

---

# DiP-GNN: Discriminative Pre-Training of Graph Neural Networks

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Graph neural network (GNN) pre-training methods have been proposed to enhance  
2 the power of GNNs. Specifically, a GNN is first pre-trained on a large-scale unlabeled  
3 graph and then fine-tuned on a separate small labeled graph for downstream  
4 applications, such as node classification. One popular pre-training method is to  
5 mask out a proportion of the edges, and a GNN is trained to recover them. How-  
6 ever, such a generative method suffers from graph mismatch. That is, the masked  
7 graph input to the GNN deviates from the original graph. To alleviate this issue,  
8 we propose DiP-GNN (Discriminative Pre-training of Graph Neural Networks).  
9 Specifically, we train a generator to recover identities of the masked edges, and  
10 simultaneously, we train a discriminator to distinguish the generated edges from  
11 the original graph’s edges. The discriminator is subsequently used for downstream  
12 fine-tuning. In our pre-training framework, the graph seen by the discriminator  
13 better matches the original graph because the generator can recover a proportion  
14 of the masked edges. Extensive experiments on large-scale homogeneous and  
15 heterogeneous graphs demonstrate the effectiveness of DiP-GNN. Our code will  
16 be publicly available.

## 17 1 Introduction

18 Graph neural networks (GNNs) have achieved superior performance in various applications, such  
19 as node classification (Kipf and Welling, 2017), knowledge graph modeling (Schlichtkrull et al.,  
20 2018) and recommendation systems (Ying et al., 2018). To enhance the power of GNNs, generative  
21 pre-training methods are developed (Hu et al., 2020b). During the pre-training stage, a GNN  
22 incorporates topological information by training on a large-scale unlabeled graph in a self-supervised  
23 manner. Then, the pre-trained model is fine-tuned on a separate small labeled graph for downstream  
24 applications. Generative GNN pre-training is akin to masked language modeling in language model  
25 pre-training (Devlin et al., 2019). That is, for an input graph, we first randomly mask out a proportion  
26 of the edges, and then a GNN is trained to recover the original identity of the masked edges.

27 One major drawback with the abovementioned approach is *graph mismatch*. That is, the input graph  
28 to the GNN deviates from the original one since a considerable amount of edges are dropped. This  
29 causes changes in topological information, e.g., node connectivity. Consequently, the learned node  
30 embeddings may not be desirable.

31 To mitigate the above issues, we propose DiP-GNN ( **D**iscriminative **P**re-training of **G**raph **N**eural  
32 **N**etworks). In DiP-GNN, we simultaneously train a generator and a discriminator. The generator  
33 is trained similar to existing generative pre-training approaches, where the model seeks to recover  
34 the masked edges and outputs a reconstructed graph. Subsequently, the reconstructed graph is fed to  
35 the discriminator, which predicts whether each edge resides in the original graph (i.e., a true edge)  
36 or is wrongly constructed by the generator (i.e., a fake edge). After pre-training, we fine-tune the  
37 discriminator on downstream tasks. Figure 1 illustrates our training framework. Note that our work

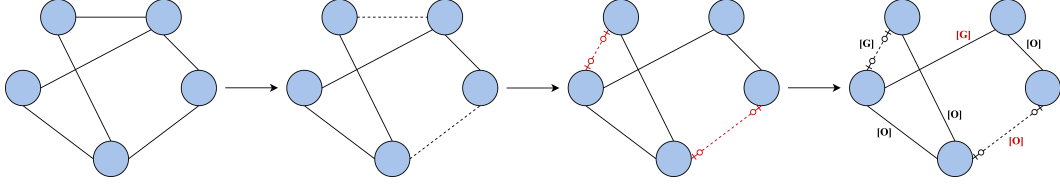


Figure 1: Illustration of DiP-GNN. From left to right: Original graph; Graph with two masked edges (dashed lines); Reconstructed graph created by the generator (generated edges are the dashed red lines); Discriminator labels each edge as  $[G]$  (generated) or  $[O]$  (original), where there are two wrong labels (shown in red).

38 is related to Generative Adversarial Nets (GAN, Goodfellow et al. 2014), and detailed discussions  
 39 are presented in Section 3.4. We remark that similar approaches have been used in natural language  
 40 processing (Clark et al., 2020). However, we identify the graph mismatch problem (see Section 4.5),  
 41 which is specific to graph-related applications and is not observed in natural language processing.

42 The proposed framework is more advantageous than generative pre-training. This is because the  
 43 reconstructed graph fed to the discriminator better matches the original graph compared with the  
 44 masked graph fed to the generator. Consequently, the discriminator can learn better node embeddings.  
 45 Such a better alignment is because the generator recovers the masked edges during pre-training, i.e.,  
 46 we observe that nearly 40% of the missing edges can be recovered. We remark that in our framework,  
 47 the graph fed to the generator has missing edges, while the graph fed to the discriminator contains  
 48 wrong edges since the generator may make erroneous predictions. However, empirically we find that  
 49 missing edges hurt more than wrong ones, making discriminative pre-training more desirable (see  
 50 Section 4.5 in the experiments).

51 We demonstrate effectiveness of DiP-GNN on large-scale homogeneous and heterogeneous graphs.  
 52 Results show that the proposed method significantly outperforms existing generative pre-training and  
 53 self-supervised learning approaches. For example, on the homogeneous Reddit dataset (Hamilton  
 54 et al., 2017) that contains 230k nodes, we obtain an improvement of 1.1 in terms of F1 score; and  
 55 on the heterogeneous OAG-CS graph (Tang et al., 2008) that contains 1.1M nodes, we obtain an  
 56 improvement of 2.8 in terms of MRR score in the paper field prediction task.

## 57 2 Background

58  $\diamond$  **Graph Neural Networks.** Graph neural networks compute a node’s representation by aggregating  
 59 information from the node’s neighbors. Concretely, for a multi-layer GNN, the feature vector  $h_v^{(k)}$  of  
 60 node  $v$  at the  $k$ -th layer is

$$h_v^{(k)} = \text{Combine} \left( a_v^{(k)}, h_v^{(k-1)} \right), \quad a_v^{(k)} = \text{Aggregate} \left( \left\{ h_u^{(k-1)} \mid \forall u \in \text{Neighbor}(v) \right\} \right),$$

61 where  $\text{Neighbor}(v)$  denotes all the neighbor nodes of  $v$ . Various implementations of  $\text{Aggregate}(\cdot)$   
 62 and  $\text{Combine}(\cdot)$  are proposed for both homogeneous (Defferrard et al., 2016; Kipf and Welling,  
 63 2017; Velickovic et al., 2018; Xu et al., 2019) and heterogeneous graphs (Schlichtkrull et al., 2018;  
 64 Wang et al., 2019; Zhang et al., 2019; Hu et al., 2020c).

