augmentations or perturbations are in need to produce robust  
 261 representations; (3) Our proposed GraphCV and DGCL consistently outperform other baselines,  
 262 proving the advantage of disentangled representation. More importantly, ADGCL achieves state-of-  
 263 the-art results on most of the datasets, which further demonstrate the success of our model to learn  
 264 minimal yet sufficient representations.

#### 265 Semi-supervised representation learning.

266 The semi-supervised representation learning  
 267 results for ogbg-molhiv are shown in Fig-  
 268 ure 4. It is obvious that our model gains sig-  
 269 nificant improvements under the three label-  
 270 rate fine-tuning settings. We also notice that  
 271 as the label rate increases, the amount of im-  
 272 provement increases as well (1%, 1.8%, and  
 273 4.4% for label rate 1%, 10%, and 20%, re-  
 274 spectively). A possible explanation could be  
 275 that as more trainable data is included in the  
 276 process of fine-tuning when the label rate

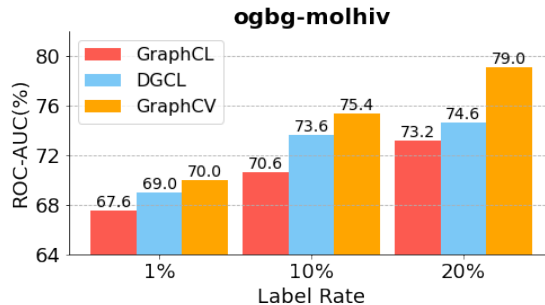


Figure 4: Performance comparison of semi-supervised learning on ogbg-molhiv.

Table 1: Overall comparison on multiple graph classification benchmarks. Results are reported as mean $\pm$ std%, the best performance is bolded and runner-ups are underlined. "-" indicates the result is not reported in original papers.

	MUTAG	PTC-MR	COLLAB	NCII	PROTEINS	IMDB-B	RDT-B	IMDB-M	RDT-M5K	DD
node2vec	72.6 $\pm$ 10.2	58.6 $\pm$ 8.0	-	54.9 $\pm$ 1.6	57.5 $\pm$ 3.6	-	-	-	-	-
sub2vec	61.1 $\pm$ 15.8	60.0 $\pm$ 6.4	-	52.8 $\pm$ 1.5	53.0 $\pm$ 5.6	55.3 $\pm$ 1.5	71.5 $\pm$ 0.4	36.7 $\pm$ 0.8	36.7 $\pm$ 0.4	-
graph2vec	83.2 $\pm$ 9.3	60.2 $\pm$ 6.9	-	73.2 $\pm$ 1.8	73.3 $\pm$ 2.1	71.1 $\pm$ 0.5	75.8 $\pm$ 1.0	50.4 $\pm$ 0.9	47.9 $\pm$ 0.3	-
InfoGraph	89.0 $\pm$ 1.1	61.7 $\pm$ 1.4	70.7 $\pm$ 1.1	76.2 $\pm$ 1.1	74.4 $\pm$ 0.3	73.0 $\pm$ 0.9	82.5 $\pm$ 1.4	49.7 $\pm$ 0.5	53.5 $\pm$ 1.0	72.9 $\pm$ 1.8
VGAE	87.7 $\pm$ 0.7	61.2 $\pm$ 1.8	-	-	-	70.7 $\pm$ 0.7	87.1 $\pm$ 0.1	49.3 $\pm$ 0.4	52.8 $\pm$ 0.2	-
MVGRL	89.7 $\pm$ 1.1	62.5 $\pm$ 1.7	-	-	-	74.2 $\pm$ 0.7	84.5 $\pm$ 0.6	51.2 $\pm$ 0.5	-	-
GraphCL	86.8 $\pm$ 1.3	63.6 $\pm$ 1.8	71.4 $\pm$ 1.2	77.9 $\pm$ 0.4	74.4 $\pm$ 0.5	71.1 $\pm$ 0.4	89.5 $\pm$ 0.8	-	56.0 $\pm$ 0.3	<b>78.6<math>\pm</math>0.4</b>
InfoGCL	91.2 $\pm$ 1.3	63.5 $\pm$ 1.5	80.0 $\pm$ 1.3	80.2 $\pm$ 0.6	-	75.1 $\pm$ 0.9	-	51.4 $\pm$ 0.8	-	-
DGCL	92.1 $\pm$ 0.8	<u>65.8<math>\pm</math>1.5</u>	<b>81.2<math>\pm</math>0.3</b>	81.9 $\pm$ 0.2	76.4 $\pm$ 0.5	75.9 $\pm$ 0.7	91.8 $\pm$ 0.2	51.9 $\pm$ 0.4	56.1 $\pm$ 0.2	-
AD-GCL	89.7 $\pm$ 1.0	-	73.3 $\pm$ 0.6	69.7 $\pm$ 0.5	73.8 $\pm$ 0.5	72.3 $\pm$ 0.6	85.5 $\pm$ 0.8	49.9 $\pm$ 0.7	54.9 $\pm$ 0.4	75.1 $\pm$ 0.4
GASSL	90.9	64.6 $\pm$ 6.1	78	80.2	-	74.2	-	51.7	-	-
<b>GraphCV</b>	<b>92.6<math>\pm</math>0.9</b>	<b>67.4<math>\pm</math>1.3</b>	<u>80.5<math>\pm</math>0.5</u>	<b>82.0<math>\pm</math>1.0</b>	<b>77.3<math>\pm</math>0.4</b>	<b>76.7<math>\pm</math>0.5</b>	<b>92.4<math>\pm</math>0.9</b>	<b>52.2<math>\pm</math>0.5</b>	<b>57.2<math>\pm</math>0.4</b>	<b>80.5<math>\pm</math>0.5</b>

277 increases, so does the affiliated redundant information, which as a result, deteriorate the performance  
 278 even more. Therefore, removing redundant information causes a higher performance boost.

### 279 4.3 Ablation Study

280 To further verify the effectiveness of different modules in GraphCV, we perform ablation studies on  
 281 each one of the module by creating the model variants illustrated below. The comparison results are  
 282 shown in Table 2.

- 283 • **w/o Intra-view Recon.** Reconstruction is only executed within the cross view i.e.,  $w \neq w'$ .
- 284 • **w/o Inter-view Recon.** Reconstruction is only executed within the same view i.e.,  $w = w'$ .
- 285 • **w/o Adv. Training.** Adversarial view is discarded in the contrastive loss.