65  $\diamond$  **Graph Neural Network Pre-Training.** Previous unsupervised learning methods leverage the  
 66 graph’s proximity (Tang et al., 2015) or information gathered by random walks (Perozzi et al., 2014;  
 67 Grover and Leskovec, 2016; Dong et al., 2017; Qiu et al., 2018). However, the learned embeddings  
 68 cannot be transferred to unseen nodes, limiting the methods’ applicability. Other unsupervised  
 69 learning algorithms adopt contrastive learning (Hassani and Ahmadi, 2020; Qiu et al., 2020; Zhu  
 70 et al., 2020, 2021; You et al., 2020, 2021). That is, we generate two views of the same graph, and  
 71 then maximize agreement of node presentations in the two views. However, our experiments reveal  
 72 that these methods do not scale well to extremely large graphs with millions of nodes.

73 Many GNN pre-training methods focus on generative objectives. For example, GAE (Graph Auto-  
 74 Encoder, Kipf and Welling 2016) proposes to reconstruct the graph structure; GraphSAGE (Hamilton  
 75 et al., 2017) optimizes an unsupervised loss derived from a random-walk-based metric; and DGI  
 76 (Deep Graph Infomax, Velickovic et al. 2019) maximizes the mutual information between node  
 77 representations and a graph summary representation.

78 There are also pre-training methods that extract graph-level representations, i.e., models are trained on  
79 a large amount of small graphs instead of a single large graph. For example, [Hu et al. 2020a](#) propose  
80 pre-training methods that operate on both graph and node level; and [InfoGraph \(Sun et al., 2020\)](#)  
81 proposes to maximize the mutual information between graph representations and representations  
82 of the graphs’ sub-structures. In this work, we focus on pre-training GNNs on a single large graph  
83 instead of multiple small graphs.

### 84 3 Method

85 We formally introduce the proposed discriminative GNN pre-training framework DiP-GNN. The  
86 algorithm contains two ingredients that operate on edges and features.

#### 87 3.1 Edge Generation and Discrimination

88 Suppose we have a graph  $\mathcal{G} = (\mathcal{N}, \mathcal{E})$ , where  $\mathcal{N}$  denotes all the nodes and  $\mathcal{E}$  denotes all the edges.  
89 We randomly mask out a proportion of the edges, such that  $\mathcal{E} = \mathcal{E}_u \cup \mathcal{E}_m$ , where  $\mathcal{E}_u$  is the unmasked  
90 set of edges and  $\mathcal{E}_m$  is the set of edges that are masked out.

91 For a masked edge  $e = (n_1, n_2) \in \mathcal{E}_m$ , where  $n_1$  and  $n_2$  are the two nodes connected by  $e$ , the  
92 generator’s goal is to predict  $n_1$  given  $n_2$  and the unmasked edges  $\mathcal{E}_u$ . For each node  $n$ , we compute  
93 its representation  $h_g(n) = f_g^e(n, \theta_g^e)$  using the generator  $f_g^e(\cdot, \theta_g^e)$ , which is parameterized by  $\theta_g^e$ . We  
94 remark that the computation of  $h_g(\cdot)$  only relies on the unmasked edges  $\mathcal{E}_u$ . We assume that the  
95 generation process of each edge is independent. Then, we have the prediction probability

$$p(n_1|n_2, \mathcal{E}_u) = \frac{\exp(d(h_g(n_1), h_g(n_2)))}{\sum_{n' \in \mathcal{C}} \exp(d(h_g(n'), h_g(n_2)))}, \text{ where } \mathcal{C} = \{n_1\} \cup (\mathcal{N} \setminus \text{Neighbor}(n_2)). \quad (1)$$

96 Here,  $\mathcal{C}$  is the candidate set for  $n_1$ , which contains all the nodes that are not connected to  $n_2$  except  
97  $n_1$  itself. Moreover, the distance function  $d(\cdot, \cdot)$  is chosen as a trainable cosine similarity, i.e.,

$$d(u, v) = \frac{(W^{\cos u})^\top v}{\|W^{\cos u}\| \cdot \|v\|}, \quad (2)$$

98 where  $W^{\cos}$  is a trainable weight. The training loss for the generator is defined as

$$\mathcal{L}_g^e(\theta_g^e) = \sum_{(n_1, n_2) \in \mathcal{E}_m} -\log p(n_1|n_2, \mathcal{E}_u), \quad (3)$$

99 which is equivalent to maximizing the likelihood of correct predictions.

100 The goal of the generator is to recover the masked edges in  $\mathcal{E}_m$ . Therefore, after we train the  
101 generator, we use the trained model to generate  $\mathcal{E}_g = \{(\hat{n}_1, n_2)\}_{(n_1, n_2) \in \mathcal{E}_m}$ , where each  $\hat{n}_1$  is the  
102 model’s prediction as  $\hat{n}_1 = \operatorname{argmax}_{n' \in \mathcal{C}} p(n'|n_2, \mathcal{E}_u)$ . Because the generator cannot correctly predict  
103 every edge, some edges in  $\mathcal{E}_g$  are wrongly generated (i.e., not in  $\mathcal{E}_m$ ). We refer to such edges as *fake*  
104 edges, and the rest as *true* edges. Concretely, we denote the true edges  $\mathcal{E}_{\text{true}} = \mathcal{E}_u \cup (\mathcal{E}_m \cap \mathcal{E}_g)$ , i.e.,  
105 the unmasked edges and the edges correctly generated by the generator. Correspondingly, we denote  
106 the fake edges  $\mathcal{E}_{\text{fake}} = \mathcal{E} \setminus \mathcal{E}_{\text{true}}$ .

107 The discriminator is trained to distinguish edges that are from the original graph (i.e., the true edges)  
108 and edges that are not (i.e., fake edges). Specifically, given the true edges  $\mathcal{E}_{\text{true}}$  and the fake ones  
109  $\mathcal{E}_{\text{fake}}$ , we first compute  $h_d(n) = f_d^e(n, \theta_d^e)$  for every node  $n \in \mathcal{N}$ , where  $f_d^e(\cdot, \theta_d^e)$  is the discriminator  
110 model parameterized by  $\theta_d^e$ . We highlight that different from computing  $h_g(\cdot)$ , the computation of  
111  $h_d(\cdot)$  relies on all the edges, such that the discriminator can separate a fake edge from a true one.  
112 Then, for each edge  $e = (n_1, n_2) \in \mathcal{E}_{\text{true}} \cup \mathcal{E}_{\text{fake}}$ , the discriminator outputs

$$p_{\text{fake}} = p(e \in \mathcal{E}_{\text{fake}} | \mathcal{E}_{\text{true}}, \mathcal{E}_{\text{fake}}) = \operatorname{sigmoid}(d(h_d(n_1), h_d(n_2))), \quad (4)$$

113 where  $d(\cdot, \cdot)$  is the distance function in Eq. 2. The training loss for the discriminator is the binary  
114 cross-entropy loss of predicting whether an edge is fake or not, defined as

$$\mathcal{L}_d^e(\theta_d^e) = \sum_{e \in \mathcal{E}_{\text{true}} \cup \mathcal{E}_{\text{fake}}} -\mathbf{1}\{e \in \mathcal{E}_{\text{fake}}\} \log(p_{\text{fake}}) - \mathbf{1}\{e \in \mathcal{E}_{\text{true}}\} \log(1 - p_{\text{fake}}), \quad (5)$$

115 where  $\mathbf{1}\{\cdot\}$  is the indicator function.

116 The edge loss is the weighted sum of the generator’s and the discriminator’s loss

$$\mathcal{L}^e(\theta_g^e, \theta_d^e) = \mathcal{L}_g^e(\theta_g^e) + \lambda \mathcal{L}_d^e(\theta_d^e), \quad (6)$$

117 where  $\lambda$  is a hyper-parameter. Note that structures of the generator  $f_g^e$  and the discriminator  $f_d^e$  are  
118 flexible, e.g., they can be graph convolutional networks (GCN) or graph attention networks (GAT).

### 119 3.2 Feature Generation and Discrimination

120 In real-world applications, nodes are often associated with features. For example, in the Reddit dataset  
 121 (Hamilton et al., 2017), a node’s feature is a vectorized representation of the post corresponding to  
 122 the node. As another example, in citation networks (Tang et al., 2008), a paper’s title can be treated  
 123 as a node’s feature. Previous work (Hu et al., 2020b) has demonstrated that generating features and  
 124 edges simultaneously can improve the GNN’s representation power.

125 Node features can be either texts (e.g., in citation networks) or vectors (e.g., in recommendation  
 126 systems). In this section, we develop feature generation and discrimination procedures for texts.  
 127 Vector features are akin to encoded text features, and we can use linear layers to generate and  
 128 discriminate them. Details about vector features are deferred to Appendix B.

129 For text features, we parameterize both the feature generator and discriminator using bi-directional  
 130 Transformer models (Vaswani et al., 2017), similar to BERT (Devlin et al., 2019). Denote  $f_g^f(\cdot, \theta_g^f) =$   
 131  $\text{trm}_g \circ \text{emb}_g(\cdot)$  the generator parameterized by  $\theta_g^f$ , where  $\text{emb}_g$  is the word embedding function and  
 132  $\text{trm}_g$  denotes subsequent Transformer layers. For an input text feature  $\mathbf{x} = [x_1, \dots, x_L]$  where  $L$  is  
 133 the sequence length, we randomly select indices to mask out, i.e., we randomly select an index set  
 134  $\mathcal{M} \subset \{1, \dots, L\}$ . For a masked position  $i \in \mathcal{M}$ , the prediction probability is given by

$$p(x_i|\mathbf{x}) = \frac{\exp(\text{emb}_g(x_i)^\top v_g(x_i))}{\sum_{x' \in \text{vocab}} \exp(\text{emb}_g(x')^\top v_g(x'))}, \quad v_g(x_i) = \text{trm}_g(W_g^{\text{proj}}[h_g(n_{\mathbf{x}}), \text{emb}_g(x_i)]).$$

135 Here  $W_g^{\text{proj}}$  is a trainable weight and  $h_g(n_{\mathbf{x}})$  is the representation of the node corresponding to  $\mathbf{x}$   
 136 computed by the edge generation GNN. Note that we concatenate the text embedding  $\text{emb}_g(x_i)$  and  
 137 the feature node’s embedding  $h_g(n_{\mathbf{x}})$ , such that the feature generator can aggregate information from  
 138 the graph structure. We train the generator by maximizing the probability of predicting the correct  
 139 token, i.e., by minimizing the loss

$$\mathcal{L}_g^f(\theta_g^e, \theta_g^f) = \sum_{\mathbf{x}} \sum_{i \in \mathcal{M}} -\log p(x_i|\mathbf{x}). \quad (7)$$

140 After we train the generator, we use the trained model to predict all the masked tokens, after which  
 141 we obtain a new text feature  $\mathbf{x}^{\text{corr}}$ . Here, we set  $x_i^{\text{corr}} = x_i$  for  $i \notin \mathcal{M}$  and  $x_i^{\text{corr}} = \hat{x}_i$  for  $i \in \mathcal{M}$ ,  
 142 where  $\hat{x}_i = \text{argmax}_{x' \in \text{vocab}} p(x_i|\mathbf{x})$  is the generator’s prediction.

143 The discriminator is trained to distinguish the fake tokens (i.e., wrongly generated tokens) from the  
 144 true ones (i.e., the unmasked and correctly generated tokens) in  $\mathbf{x}^{\text{corr}}$ . Similar to the generator, we  
 145 denote  $f_d^f(\cdot, \theta_d^f) = \text{trm}_d \circ \text{emb}_d(\cdot)$  as the discriminator parameterized by  $\theta_d^f$ . For each position  $i$ ,  
 146 the discriminator’s prediction probability is defined as

$$p(x_i^{\text{corr}} = x_i) = \text{sigmoid}(w^\top v_d(x_i^{\text{corr}})), \quad v_d(x_i^{\text{corr}}) = \text{trm}_d(W_d^{\text{proj}}[h_d(n_{\mathbf{x}}), \text{emb}_d(x_i^{\text{corr}})]).$$

147 Here  $w$  and  $W_d^{\text{proj}}$  are trainable weights and  $h_d(n_{\mathbf{x}})$  is the representation of the node corresponding  
 148 to  $\mathbf{x}$  computed by the edge discriminator GNN. The training loss for the discriminator is

$$\mathcal{L}_d^f(\theta_d^e, \theta_d^f) = \sum_{\mathbf{x}} \sum_{i=1}^L -\mathbf{1}\{x_i^{\text{corr}} = x_i\} \log(p_{\text{true}}) - \mathbf{1}\{x_i^{\text{corr}} \neq x_i\} \log(1 - p_{\text{true}}), \quad (8)$$

149 where  $p_{\text{true}} = p(x_i^{\text{corr}} = x_i)$ .

150 The text feature loss is defined as

$$\mathcal{L}^f(\theta_g^e, \theta_g^f, \theta_d^e, \theta_d^f) = \mathcal{L}_g^f(\theta_g^e, \theta_g^f) + \lambda \mathcal{L}_d^f(\theta_d^e, \theta_d^f), \quad (9)$$

151 where  $\lambda$  is a hyper-parameter.

### 152 3.3 Model Training

153 We jointly minimize the edge loss and the feature loss, where the loss function is

$$\begin{aligned} \mathcal{L}(\theta_g^e, \theta_g^f, \theta_d^e, \theta_d^f) &= \mathcal{L}^e(\theta_g^e, \theta_d^e) + \mathcal{L}^f(\theta_g^e, \theta_g^f, \theta_d^e, \theta_d^f) \\ &= (\mathcal{L}_g^e(\theta_g^e) + \mathcal{L}_g^f(\theta_g^e, \theta_g^f)) + \lambda (\mathcal{L}_d^e(\theta_d^e) + \mathcal{L}_d^f(\theta_d^e, \theta_d^f)). \end{aligned} \quad (10)$$

154 Here,  $\lambda$  is the weight of the discriminator’s loss. We remark that our framework is flexible because  
 155 the generator’s loss ( $\mathcal{L}_g^e$  and  $\mathcal{L}_g^f$ ) is decoupled from the discriminator’s ( $\mathcal{L}_d^e$  and  $\mathcal{L}_d^f$ ). As such,

156 existing generative pre-training methods can be applied to train the generator. In DiP-GNN, the  
157 discriminator has a better quality than the generator because of the graph mismatch issue (see  
158 Section 4.5). Therefore, after pre-training, we discard the generator and fine-tune the *discriminator*  
159 on downstream tasks. A detailed training pipeline is presented in Appendix A.