286 From Table 2 we can see that our model with the combination of cross-view reconstruction and  
 287 adversarial training module outperforms all of the variants. Discarding any reconstruction view  
 288 could cause the failure to reach the optimal situation illustrated in Figure 3. We can not guarantee  
 289 the representation disentanglement assumption if we skip the inter-view reconstruction, and the  
 290 sufficiency assumption may not hold if we abandon intra-view reconstruction. Either way, the  
 291 predictive representations may suffer from enormous information loss during the contrastive learning  
 292 and further lead to the performance deterioration. Compared with our model, the variant w/o Adv.  
 293 Training may bring too much redundant information to the downstream classifier, therefore creating  
 294 more confusions. The relatively larger performance deterioration for the two variants w/o Intra-  
 295 view Recon and w/o Inter-view suggests the rule "better than nothing". That is, having redundant  
 296 information is better than having it partially.

Table 2: Overall comparison of the model variants' performance. Results are reported as mean $\pm$ std%, the best performance is bolded.

	MUTAG	PTC-MR	COLLAB	NCII	PROTEINS	IMDB-B	RDT-B	IMDB-M	RDT-M5K	DD
w/o Intra Recon	91.5 $\pm$ 1.2	65.8 $\pm$ 1.3	78.4 $\pm$ 0.7	79.6 $\pm$ 0.7	75.6 $\pm$ 0.5	75.4 $\pm$ 0.8	92.0 $\pm$ 0.4	51.5 $\pm$ 0.4	55.8 $\pm$ 0.6	79.3 $\pm$ 0.7
w/o Inter Recon	91.0 $\pm$ 0.9	64.7 $\pm$ 1.4	78.0 $\pm$ 0.8	78.7 $\pm$ 1.2	74.9 $\pm$ 0.7	75.0 $\pm$ 0.6	91.1 $\pm$ 0.7	50.8 $\pm$ 0.2	55.6 $\pm$ 0.4	79.0 $\pm$ 0.8
w/o Adv. Training	92.1 $\pm$ 0.6	66.8 $\pm$ 0.5	80.3 $\pm$ 0.5	81.2 $\pm$ 0.9	77.0 $\pm$ 0.3	76.4 $\pm$ 0.6	92.2 $\pm$ 1.0	52.0 $\pm$ 0.4	56.8 $\pm$ 0.5	80.1 $\pm$ 0.6
<b>GraphCV</b>	<b>92.6<math>\pm</math>0.9</b>	<b>67.4<math>\pm</math>0.5</b>	<b>80.5<math>\pm</math>0.5</b>	<b>82.0<math>\pm</math>1.0</b>	<b>77.3<math>\pm</math>0.4</b>	<b>76.7<math>\pm</math>0.5</b>	<b>92.5<math>\pm</math>0.9</b>	<b>52.2<math>\pm</math>0.5</b>	<b>57.2<math>\pm</math>0.4</b>	<b>80.5<math>\pm</math>0.5</b>

### 297 4.4 Robustness Analysis

298 In this section, we conduct extra experiments on ogbg-molhiv dataset to evaluate the effectiveness of  
 299 our design in ensuring the representation robustness under aggressive augmentation and perturbation.  
 300 The results are shown in Figure 5. In the left two subplots, we plot accuracy verses edge perturbation  
 301 and attribute masking strengths, respectively. Specifically, we keep the GraphCL and our proposed  
 302 GraphCV under the same hyper-parameter setting and set the  $\lambda_r$  and  $\lambda_a$  of GraphCV as 5.0 and  
 303 0.5, respectively. From the results we can see that GraphCV not only consistently outperforms  
 304 GraphCL but also is less affected by larger augmentation strengths. Similar observation can be find  
 305 in the right two subplots, where we compare our method with GASSL under different perturbation  
 306 bounds and attack steps to demonstrate its robustness against adversarial attacks. Since both our  
 307 model and GASSL use GIN as the backbone network, we hereby add the performance of GIN as the  
 308 compared baseline. Although aggressive adversarial attacks can largely deteriorate the performance,  
 309 our proposed GraphCV still achieves more robust performance than GASSL.

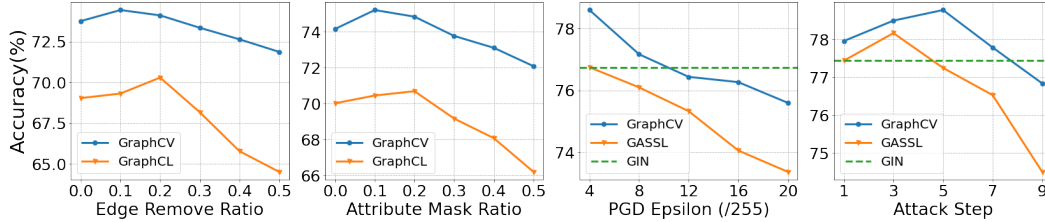


Figure 5: Performance versus augmentation strengths, perturbation bound and attack step.

## 310 5 Related Work

311 **Graph contrastive learning.** Contrastive learning is firstly proposed in the compute vision field [5]  
 312 and raises a surge of interests in the area of self-supervised graph representation learning for the past  
 313 few years. The principle behind contrastive learning is to utilize the instance-level identity as supervi-  
 314 sion and maximize the consistency between positive pairs in hidden space through designed contrast  
 315 mode. Previous graph contrastive learning works generally rely on various graph augmentation  
 316 (transformation) techniques [36, 26, 13, 44, 30] to generate positive pair from original data as similar  
 317 samples. Recent works in this field try to improve the effectiveness of graph contrastive learning by  
 318 finding more challenge view [31, 40, 43] or adding adversarial perturbation [42]. However, most of  
 319 the existing methods contrast over entangled embeddings, where the complex intertwined information  
 320 may pose obstacles to extracting useful information for downstream tasks. Our model is spared from  
 321 the issue by contrasting over disentangled representations.

322 **Disentangled representation learning on graphs.** Disentangled representation learning arises from  
 323 the computer vision field [14, 46] to disentangle the heterogeneous latent factors of the representations,  
 324 and therefore making the representations more robust and interpretable [4]. This idea has now been  
 325 widely adopted in graph representation learning. [20, 21] utilizes neighborhood routing mechanism  
 326 to identify the latent factors in the node representations. Some other generative models [16, 29]  
 327 utilize Variational Autoencoders to balance reconstruction and disentanglement. Recent work [18]  
 328 outspreads the application of disentangled representations learning in self-supervised graph learning  
 329 by contrasting the factorized representations. Although these methods gain significant benefit from  
 330 the representation disentanglement, the underlined excessive information could still overload the  
 331 model, thus resulting in limited capacities. Our model targets the issue by removing the redundant  
 332 information that is considered irrelevant to the graph property.