### 160 3.4 Comparison with GAN

161 We remark that our framework is different from Generative Adversarial Nets (GAN, Goodfellow et al.  
162 2014). In GAN, the generator-discriminator training framework is formulated as a min-max game,  
163 where the generator is trained adversarially to fool the discriminator. The two models are updated  
164 using alternating gradient descent/ascent.

165 However, the min-max game formulation of GAN is not applicable to our framework. This is because  
166 in GNN pre-training, the generator generates discrete edges, unlike continuous pixel values in the  
167 image domain. Such a property prohibits back-propagation from the discriminator to the generator.  
168 Existing works (Wang et al., 2018) use reinforcement learning (specifically policy gradient) to  
169 circumvent the non-differentiability issue. However, reinforcement learning introduces extensive  
170 hyper-parameter tuning and suffers from scalability issues. For example, the largest graph used in  
171 Wang et al. 2018 only contains 18k nodes, whereas the smallest graph used in our experiments has  
172 about 233k nodes.

173 Additionally, the goal of GAN is to train good-quality generators, which is different from our focus.  
174 In our discriminative pre-training framework, we focus on the discriminator because of better graph  
175 alignments. In practice, we find that accuracy of the generator is already high even without the  
176 discriminator, e.g., the accuracy is higher than 40% with 255 negative samples. And we observe that  
177 further improving the generator does not benefit downstream tasks.

## 178 4 Experiments

179 We implement all the algorithms using PyTorch (Paszke et al., 2019) and PyTorch Geometric (Fey  
180 and Lenssen, 2019). Experiments are conducted on NVIDIA A100 GPUs. By default, we use  
181 Heterogeneous Graph Transformer (HGT, Hu et al. 2020c) as the backbone GNN. We also discuss  
182 other choices in the experiments. Training and implementation details are deferred to Appendix C.

### 183 4.1 Settings and Datasets

184  $\diamond$  **Settings.** We consider a *node transfer* setting in the experiments. In practice we often work with  
185 a single large-scale graph, on which labels are sparse. In this case, we can use the large amount  
186 of unlabeled data as the pre-training dataset, and the rest are treated as labeled fine-tuning nodes.  
187 Correspondingly, edges between pre-training nodes are added to the pre-training data, and edges  
188 between fine-tuning nodes are added to the fine-tuning data. In this way, the model cannot see the  
189 fine-tuning data during pre-training, and vice versa.

190 We remark that our setting is different from conventional self-supervised learning, namely we pre-  
191 train and fine-tune on two separate graphs. This meets the practical need of transfer learning, e.g., a  
192 trained GNN needs to transfer across locales and time spans in recommendation systems.

193  $\diamond$  **Homogeneous Graph.** We use the Reddit dataset (Hamilton et al., 2017), which is a publicly  
194 available large-scale graph. In this graph, each node corresponds to a post, and is labeled with a  
195 “subreddit”. Each node has a 603-dimensional feature vector constructed from the corresponding  
196 post. Two nodes (posts) are connected if the same user commented on both. The dataset contains  
197 posts from 50 subreddits sampled from posts initiated in September 2014. In total, there are 232,965  
198 posts with an average node degree of 492. We use 70% of the data as the pre-training data, and the  
199 rest as the fine-tuning data, which are further split into training, validation, and test sets equally. We  
200 consider node classification as the downstream fine-tuning task.

201  $\diamond$  **Product Recommendation Graph.** We collect in-house product recommendation data from an  
202 e-commerce website. We build a bi-partite graph with two node types: search queries and product  
203 ids. The dataset contains about 633k query nodes, 2.71M product nodes, and 228M edges. We  
204 sample 70% of the nodes (and corresponding edges) for pre-training, and the rest are evenly split  
205 for fine-tuning training, validation and testing. We consider link prediction as the downstream task,  
206 where for each validation and test query node, we randomly mask out 20% of its edges to recover.

207 For each masked edge that corresponds to a query node and a positive product node, we randomly  
 208 sample 255 negative products. The task is to find the positive product out of the total 256 products.

209  $\diamond$  **Heterogeneous Graph.** We use the OAG-CS dataset (Tang et al., 2008; Sinha et al., 2015), which  
 210 is a publicly available heterogeneous graph containing computer science papers. The dataset contains  
 211 over 1.1M nodes and 28.4M edges. In this graph, there are five node types (institute, author, venue,  
 212 paper and field) and ten edge types. The “field” nodes are further categorized into six levels from L0  
 213 to L5, which are organized using a hierarchical tree. Details are shown in Figure 2.

214 We use papers published before 2014 as the  
 215 pre-training dataset (63%), papers published between 2014 (inclusive) and 2016 (inclusive) as  
 216 the fine-tuning training set (20%), papers published in 2017 as the fine-tuning validation set  
 217 (7%), and papers published after 2017 as the fine-tuning test set (10%). During fine-tuning,  
 218 by default we only use 10% of the fine-tuning training data (i.e., 2% of the overall data) be-  
 219 cause in practice labeled data are often scarce. We consider three tasks for fine-tuning: author  
 220 name disambiguation (AD), paper field classification (PF) and paper venue classification (PV). For paper field classification, we only consider L2  
 221 fields. In the experiments, we use the pre-processed graph from Hu et al. 2020b.

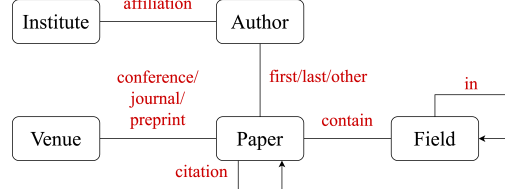


Figure 2: Details of OAG-CS. There are 5 node types (in black) and 10 edge types (in red).

## 228 4.2 Implementation Details

229  $\diamond$  **Graph subsampling.** In practice, graphs are often too large to fit in the hardware, e.g., the Reddit  
 230 dataset (Hamilton et al., 2017) contains over 230k nodes. Therefore, we sample a dense subgraph  
 231 from the large-scale graph in each training iteration. For homogeneous graphs, we apply the LADIES  
 232 algorithm (Zou et al., 2019), which theoretically guarantees that the sampled nodes are highly inter-  
 233 connected with each other and can maximally preserve the graph structure. For heterogeneous graphs,  
 234 we use the HGSampling algorithm (Hu et al., 2020b), which is a heterogeneous version of LADIES.

235  $\diamond$  **Node sampling for the edge generator.** In the edge generator, for a masked edge  $(s, t)$ , we fix the  
 236 node  $t$  and seek to identify the other node  $s$ . One approach is to identify  $s$  from all the graph nodes,  
 237 i.e., by setting  $\mathcal{C} = \mathcal{N}$  in Eq. 1. However, this task is computationally intractable when the number  
 238 of nodes is large, i.e., the model needs to find  $s$  out of hundreds of thousands of nodes. Therefore,  
 239 we sample some negative nodes  $\{s_i^g\}_{i=1}^{n_{\text{neg}}}$  such that  $(s_i^g, t) \notin \mathcal{E}$ . Then, the candidate set to generate  
 240 the source node becomes  $\{s, s_1^g, \dots, s_{n_{\text{neg}}}^g\}$  instead of all the graph nodes  $\mathcal{N}$ . We remark that such a  
 241 sampling approach is standard for GNN pre-training and link prediction (Hamilton et al., 2017; Sun  
 242 et al., 2020; Hu et al., 2020b).

243  $\diamond$  **Edge sampling for the edge discriminator.** In computing the loss for the discriminator, the number  
 244 of edges in  $\mathcal{E}_u$  is significantly larger than those in  $\mathcal{E}_g$ , i.e., we only mask a small proportion of the  
 245 edges. To avoid the discriminator from outputting trivial predictions (i.e., all the edges belong to  $\mathcal{E}_u$ ),  
 246 we balance the two loss terms in  $\mathcal{L}_d^e$ . Specifically, we sample  $\mathcal{E}_u^d \subset \mathcal{E}_u$  such that  $|\mathcal{E}_u^d| = \alpha|\mathcal{E}_g|$ , where  
 247  $\alpha$  is a hyper-parameter. Then, we compute  $\mathcal{L}_d^e$  on  $\mathcal{E}_g$  and  $\mathcal{E}_u^d$ . Note that the node representations  $h_d$   
 248 are still computed using all the generated and unmasked edges  $\mathcal{E}_g$  and  $\mathcal{E}_u$ .

## 249 4.3 Baselines

250 We compare our method with several baselines in the experiments. For fair comparison, all the  
 251 methods are trained for the same number of GPU hours.

252  $\diamond$  **GAE** (Graph Auto-Encoder, Kipf and Welling 2016) adopts an auto-encoder for unsupervised  
 253 learning on graphs. In GAE, node embeddings are learnt using a GNN, and we minimize the  
 254 discrepancy between the original and the reconstructed adjacency matrix.

255  $\diamond$  **GraphSAGE** (Hamilton et al., 2017) encourages embeddings of neighboring nodes to be similar.  
 256 For each node, the method learns a function that generates embeddings by sampling and aggregating  
 257 features from the node’s neighbors.

Table 1: Experimental results on homogeneous graphs. We report F1 averaged over 10 runs for the Reddit data and MRR over 10 runs for the product recommendation data. The best results are shown in **bold**.

	Reddit	Recomm.
w/o pre-train	87.3	46.3
GAE	88.5	56.7
GraphSAGE	88.0	53.0
DGI	87.7	53.3
GPT-GNN	89.6	58.6
GRACE	89.0	51.5
GraphCL	88.6	—
JOAOv2	89.1	—
DiP-GNN	<b>90.7</b>	<b>60.1</b>

Table 2: Experimental results on OAG-CS (heterogeneous). Left to right: paper-field, paper-venue, author-name-disambiguation. We report MRR over 10 runs. The best results are shown in **bold**.

	PF	PV	AD
w/o pre-train	32.7	19.6	60.0
GAE	40.3	24.5	62.5
GraphSAGE	37.8	22.1	62.9
DGI	38.1	22.5	63.0
GPT-GNN	41.6	25.6	63.1
GRACE	38.0	21.5	62.0
GraphCL	38.0	22.0	61.5
JOAOv2	38.6	23.5	62.8
DiP-GNN	<b>44.1</b>	<b>27.7</b>	<b>65.6</b>

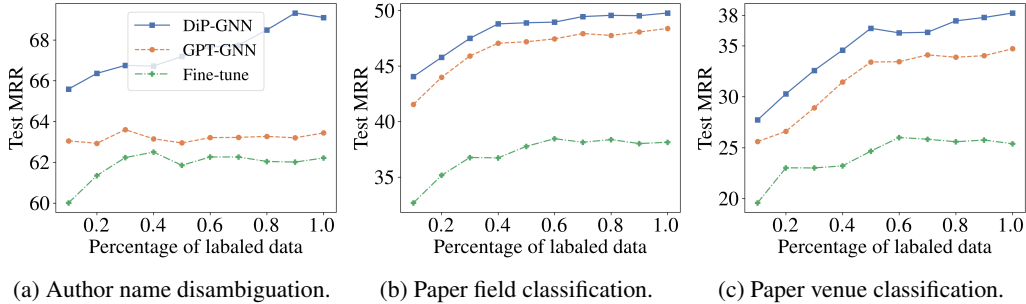


Figure 3: Model performance vs. amount of labeled data on OAG-CS.

258  $\diamond$  **DGI** (Deep Graph Infomax, [Velickovic et al. 2019](#)) maximizes information between node represen-  
 259 tations and corresponding high-level summaries of graphs. Thus, a node’s embedding summarizes a  
 260 sub-graph centered around it.

261  $\diamond$  **GPT-GNN** ([Hu et al., 2020b](#)) adopts a generative pre-training objective. The method generates  
 262 edges by minimizing a link prediction objective, and incorporates node features in the framework.

263  $\diamond$  **GRACE** (Graph Contrastive Representation, [Zhu et al. 2020](#)) leverages a contrastive objective.  
 264 The algorithm generates two views of the same graph through node and feature corruption, and then  
 265 maximize agreement of node representations in the two views.

266  $\diamond$  **GraphCL** ([You et al., 2020](#)) is another graph contrastive learning approach that adopts node and  
 267 edge augmentation techniques, such as node dropping and edge perturbation.

268  $\diamond$  **JOAO** (Joint Augmentation Optimization, [You et al. 2021](#)) improves GraphCL by deigning a  
 269 bi-level optimization objective to automatically and dynamically selects augmentation methods.

## 270 4.4 Main Results

271 In Table 1 and Table 2, *w/o pre-train* means direct training on the fine-tuning dataset without pre-  
 272 training. Results on the Reddit dataset are F1 scores averaged over 10 runs, and results on the product  
 273 recommendation graph are MRR scores averaged over 10 runs. All the performance gain have passed a  
 274 hypothesis test with p-value  $< 0.05$ .

275 Table 1 summarizes experimental results on the homogeneous graphs: Reddit and Recommendation.  
 276 We see that pre-training indeed benefits downstream tasks. For example, performance of GNN  
 277 improves by at  $\geq 0.4$  F1 on Reddit (DGI) and  $\geq 5.2$  MRR on Recommendation (GRACE). Also,  
 278 notice that among the baselines, generative approaches (GAE and GPT-GNN) yield promising  
 279 performance. On the other hand, the contrastive method (GRACE, GraphCL and JOAO) does not  
 280 scale well to large graphs, e.g., the OAG-CS graph which contains 1.1M nodes and 28.4M edges. By

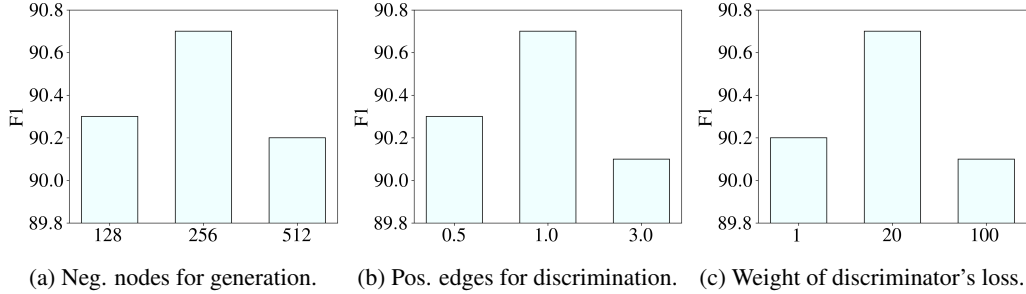


Figure 4: Ablation experiments on Reddit. By default, we set the number of negative nodes to 256, the factor of positive edges to 1.0, and weight of the discriminator's loss to 20.