333 **Graph information bottleneck.** The Information bottleneck (IB) [32] has been widely adopted  
 334 as a critical principle of representation learning. A representation contains minimal yet sufficient  
 335 information is considered to be in compliance with the IB principle and many works [3, 28, 7] have  
 336 empirically and theoretically proved that representation agree with IB principle is both informative  
 337 and robust. Recently, IB principle is also borrowed to guide the representation learning of graph  
 338 structure data. Current methods [37, 40, 31] usually propose different regularization designs to  
 339 learn compressed yet informative representations in accordance with IB principle. We follow the  
 340 information bottleneck to learn the expressive and robust representation from disentangled information  
 341 in this work.

## 342 6 Conclusion

343 In this paper, we study graph representation learning in light of information bottleneck. To reach the  
 344 optimum we illustrate, we propose a novel model, namely GraphCV, which is designed to disentangle  
 345 the essential factors from augmented graph through a cross-view reconstruction mechanism so that the  
 346 information entanglement brought by augmentations will not cause the loss of predictive information  
 347 during contrastive learning. We also add an adversarial view as the third view of the contrastive  
 348 learning to further remove redundant information and enhance representation robustness. In addition,  
 349 we theoretically analyze the effectiveness of each component in our model and derive the objective  
 350 based on the analysis. Extensive experiments on multiple graph benchmark datasets and different  
 351 settings prove the ability of GraphCV to learn robust and informative graph representation. In the  
 352 future, we can explore how to come up with a practical objective to further decrease the upper bound  
 353 of the mutual information between the disentangled representations and try to utilize more efficient  
 354 training strategy to make the proposed model more time-saving on large-scale graphs.

## 355 References

- 356 [1] Alessandro Achille and Stefano Soatto. Emergence of invariance and disentanglement in deep  
357 representations. *JMLR*, 2018.
- 358 [2] Bijaya Adhikari, Yao Zhang, Naren Ramakrishnan, and B. Aditya Prakash. Sub2Vec: Feature  
359 Learning for Subgraphs. In *KDD*, 2018.
- 360 [3] Alexander A. Alemi, Ian Fischer, Joshua V. Dillon, and Kevin Murphy. Deep Variational  
361 Information Bottleneck. *ICLR*, 2017.
- 362 [4] Yoshua Bengio, Aaron Courville, and Pascal Vincent. Representation Learning: A Review and  
363 New Perspectives. *TPAMI*, 2013.
- 364 [5] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A Simple Framework  
365 for Contrastive Learning of Visual Representations. In *ICML*, 2020.
- 366 [6] Thomas M Cover. *Elements of information theory*. John Wiley & Sons, 1999.
- 367 [7] Marco Federici, Anjan Dutta, Patrick Forré, Nate Kushman, and Zeynep Akata. Learning robust  
368 representations via multi-view information bottleneck. In *ICLR*, 2020.
- 369 [8] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversar-  
370 ial examples. In *ICLR*, 2014.
- 371 [9] Aditya Grover and Jure Leskovec. node2vec: Scalable Feature Learning for Networks. In *KDD*,  
372 2016.
- 373 [10] Will Hamilton, Zhitao Ying, and Jure Leskovec. Inductive Representation Learning on Large  
374 Graphs. In *NeurIPS*, 2017.
- 375 [11] William L Hamilton, Rex Ying, and Jure Leskovec. Representation learning on graphs: Methods  
376 and applications. *arXiv preprint arXiv:1709.05584*, 2017.
- 377 [12] Kaveh Hassani and Amir Hosein Khasahmadi. Contrastive Multi-View Representation Learning  
378 on Graphs. In *ICML*, 2020.
- 379 [13] Kaveh Hassani and Amir Hosein Khasahmadi. Contrastive multi-view representation learning  
380 on graphs. In *ICML*, 2020.
- 381 [14] Jun-Ting Hsieh, Bingbin Liu, De-An Huang, Li F Fei-Fei, and Juan Carlos Niebles. Learning  
382 to Decompose and Disentangle Representations for Video Prediction. In *NeurIPS*, 2018.
- 383 [15] Weihua Hu, Matthias Fey, Marinka Zitnik, Yuxiao Dong, Hongyu Ren, Bowen Liu, Michele  
384 Catasta, and Jure Leskovec. Open Graph Benchmark: Datasets for Machine Learning on Graphs.  
385 In *NeurIPS*, 2020.
- 386 [16] Thomas N. Kipf and Max Welling. Variational Graph Auto-Encoders. In *NeurIPS*, 2016.
- 387 [17] Thomas N. Kipf and Max Welling. Semi-Supervised Classification with Graph Convolutional  
388 Networks. In *ICLR*, 2017.
- 389 [18] Haoyang Li, Xin Wang, Ziwei Zhang, Zehuan Yuan, Hang Li, and Wenwu Zhu. Disentangled  
390 Contrastive Learning on Graphs. In *NeurIPS*, 2021.
- 391 [19] Sihang Li, Xiang Wang, An Zhang, Xiangnan He, and Tat-Seng Chua. Let invariant rationale  
392 discovery inspire graph contrastive learning. In *ICML*, 2022.
- 393 [20] Yanbei Liu, Xiao Wang, Shu Wu, and Zhitao Xiao. Independence promoted graph disentangled  
394 networks. In *AAAI*, 2020.
- 395 [21] Jianxin Ma, Peng Cui, Kun Kuang, Xin Wang, and Wenwu Zhu. Disentangled Graph Convolu-  
396 tional Networks. In *ICML*, 2019.
- 397 [22] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu.  
398 Towards deep learning models resistant to adversarial attacks. In *ICLR*, 2018.

- 399 [23] Christopher Morris, Nils M. Kriege, Franka Bause, Kristian Kersting, Petra Mutzel, and Marion  
400 Neumann. Tudataset: A collection of benchmark datasets for learning with graphs. In *ICML*  
401 *2020 Workshop on Graph Representation Learning and Beyond (GRL+ 2020)*, 2020.
- 402 [24] Daniel Moyer, Shuyang Gao, Rob Brekelmans, Aram Galstyan, and Greg Ver Steeg. Invariant  
403 representations without adversarial training. *NeurIPS*, 2018.
- 404 [25] A. Narayanan, Mahinthan Chandramohan, R. Venkatesan, Lihui Chen, Yang Liu, and Shantanu  
405 Jaiswal. graph2vec: Learning Distributed Representations of Graphs. *ArXiv*, 2017.
- 406 [26] Jiezhong Qiu, Qibin Chen, Yuxiao Dong, Jing Zhang, Hongxia Yang, Ming Ding, Kuansan  
407 Wang, and Jie Tang. GCC: Graph Contrastive Coding for Graph Neural Network Pre-Training.  
408 In *KDD*, 2020.
- 409 [27] Joshua Robinson, Li Sun, Ke Yu, Kayhan Batmanghelich, Stefanie Jegelka, and Suvrit Sra. Can  
410 contrastive learning avoid shortcut solutions? In *NeurIPS*, 2021.
- 411 [28] Ravid Shwartz-Ziv and Naftali Tishby. Opening the Black Box of Deep Neural Networks via  
412 Information. *arXiv:1703.00810 [cs]*, April 2017.
- 413 [29] Martin Simonovsky and Nikos Komodakis. GraphVAE: Towards generation of small graphs  
414 using variational autoencoders. In *ICLR*, 2018.
- 415 [30] Fan-Yun Sun, Jordan Hoffmann, Vikas Verma, and Jian Tang. InfoGraph: Unsupervised and  
416 Semi-supervised Graph-Level Representation Learning via Mutual Information Maximization.  
417 In *ICLR*, 2019.
- 418 [31] Susheel Suresh, Pan Li, Cong Hao, and Jennifer Neville. Adversarial Graph Augmentation to  
419 Improve Graph Contrastive Learning. In *NeurIPS*, 2021.
- 420 [32] Naftali Tishby, Fernando C Pereira, and William Bialek. The information bottleneck method.  
421 *arXiv preprint physics/0004057*, 2000.
- 422 [33] Michael Tschannen, Josip Djolonga, Paul K Rubenstein, Sylvain Gelly, and Mario Lucic. On  
423 mutual information maximization for representation learning. In *ICLR*, 2019.
- 424 [34] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation Learning with Contrastive  
425 Predictive Coding. *arXiv e-prints*, 2018.
- 426 [35] Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua  
427 Bengio. Graph Attention Networks. In *ICLR*, 2018.
- 428 [36] Petar Veličković, William Fedus, William L. Hamilton, Pietro Liò, Yoshua Bengio, and R. Devon  
429 Hjelm. Deep Graph Infomax. In *ICLR*, 2019.
- 430 [37] Tailin Wu, Hongyu Ren, Pan Li, and Jure Leskovec. Graph Information Bottleneck. In *NeurIPS*.  
431 Curran Associates, Inc., 2020.
- 432 [38] Zhirong Wu, Yuanjun Xiong, Stella X. Yu, and Dahua Lin. Unsupervised Feature Learning via  
433 Non-Parametric Instance Discrimination. In *CVPR*, 2018.
- 434 [39] Cihang Xie, Mingxing Tan, Boqing Gong, Jiang Wang, Alan L Yuille, and Quoc V Le. Adver-  
435 sarial examples improve image recognition. In *CVPR*, 2020.
- 436 [40] Dongkuan Xu, Wei Cheng, Dongsheng Luo, Haifeng Chen, and Xiang Zhang. InfoGCL:  
437 Information-Aware Graph Contrastive Learning. In *NeurIPS*, 2021.
- 438 [41] Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How Powerful are Graph Neural  
439 Networks? In *ICLR*, 2019.
- 440 [42] Longqi Yang, Liangliang Zhang, and Wenjing Yang. Graph Adversarial Self-Supervised  
441 Learning. In *NeurIPS*, 2021.
- 442 [43] Yuning You, Tianlong Chen, Yang Shen, and Zhangyang Wang. Graph contrastive learning  
443 automated. In *ICLR*, 2021.

- 444 [44] Yuning You, Tianlong Chen, Yongduo Sui, Ting Chen, Zhangyang Wang, and Yang Shen.  
445 Graph Contrastive Learning with Augmentations. In *NeurIPS*, 2020.
- 446 [45] Junchi Yu, Tingyang Xu, Yu Rong, Yatao Bian, Junzhou Huang, and Ran He. Graph Information  
447 Bottleneck for Subgraph Recognition. In *ICLR*, 2020.
- 448 [46] Long Zhao, Yuxiao Wang, Jiaping Zhao, Liangzhe Yuan, Jennifer J. Sun, Florian Schroff,  
449 Hartwig Adam, Xi Peng, Dimitris Metaxas, and Ting Liu. Learning View-Disentangled Human  
450 Pose Representation by Contrastive Cross-View Mutual Information Maximization. In *CVPR*,  
451 2021.
- 452 [47] Chen Zhu, Yu Cheng, Zhe Gan, Siqi Sun, Tom Goldstein, and Jingjing Liu. FreeLb: Enhanced  
453 adversarial training for natural language understanding. In *ICLR*, 2020.
- 454 [48] Yanqiao Zhu, Yichen Xu, Qiang Liu, and Shu Wu. An Empirical Study of Graph Contrastive  
455 Learning. *arXiv.org*, 2021.

## 456 A Summary of Datasets

457 In this work, we use ten datasets from TU Benchmark Datasets [23] to evaluate our proposed  
458 GraphCV under unsupervised setting, where five of them are biochemical datasets and the other five  
459 belong to social network datasets. Besides, we utilize the ogng-molhiv dataset from Open Graph  
460 Benchmark (OGB) [15] to further evaluate GraphCV under semi-supervised setting. The statistics of  
461 these datasets are provided in Table 3

Dataset	#Graphs	Avg #Nodes	Avg #Edges	#Class	Metric	Category
MUTAG	188	17.93	19.79	2	Accuracy	biochemical
PTC-MR	344	14.29	14.69	2	Accuracy	biochemical
PROTEINS	1,113	39.06	72.82	2	Accuracy	biochemical
NCII	4,110	29.87	32.30	2	Accuracy	biochemical
DD	1,178	284.32	715.66	2	Accuracy	biochemical
COLLAB	5,000	74.49	2457.78	3	Accuracy	social network
IMDB-B	1,000	19.77	96.53	2	Accuracy	social network
RDT-B	2,000	429.63	497.75	2	Accuracy	social network
IMDB-M	1,500	13.00	65.94	3	Accuracy	social network
RDT-M5K	4,999	508.52	594.87	5	Accuracy	social network
ogbg-molhiv	41,127	25.50	27.50	2	ROC-AUC	MoleculeNet

Table 3: Summary of datasets statistics.