Table 3: Test F1 score of model variants on Reddit.

Model	F1
Edges+Features	90.7
Edges	90.4
Features	90.2
RandomEdges	89.8

Table 4: Test F1 score of models with different backbone GNNs on Reddit.

Model	HGT	GAT
w/o pretrain	87.3	86.4
GPT-GNN	89.6	87.5
DiP-GNN	<b>90.7</b>	<b>88.5</b>

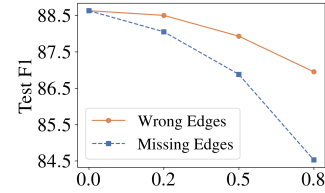


Figure 5: F1 vs. proportion of manipulated edges on Reddit.

281 using the proposed discriminative pre-training framework, our method significantly outperforms all  
 282 the baseline approaches. For example, DiP-GNN outperforms GPT-GNN by 1.1 on Reddit and 1.5  
 283 on Recommendation.

284 Experimental results on the heterogeneous OAG-CS dataset are summarized in Table 2. Similar to  
 285 the homogeneous graphs, notice that pre-training improves model performance by large margins.  
 286 For example, pre-training improves MRR by at least 5.1, 2.5 and 2.5 on the PF, PV and AD tasks,  
 287 respectively. Moreover, by using the proposed training framework, models can learn better node  
 288 embeddings and yield consistently better performance compared with all the baselines.

289 Recall that during fine-tuning on OAG-CS, we only use 10% of the labeled fine-tuning data (about  
 290 2% of the overall data). In Figure 3, we examine the effect of the amount of labeled data. We see  
 291 that model performance improves when we increase the amount of labeled data. Also, notice that  
 292 DiP-GNN consistently outperforms GPT-GNN in all the three tasks under all the settings.

#### 293 4.5 Analysis

294  $\diamond$  **Comparison with semi-supervised learning.** We compare  
 295 DiP-GNN with a semi-supervised learning method: C&S (Cor-  
 296 rect&Smooth, Huang et al. 2020). Figure 7 summarizes the re-  
 297 sults. We see that C&S yields a 0.5 improvement compared with  
 298 the supervised learning method (i.e., w/o pre-train). However,  
 299 performance of C&S is significantly lower than both DiP-GNN  
 300 and other pre-training methods such as GPT-GNN.

301  $\diamond$  **Hyper-parameters.** There are several hyper-parameters that  
 302 we introduce in DiP-GNN: the number of negative nodes that  
 303 are sampled for generating edges (Section 4.2); the number  
 304 of positive edges that are sampled for the discriminator's task  
 305 (Section 4.2); and the weight of the discriminator's loss (Eq. 10).  
 306 Figure 4 illustrate ablation experimental results on the Reddit  
 307 dataset. From the results, we see that DiP-GNN is robust to  
 308 these hyper-parameters. We remark that under all the settings,  
 309 ours model behaves better than the best-performing baseline (89.6 for GPT-GNN).

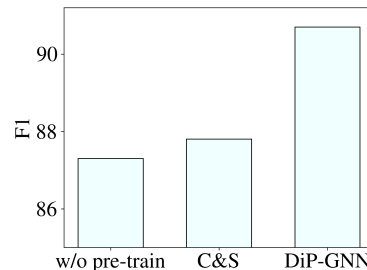


Figure 7: Comparison with semi-supervised learning methods. We report test F1 score on Reddit.

310  $\diamond$  **Model variants.** We also examine variants of DiP-GNN. Recall that the generator and the  
 311 discriminator operate on both edges and node features. We first check the contribution of these two



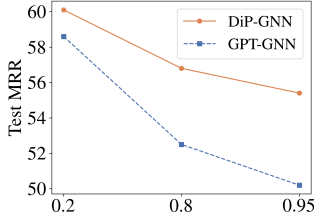


Figure 6: Performance vs. proportion of masked edges on product recommendation.

Table 5: Generator and discriminator performance vs. proportion of masked edges during pre-training. *Coverage* is the proportion of true edges input to the models.

Masked%	Acc		Coverage		
	Gen.	Dis.	Gen.	Dis.	Ratio
20	0.50	0.87	0.80	0.90	×1.13
80	0.33	0.84	0.20	0.46	×2.30
95	0.20	0.80	0.05	0.24	×4.80

312 factors. We also investigate the scenario where edges are randomly generated, and the discriminator  
 313 still seeks to find the generated edges. Table 3 summarizes results on the Reddit dataset.

314 We see that by only using edges, model performance drops by 0.3; and by only using node features,  
 315 performance drops by 0.5. This indicates that the graph structure plays a more important role  
 316 in the proposed framework than the features. Also notice that performance of *RandomEdges* is  
 317 unsatisfactory. This is because implausible edges are generated when using a random generator,  
 318 making the discriminator’s task significantly easier. We remark that performance of all the model  
 319 variants is better than the best-performing baseline, which is 89.6 for GPT-GNN.

320 Table 4 examines performance of our method and GPT-GNN using different backbone GNNs. Recall  
 321 that by default, we use HGT (Hu et al., 2020c) as the backbone. We see that when GAT (Velickovic  
 322 et al., 2018) is used, performance of DiP-GNN is still significantly better than GPT-GNN.

323 **◇ Missing edges hurt more than wrong edges.** In our pre-training framework, the generator is  
 324 trained to reconstruct the masked graph, after which the reconstructed graph is fed to the discriminator.  
 325 During this procedure, the graph input to the generator has *missing edges*, and the graph input to  
 326 the discriminator has *wrong edges*. From Figure 5, we see that wrong edges hurt less than missing  
 327 ones. For example, model performance drops by 0.7% when 50% of wrong edges are added to the  
 328 original graph, and performance decreases by 1.8% when 50% of original edges are missing. This  
 329 indicates that performance relies on the amount of original edges seen by the models. Intuitively,  
 330 wrong edges add noise to the graph, but they do not affect information flow. On the contrary, missing  
 331 edges cut information flow. Moreover, in practice we work with graph attention models, and the  
 332 attention mechanism can alleviate the wrong edges by assigning low attention scores to them.

333 **◇ Why is discriminative pre-training better?** Figure 6 illustrates effect of the proportion of masked  
 334 edges during pre-training. We see that when we increase the proportion from 0.2 to 0.8, performance  
 335 of GPT-GNN drops by 6.1, whereas performance of DiP-GNN only drops by 3.3. This indicates that  
 336 the generative pre-training method is more sensitive to the masking proportion.