462 All of the eleven datasets are public available, we attach their links as follow:

- 463 • TU datasets: <https://chrsmrrs.github.io/datasets/docs/datasets/>  
464 • ogbg-molhiv dataset: <https://ogb.stanford.edu/docs/graphprop/#ogbg-mol>

## 465 B Implementation Details

466 All experiments are conducted with the following settings:

- 467 • Operating System: Ubuntu 18.04.5 LTS  
468 • CPU: AMD(R) Ryzen 9 3900x  
469 • GPU: NVIDIA GeForce RTX 2080ti  
470 • Software: Python 3.8.5; Pytorch 1.10.1; PyTorch Geometric 2.0.4; PyGCL 0.1.2; Numpy  
471 1.20.1; scikit-learn 0.24.1.

472 We choose GIN [41] as the backbone graph encoder and the model is optimized through Adam  
473 optimizer. We follow [44, 42, 18] to employ a linear SVM classifier for downstream task-specific  
474 classification. The graph augmentation operations used in our work are same as [44], including  
475 node dropping, edge perturbation, attribute masking and subgraph sampling, all of them are bor-  
476 rowed from the implementation of [48]. There are two specific hyper-parameters in our model,  
477 namely  $\lambda_r$  and  $\lambda_a$ , the search space of them are  $\{0.0, 1.0, 3.0, 5.0, 10.0\}$  and  $\{0.0, 0.5, 1.0, 1.5, 2.0\}$ ,  
478 respectively. For other important hyper-parameters, we find the best value of pre-training epoch  
479 from  $\{20, 50, 100, 200, 300\}$ , learning rate from  $\{0.01, 0.005, 0.001, 0.0005, 0.0001\}$ , embedding  
480 dimension from  $\{32, 64, 128, 256, 512\}$ , number of GNN layers from  $\{2, 3, 4, 5\}$ , batch size from  
481  $\{32, 64, 128, 256, 512\}$  (except for ogbg-molhiv  $\{64, 128, 256, 512, 1024\}$ ). Besides, we fix the per-  
482 turbation bound  $\epsilon$ , ascent step size  $\alpha$  and ascent step  $T$  as 0.008, 0.008 and 5 during hyper-parameter  
483 fine-tuning.

## 484 C Proof

### 485 C.1 Proof of Theorem 1

486 We repeat Theorem 1 as follows.

487 **Theorem 3** Suppose  $f(\cdot)$  is a GNN encoder as powerful as 1-WL test. Let  $g_{pre}(\cdot)$  elicits only the  
488 augmentation information from  $\mathbf{z}$  meanwhile  $g_{com}(\cdot)$  extracts the essential factors of  $G$  from  $\mathbf{z}_1$  and  
489  $\mathbf{z}_2$ . Then we have:

$$I(t_1(G); \mathbf{z}_2^{com}, \mathbf{z}_2^{pre}) \geq I(\mathbf{z}_1^{pre}; \mathbf{z}_2^{pre}) \text{ where } G \in \mathcal{G} \text{ and } t_1(\cdot), t_2(\cdot) \in \mathcal{T}.$$

490 **Proof.** According to the assumption in Theorem 1, for any two graphs  $G, G' \in \mathcal{G}$ , if  $G \cong G'$  then we  
491 have  $\mathbf{z} = \mathbf{z}'$ , where  $\mathbf{z} = f(G)$  and  $\mathbf{z}' = f(G')$ .

492 Besides,  $\mathbf{z}^{pre} = g_{pre}(\mathbf{z})$  is specific to the predictive factors and  $\mathbf{z}^{com} = g_{com}(\mathbf{z})$  is particular to the  
493 non-predictive factors, which means  $\mathbf{z}^{pre}$  and  $\mathbf{z}^{com}$  are mutually excluded and  $\mathbf{z}^{pre} \sim G$ . So we  
494 have,

$$\begin{aligned} p(\mathbf{z}^{pre}, \mathbf{z}^{com}) &= p(\mathbf{z}^{pre})p(\mathbf{z}^{com}) \\ p(\mathbf{z}^{pre}, \mathbf{z}^{com} | t(G)) &= p(\mathbf{z}^{pre} | t(G))p(\mathbf{z}^{com} | t(G)). \end{aligned} \quad (12)$$

495 Then, we want to prove that given three random variables  $a, b$  and  $c$ , if they satisfy  $p(b, c) = p(b)p(c)$   
496 and  $p(b, c | a) = p(b | a)p(c | a)$ , we have  $I(a, b | c) = I(a, b)$ . According to the definition of  
497 mutual information, we have that,

$$\begin{aligned} I(a; b | c) &= \\ &= \sum_a \sum_b \sum_c p(a, b, c) \log \frac{p(a, b, c)p(c)}{p(a, c)p(b, c)} \\ &= \sum_a \sum_b \sum_c p(a)p(b, c | a) \log \frac{p(b, c | a)p(c)}{p(c | a)p(b)p(c)} \\ &= \sum_a \sum_b \sum_c p(a)p(b | a)p(c | a) \log \frac{p(b | a)p(c | a)}{p(c | a)p(b)} \\ &= \sum_a \sum_b p(a)p(b | a) \log \frac{p(b | a)}{p(b)} \\ &= \sum_a \sum_b p(a, b) \log \frac{p(b | a)}{p(b)} \\ &= I(a; b). \end{aligned} \quad (13)$$

498 After that, by applying the chain rule to  $I(t_1(G); \mathbf{z}_2^{pre}, \mathbf{z}_2^{com})$ , we have,

$$\begin{aligned}
I(t_1(G); \mathbf{z}_2^{pre}, \mathbf{z}_2^{com}) &= I(t_1(G); \mathbf{z}_2^{pre} | \mathbf{z}_2^{com}) + I(t_1(G); \mathbf{z}_2^{com}) \\
&\stackrel{(2)}{=} I(t_1(G); \mathbf{z}_2^{pre}) + I(t_1(G); \mathbf{z}_2^{com}) \\
&\stackrel{(a)}{\geq} I(t_1(G); \mathbf{z}_2^{pre}) \\
&\stackrel{(b)}{\geq} I(\mathbf{z}_1^{com}, \mathbf{z}_1^{pre}; \mathbf{z}_2^{pre}) \\
&\stackrel{(2)}{=} I(\mathbf{z}_1^{com}; \mathbf{z}_2^{pre}) + I(\mathbf{z}_1^{pre}; \mathbf{z}_2^{pre}) \\
&\stackrel{(a)}{\geq} I(\mathbf{z}_1^{pre}; \mathbf{z}_2^{pre}),
\end{aligned} \tag{14}$$