337 Table 5 summarizes pre-training quality. First, the generative task (i.e., the generator) is more difficult  
 338 than the discriminative task (i.e., the discriminator). For example, when we increase the proportion  
 339 of masked edges from 20% to 80%, accuracy of the generator drops by 17% while accuracy of the  
 340 discriminator only decreases by 3%. Second, the graph input to the discriminator better aligns with  
 341 the original graph. For example, when 80% of the edges are masked, the discriminator sees 2.3 times  
 342 more original edges than the generator. Therefore, the discriminative task is more advantageous  
 343 because model quality relies on the number of observed original edges (Figure 5).

## 344 5 Conclusion and Discussions

345 We propose Discriminative Pre-Training of Graph Neural Networks (DiP-GNN), where we simulta-  
 346 neously train a generator and a discriminator. During pre-training, we mask out some edges in the  
 347 graph, and a generator is trained to recover the masked edges. Subsequently, a discriminator seeks to  
 348 distinguish the generated edges from the original ones. We conduct extensive experiments to validate  
 349 the effectiveness of DiP-GNN.

350 In this work, we focus on node-level tasks, such as node classification. The proposed framework is  
 351 generic and can be extended to graph-level applications, e.g., graph classification. The authors do not  
 352 find any immediate negative societal impact.

353 **References**

- 354 CLARK, K., LUONG, M., LE, Q. V. and MANNING, C. D. (2020). ELECTRA: pre-training text  
355 encoders as discriminators rather than generators. In *8th International Conference on Learning*  
356 *Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net.
- 357 DEFFERRARD, M., BRESSON, X. and VANDERGHEYNST, P. (2016). Convolutional neural networks  
358 on graphs with fast localized spectral filtering. In *Advances in Neural Information Processing*  
359 *Systems 29: Annual Conference on Neural Information Processing Systems 2016, December 5-10,*  
360 *2016, Barcelona, Spain* (D. D. Lee, M. Sugiyama, U. von Luxburg, I. Guyon and R. Garnett, eds.).
- 361 DEVLIN, J., CHANG, M.-W., LEE, K. and TOUTANOVA, K. (2019). BERT: Pre-training of deep  
362 bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of*  
363 *the North American Chapter of the Association for Computational Linguistics: Human Language*  
364 *Technologies, Volume 1 (Long and Short Papers)*. Association for Computational Linguistics,  
365 Minneapolis, Minnesota.
- 366 DONG, Y., CHAWLA, N. V. and SWAMI, A. (2017). metapath2vec: Scalable representation learning  
367 for heterogeneous networks. In *Proceedings of the 23rd ACM SIGKDD International Conference*  
368 *on Knowledge Discovery and Data Mining, Halifax, NS, Canada, August 13 - 17, 2017*. ACM.
- 369 FEY, M. and LENSSEN, J. E. (2019). Fast graph representation learning with PyTorch Geometric. In  
370 *ICLR Workshop on Representation Learning on Graphs and Manifolds*.
- 371 GOODFELLOW, I. J., POUGET-ABADIE, J., MIRZA, M., XU, B., WARDE-FARLEY, D., OZAIR,  
372 S., COURVILLE, A. C. and BENGIO, Y. (2014). Generative adversarial nets. In *Advances in*  
373 *Neural Information Processing Systems 27: Annual Conference on Neural Information Processing*  
374 *Systems 2014, December 8-13 2014, Montreal, Quebec, Canada* (Z. Ghahramani, M. Welling,  
375 C. Cortes, N. D. Lawrence and K. Q. Weinberger, eds.).
- 376 GROVER, A. and LESKOVEC, J. (2016). node2vec: Scalable feature learning for networks. In  
377 *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and*  
378 *Data Mining, San Francisco, CA, USA, August 13-17, 2016* (B. Krishnapuram, M. Shah, A. J.  
379 Smola, C. C. Aggarwal, D. Shen and R. Rastogi, eds.). ACM.
- 380 HAMILTON, W. L., YING, Z. and LESKOVEC, J. (2017). Inductive representation learning on large  
381 graphs. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural*  
382 *Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA* (I. Guyon,  
383 U. von Luxburg, S. Bengio, H. M. Wallach, R. Fergus, S. V. N. Vishwanathan and R. Garnett, eds.).
- 384 HASSANI, K. and AHMADI, A. H. K. (2020). Contrastive multi-view representation learning on  
385 graphs. In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020,*  
386 *13-18 July 2020, Virtual Event*, vol. 119 of *Proceedings of Machine Learning Research*. PMLR.
- 387 HU, W., LIU, B., GOMES, J., ZITNIK, M., LIANG, P., PANDE, V. S. and LESKOVEC, J. (2020a).  
388 Strategies for pre-training graph neural networks. In *8th International Conference on Learning*  
389 *Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net.
- 390 HU, Z., DONG, Y., WANG, K., CHANG, K. and SUN, Y. (2020b). GPT-GNN: generative pre-  
391 training of graph neural networks. In *KDD '20: The 26th ACM SIGKDD Conference on Knowledge*  
392 *Discovery and Data Mining, Virtual Event, CA, USA, August 23-27, 2020* (R. Gupta, Y. Liu, J. Tang  
393 and B. A. Prakash, eds.). ACM.
- 394 HU, Z., DONG, Y., WANG, K. and SUN, Y. (2020c). Heterogeneous graph transformer. In *WWW*  
395 *'20: The Web Conference 2020, Taipei, Taiwan, April 20-24, 2020* (Y. Huang, I. King, T. Liu and  
396 M. van Steen, eds.). ACM / IW3C2.
- 397 HUANG, Q., HE, H., SINGH, A., LIM, S.-N. and BENSON, A. R. (2020). Combining label propa-  
398 gation and simple models out-performs graph neural networks. *arXiv preprint arXiv:2010.13993*.
- 399 KIPF, T. N. and WELLING, M. (2016). Variational graph auto-encoders. *arXiv preprint*  
400 *arXiv:1611.07308*.