499 where  $\stackrel{(2)}{=}$  is derived from the conclusion we get in Equation 13,  $\stackrel{(a)}{\geq}$  is based on the non-negativity of  
500 mutual information, i.e.,  $I(\cdot) \geq 0$ , and  $\stackrel{(b)}{\geq}$  is because data processing inequality [6]. Finally, we reach  
501 to the lower bound of  $I(t_1(G); \mathbf{z}_2^{pre}, \mathbf{z}_2^{com})$  in Equation 13, thus we can maximize the consistency  
502 between the information we capture from the two augmentation graph views by minimizing  $\mathcal{L}_{pre}$ .

## 503 C.2 Proof of Theorem 2

504 We repeat Theorem 2 as follows.

505 **Theorem 4** Assume  $q$  is a Gaussian distribution,  $g_r$  is the parameterized reconstruction model which  
506 infer  $\mathbf{z}_w$  from  $(\mathbf{z}_w^{pre}, \mathbf{z}_w^{com})$ . Then we have:

$$H(\mathbf{z}_w | \mathbf{z}_w^{pre}, \mathbf{z}_w^{com}) \leq \|\mathbf{z}_w - g_r(\mathbf{z}_w^{pre} \odot \mathbf{z}_w^{com})\|_2^2 \text{ where } w = w' \text{ or } w \neq w'.$$

507 **Proof.** To reconstruct the entangled representation  $\mathbf{z}_w$  from its corresponding non-predictive repre-  
508 sentation  $\mathbf{z}_w^{pre}$  and the predictive representation of any augmentation view  $\mathbf{z}_{w'}^{com}$  ( $w$  and  $w'$  are not  
509 necessarily equal), we need to minimize the conditional entropy:

$$H(\mathbf{z}_w | \mathbf{z}_w^{pre}, \mathbf{z}_{w'}^{com}) = -\mathbb{E}_{p(\mathbf{z}_w, \mathbf{z}_w^{pre}, \mathbf{z}_{w'}^{com})} [\log p(\mathbf{z}_w | \mathbf{z}_w^{pre}, \mathbf{z}_{w'}^{com})] \tag{15}$$

510 Since the real distribution of  $p(\mathbf{z}_w | \mathbf{z}_w^{pre}, \mathbf{z}_{w'}^{com})$  is unknown and intractable, we hereby introduce a  
511 variational distribution  $q(\mathbf{z}_w | \mathbf{z}_w^{pre}, \mathbf{z}_{w'}^{com})$  to approximate it. Therefore, we have,

$$\begin{aligned}
\mathbb{E}_{p(\mathbf{z}_w, \mathbf{z}_w^{pre}, \mathbf{z}_{w'}^{com})} [\log p(\mathbf{z}_w | \mathbf{z}_w^{pre}, \mathbf{z}_{w'}^{com})] &= \\
&\mathbb{E}_{p(\mathbf{z}_w, \mathbf{z}_w^{pre}, \mathbf{z}_{w'}^{com})} [\log q(\mathbf{z}_w | \mathbf{z}_w^{pre}, \mathbf{z}_{w'}^{com})] \\
&+ D_{\text{KL}}(p(\mathbf{z}_w | \mathbf{z}_w^{pre}, \mathbf{z}_{w'}^{com}) \| q(\mathbf{z}_w | \mathbf{z}_w^{pre}, \mathbf{z}_{w'}^{com})).
\end{aligned} \tag{16}$$

512 Due to the non-negativity of KL-divergence between any two distributions, it is safe to say  
513  $-\mathbb{E}_{p(\mathbf{z}_w, \mathbf{z}_w^{pre}, \mathbf{z}_{w'}^{com})} [\log q(\mathbf{z}_w | \mathbf{z}_w^{pre}, \mathbf{z}_{w'}^{com})]$  is the upper bound of  $H(\mathbf{z}_w | \mathbf{z}_w^{pre}, \mathbf{z}_{w'}^{com})$ . Based  
514 on the assumption of Theorem 2, let  $q(\mathbf{z}_w | \mathbf{z}_w^{pre}, \mathbf{z}_{w'}^{com})$  being a Gaussian distribution  
515  $\mathcal{N}(\mathbf{z}_w | g_r(\mathbf{z}_w^{pre} \odot \mathbf{z}_{w'}^{com}), \sigma^2 \mathbf{I})$ , where  $g_r(\cdot)$  is the reconstruct network that predict  $\mathbf{z}_w$  from  
516  $(\mathbf{z}_w^{pre}, \mathbf{z}_{w'}^{com})$  and  $\sigma$  is the variance. Thus we have,

$$\begin{aligned}
H(\mathbf{z}_w | \mathbf{z}_w^{pre}, \mathbf{z}_{w'}^{com}) &\leq -\mathbb{E}_{p(\mathbf{z}_w, \mathbf{z}_w^{pre}, \mathbf{z}_{w'}^{com})} [\log q(\mathbf{z}_w | \mathbf{z}_w^{pre}, \mathbf{z}_{w'}^{com})] \\
&= -\mathbb{E}_{p(\mathbf{z}_w, \mathbf{z}_w^{pre}, \mathbf{z}_{w'}^{com})} \left[ \log \left( \frac{1}{\sqrt{2\pi}I\sigma} e^{-\frac{1}{2} \frac{(\mathbf{z}_w - g_r(\mathbf{z}_w^{pre} \odot \mathbf{z}_{w'}^{com}))^2}{(\sigma^2 \mathbf{I})}} \right) \right] \\
&= -\mathbb{E}_{p(\mathbf{z}_w, \mathbf{z}_w^{pre}, \mathbf{z}_{w'}^{com})} \left[ \log \left( \frac{1}{\sqrt{2\pi}I\sigma} \right) - \frac{(\mathbf{z}_w - g_r(\mathbf{z}_w^{pre} \odot \mathbf{z}_{w'}^{com}))^2}{2\sigma^2 \mathbf{I}} \right]
\end{aligned} \tag{17}$$

517 Hence, we get the upper bound of  $H(\mathbf{z}_w | \mathbf{z}_w^{pre}, \mathbf{z}_{w'}^{com})$  as Equation 17. To minimize the value of the  
518 unsolvable entropy, we can instead minimize the value of its upper bound and thereby derive the  
519 objective function as follow by neglecting the constant terms,

$$\min \mathbb{E}_{p(\mathbf{z}_w, \mathbf{z}_w^{pre}, \mathbf{z}_{w'}^{com})} \|\mathbf{z}_w - g_r(\mathbf{z}_w^{pre} \odot \mathbf{z}_{w'}^{com})\|_2^2. \tag{18}$$

520 Since we adopt two augmentation views and propose the cross-view reconstruction mechanism in our  
521 method, we can minimize the entropy by minimizing  $\mathcal{L}_{recon}$  and thus guarantee the disentanglement  
522 of  $\mathbf{z}^{pre}$  and  $\mathbf{z}^{com}$ .

523 **D Effects of Representation Disentanglement**

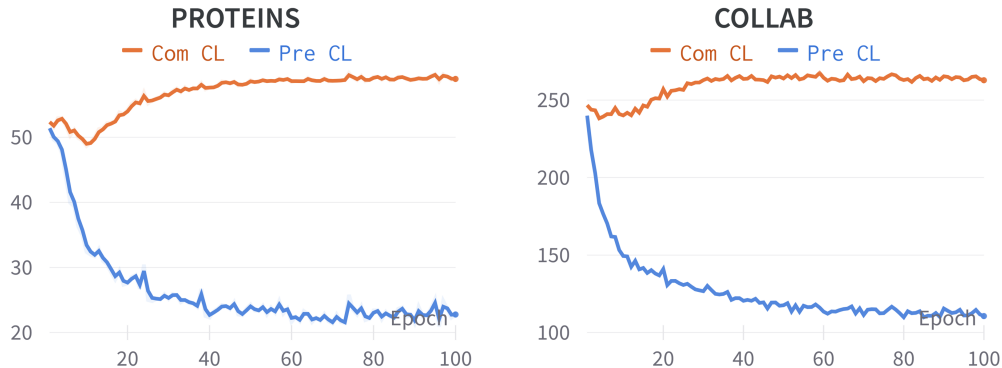


Figure 6: InfoNCE loss of the two disentangled representations between the two augmentation graph views, where orange lines are the InfoNCE loss between the two non-predictive representations and blue lines are the InfoNCE loss between the two predictive representations

524 In this section, we set experiments to further investigate the representation disentanglement of  
 525 our proposed GraphCV. Specifically, we use the InfoNCE loss [34] to dynamically measure the  
 526 representation difference between the two augmentation graph views based on the two disentangled  
 527 representations, where blue lines indicate the InfoNCE loss between  $\mathbf{z}_1^{pre}$  and  $\mathbf{z}_2^{pre}$  and orange lines  
 528 represent the InfoNCE loss between  $\mathbf{z}_1^{com}$  and  $\mathbf{z}_2^{com}$ . For simplicity, we only demonstrate the first 100  
 529 pre-training epochs of PROTEINS and COLLAB in Figure 6, we can observe similar phenomena on  
 530 other datasets. From the loss curves in Figure 6 we can find that contrastive loss between predictive  
 531 representations gradually decreases, indicating the predictive representation is optimized to capture  
 532 all the shared information between the two augmentation view. Meanwhile, we can see contrastive  
 533 loss between the non-predictive representations achieve a noticeable increases, which is consistent  
 534 with our expectation that the two independent sampled augmentation operators cause a distribution  
 535 shift between the two augmentation views.

536 **E Hyper-parameter Sensitivity**

537 In this section, we study the impacts of some important hyper-parameters in our method, including  
 538 reconstruction loss coefficient  $\lambda_r$ , adversarial loss coefficient  $\lambda_a$ , embedding dimension  $d$ , batch size  
 539  $|\mathcal{B}|$  and number of GNN layers  $L$ . Here, we select four datasets, i.e., MUTAG, PROTEINS, RDT-B  
 540 and COLLAB, to report for simplicity because the four datasets cover different domains and scales.  
 541 We illustrate the impacts of these hyper-parameters in the figures below.

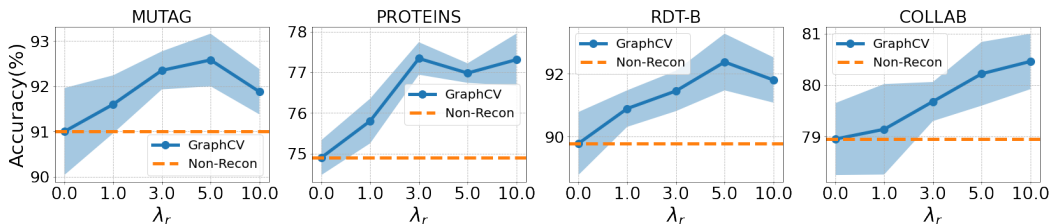


Figure 7: Impact of reconstruction loss coefficient  $\lambda_r$  on different datasets, we specify the non-reconstruction situation ( $\lambda_r = 0$ ) with the dashed line for comparison.

542 From the result demonstrated in Figure 7, we can see the optimal reconstruction loss coefficient  $\lambda_r$   
 543 is different dependent on the specific dataset, but all the values in our experiment can enhance the  
 544 performance compared with non-reconstruction variant, i.e.,  $\lambda_r = 0$ , indicating the effectiveness of  
 545 our proposed cross-view reconstruction mechanism.

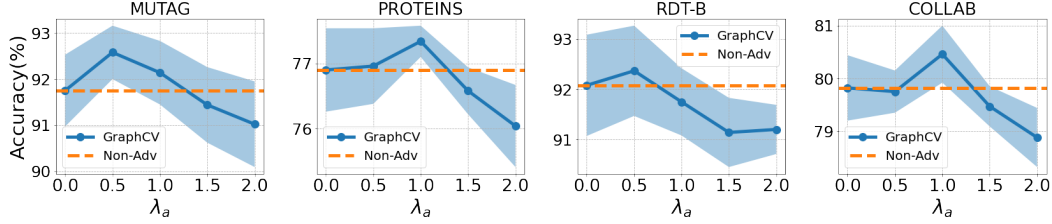


Figure 8: Impact of adversarial loss coefficient  $\lambda_a$  on different datasets, we specify the non-adversarial situation ( $\lambda_a = 0$ ) with the dashed line for comparison.

546 The Figure 8 shows that we could further raise the model performance through the adversarial  
 547 training, which proves a robust representation with less redundant information usually achieve more  
 548 performance gain compared with the brittle one. During this process, we need to choose a appropriate  
 549 adversarial loss coefficient  $\lambda_a$ , otherwise a too large  $\lambda_a$  may hurt the information sufficiency of the  
 550 learned representation.

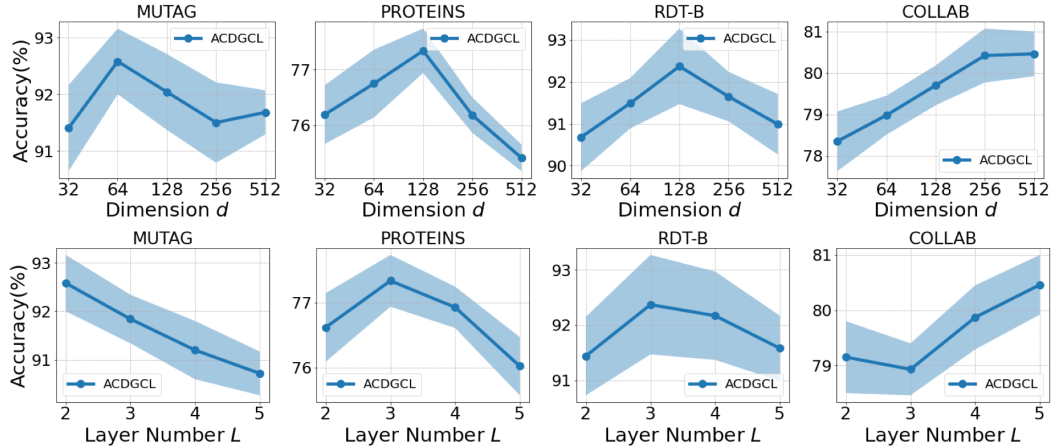


Figure 9: Impact of a embedding dimension  $d$  and GNN layer number  $L$  on different datasets.

551 We put the impacts of embedding dimension  $d$  and GNN layer number  $L$  together because we can  
 552 find a similar observation from their experimental results. From Figure 9, we observe that the optimal  
 553 values of the two hyper-parameters generally increase as the dataset scale increases. The reason  
 554 behind this phenomena could be large datasets usually contain more latent factors than the small  
 555 datasets, therefore a model with larger capacity is needed to fit the large datasets. However, such  
 556 high-capacity message-passing model will deteriorate the performance of small dataset because it  
 557 may cause the learned representation over-smoothing and hence less informative.

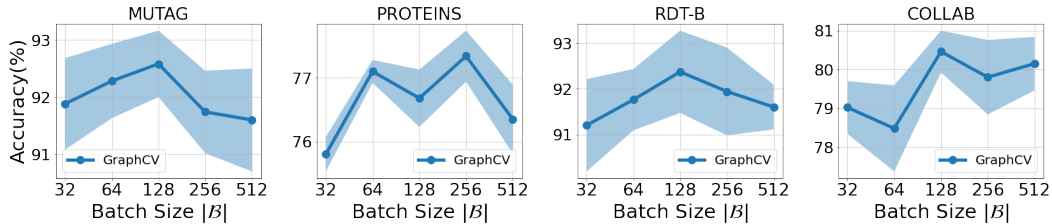


Figure 10: Impact of batch size  $|\mathcal{B}|$  on different datasets.

558 The effect of batch size  $|\mathcal{B}|$  is shown in Figure 10, we can see the performance variation of these  
 559 datasets under different batch size is relatively small. Most of the datasets achieve best performance  
 560 when the  $|\mathcal{B}|$  is set to 128.

561 **F Training Algorithm**

In this section we summarized the details of our proposed method in Algorithm 16

---

**Algorithm 1:** The training algorithm of **GraphCV**

---

**Input:** Graph dataset  $\mathcal{G} = \{G_i = (V_i, E_i)\}_{i=1}^N$ ; augmentation family  $\mathcal{T}$ ; loss coefficient  $\lambda_r, \lambda_a$ ; ascent step  $T$ ; ascent step size  $\alpha$ ; perturbation bound  $\epsilon$ .

**Output:** The disentangled predictive representations  $\mathbf{Z}^{pre} = \{\mathbf{z}_i^{pre}\}_{i=1}^N$

```

1 for each training epoch do
2   for sampled minibatch  $\mathcal{B} = \{G_i\}_{i=1}^{|\mathcal{B}|}$  do
3     for  $G_i \in \mathcal{B}$  do
4        $\mathbf{z}_{1,i} = f(t_1(G_i)), \mathbf{z}_{2,i} = f(t_2(G_i));$   $\triangleright t_1(\cdot), t_2(\cdot) \in \mathcal{T}$ 
5        $\mathbf{z}_{1,i}^{pre} = g_{pre}(\mathbf{z}_{1,i}), \mathbf{z}_{2,i}^{pre} = g_{pre}(\mathbf{z}_{2,i});$ 
6        $\mathbf{z}_{1,i}^{com} = g_{com}(\mathbf{z}_{1,i}), \mathbf{z}_{2,i}^{com} = g_{com}(\mathbf{z}_{2,i});$ 
7       Calculate  $\mathcal{L}_{pre}$  according to Equation 6;
8       Calculate  $\mathcal{L}_{recon}$  according to Equation 8;
9        $\mathcal{L} \leftarrow \mathcal{L}_{pre} + \lambda_r \mathcal{L}_{recon};$ 
10       $\delta_0 \leftarrow U(-\epsilon, \epsilon);$ 
11      for each  $t = 1$  to  $T$  do
12        Calculate the  $\mathcal{L}_{adv}$  according to Equation 10;
13         $\delta_t \leftarrow \delta_{t-1} + \alpha \nabla_{\delta} \mathcal{L}_{adv};$   $\triangleright$  Update perturbation to maximize  $\mathcal{L}_{adv}$ 
14         $\mathcal{L} \leftarrow \mathcal{L} + \frac{\lambda_a}{T} \mathcal{L}_{adv}$ 
15      Update the parameter  $\theta$  of  $f$  and  $g$  with the gradient  $\nabla_{\theta} \mathcal{L}(\theta, \mathcal{B})$  over a minibatch;
16 return  $\mathbf{Z}^{pre} = \{\mathbf{z}_i^{pre}\}_{i=1}^N$ , where  $\mathbf{z}_i^{pre} = g_{pre}(f(G_i))$ 

```

---

562