- 401 KIPF, T. N. and WELLING, M. (2017). Semi-supervised classification with graph convolutional  
402 networks. In *5th International Conference on Learning Representations, ICLR 2017, Toulon,*  
403 *France, April 24-26, 2017, Conference Track Proceedings*. OpenReview.net.
- 404 LOSHCHILOV, I. and HUTTER, F. (2019). Decoupled weight decay regularization. In *7th Interna-*  
405 *tional Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9,*  
406 *2019*. OpenReview.net.
- 407 PASZKE, A., GROSS, S., MASSA, F., LERER, A., BRADBURY, J., CHANAN, G., KILLEEN, T.,  
408 LIN, Z., GIMELSHEIN, N., ANTIGA, L., DESMAISON, A., KÖPF, A., YANG, E., DEVITO, Z.,  
409 RAISON, M., TEJANI, A., CHILAMKURTHY, S., STEINER, B., FANG, L., BAI, J. and CHINTALA,  
410 S. (2019). Pytorch: An imperative style, high-performance deep learning library. In *Advances in*  
411 *Neural Information Processing Systems 32: Annual Conference on Neural Information Processing*  
412 *Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada* (H. M. Wallach,  
413 H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. B. Fox and R. Garnett, eds.).
- 414 PEROZZI, B., AL-RFOU, R. and SKIENA, S. (2014). Deepwalk: online learning of social repre-  
415 sentations. In *The 20th ACM SIGKDD International Conference on Knowledge Discovery and*  
416 *Data Mining, KDD ’14, New York, NY, USA - August 24 - 27, 2014* (S. A. Macskassy, C. Perlich,  
417 J. Leskovec, W. Wang and R. Ghani, eds.). ACM.
- 418 QIU, J., CHEN, Q., DONG, Y., ZHANG, J., YANG, H., DING, M., WANG, K. and TANG, J. (2020).  
419 GCC: graph contrastive coding for graph neural network pre-training. In *KDD ’20: The 26th ACM*  
420 *SIGKDD Conference on Knowledge Discovery and Data Mining, Virtual Event, CA, USA, August*  
421 *23-27, 2020* (R. Gupta, Y. Liu, J. Tang and B. A. Prakash, eds.). ACM.
- 422 QIU, J., DONG, Y., MA, H., LI, J., WANG, K. and TANG, J. (2018). Network embedding as matrix  
423 factorization: Unifying deepwalk, line, pte, and node2vec. In *Proceedings of the Eleventh ACM*  
424 *International Conference on Web Search and Data Mining, WSDM 2018, Marina Del Rey, CA,*  
425 *USA, February 5-9, 2018* (Y. Chang, C. Zhai, Y. Liu and Y. Maarek, eds.). ACM.
- 426 SCHLICHTKRULL, M., KIPF, T. N., BLOEM, P., BERG, R. V. D., TITOV, I. and WELLING, M.  
427 (2018). Modeling relational data with graph convolutional networks. In *European semantic web*  
428 *conference*. Springer.
- 429 SINHA, A., SHEN, Z., SONG, Y., MA, H., EIDE, D., HSU, B.-J. and WANG, K. (2015). An  
430 overview of microsoft academic service (mas) and applications. In *Proceedings of the 24th*  
431 *international conference on world wide web*.
- 432 SUN, F., HOFFMANN, J., VERMA, V. and TANG, J. (2020). Infograph: Unsupervised and semi-  
433 supervised graph-level representation learning via mutual information maximization. In *8th*  
434 *International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April*  
435 *26-30, 2020*. OpenReview.net.
- 436 TANG, J., QU, M., WANG, M., ZHANG, M., YAN, J. and MEI, Q. (2015). LINE: large-scale infor-  
437 mation network embedding. In *Proceedings of the 24th International Conference on World Wide*  
438 *Web, WWW 2015, Florence, Italy, May 18-22, 2015* (A. Gangemi, S. Leonardi and A. Panconesi,  
439 eds.). ACM.
- 440 TANG, J., ZHANG, J., YAO, L., LI, J., ZHANG, L. and SU, Z. (2008). Arnetminer: extraction and  
441 mining of academic social networks. In *Proceedings of the 14th ACM SIGKDD international*  
442 *conference on Knowledge discovery and data mining*.
- 443 VASWANI, A., SHAZEER, N., PARMAR, N., USZKOREIT, J., JONES, L., GOMEZ, A. N., KAISER,  
444 L. and POLOSUKHIN, I. (2017). Attention is all you need. In *Advances in Neural Information*  
445 *Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017,*  
446 *December 4-9, 2017, Long Beach, CA, USA* (I. Guyon, U. von Luxburg, S. Bengio, H. M. Wallach,  
447 R. Fergus, S. V. N. Vishwanathan and R. Garnett, eds.).
- 448 VELICKOVIC, P., CUCURULL, G., CASANOVA, A., ROMERO, A., LIÒ, P. and BENGIO, Y. (2018).  
449 Graph attention networks. In *6th International Conference on Learning Representations, ICLR*  
450 *2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenRe-  
451 view.net.

- 452 VELICKOVIC, P., FEDUS, W., HAMILTON, W. L., LIÒ, P., BENGIO, Y. and HJELM, R. D. (2019).  
453 Deep graph infomax. In *7th International Conference on Learning Representations, ICLR 2019,*  
454 *New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net.
- 455 WANG, H., WANG, J., WANG, J., ZHAO, M., ZHANG, W., ZHANG, F., XIE, X. and GUO, M.  
456 (2018). Graphgan: Graph representation learning with generative adversarial nets. In *Proceedings*  
457 *of the Thirty-Second AAAI Conference on Artificial Intelligence, (AAAI-18), the 30th innovative*  
458 *Applications of Artificial Intelligence (IAAI-18), and the 8th AAAI Symposium on Educational*  
459 *Advances in Artificial Intelligence (EAAI-18), New Orleans, Louisiana, USA, February 2-7, 2018*  
460 (S. A. McIlraith and K. Q. Weinberger, eds.). AAAI Press.
- 461 WANG, X., JI, H., SHI, C., WANG, B., YE, Y., CUI, P. and YU, P. S. (2019). Heterogeneous graph  
462 attention network. In *The World Wide Web Conference, WWW 2019, San Francisco, CA, USA,*  
463 *May 13-17, 2019* (L. Liu, R. W. White, A. Mantrach, F. Silvestri, J. J. McAuley, R. Baeza-Yates  
464 and L. Zia, eds.). ACM.
- 465 XU, K., HU, W., LESKOVEC, J. and JEGELKA, S. (2019). How powerful are graph neural networks?  
466 In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA,*  
467 *May 6-9, 2019*. OpenReview.net.
- 468 YING, R., HE, R., CHEN, K., EKSOMBATCHAI, P., HAMILTON, W. L. and LESKOVEC, J. (2018).  
469 Graph convolutional neural networks for web-scale recommender systems. In *Proceedings of the*  
470 *24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, KDD*  
471 *2018, London, UK, August 19-23, 2018* (Y. Guo and F. Farooq, eds.). ACM.
- 472 YOU, Y., CHEN, T., SHEN, Y. and WANG, Z. (2021). Graph contrastive learning automated. In  
473 *International Conference on Machine Learning*. PMLR.
- 474 YOU, Y., CHEN, T., SUI, Y., CHEN, T., WANG, Z. and SHEN, Y. (2020). Graph contrastive learning  
475 with augmentations. *Advances in Neural Information Processing Systems*, **33** 5812–5823.
- 476 ZHANG, C., SONG, D., HUANG, C., SWAMI, A. and CHAWLA, N. V. (2019). Heterogeneous graph  
477 neural network. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge*  
478 *Discovery & Data Mining, KDD 2019, Anchorage, AK, USA, August 4-8, 2019* (A. Teredesai,  
479 V. Kumar, Y. Li, R. Rosales, E. Terzi and G. Karypis, eds.). ACM.
- 480 ZHU, Y., XU, Y., YU, F., LIU, Q., WU, S. and WANG, L. (2020). Deep graph contrastive  
481 representation learning. *arXiv preprint arXiv:2006.04131*.
- 482 ZHU, Y., XU, Y., YU, F., LIU, Q., WU, S. and WANG, L. (2021). Graph contrastive learning with  
483 adaptive augmentation. In *Proceedings of the Web Conference 2021*.
- 484 ZOU, D., HU, Z., WANG, Y., JIANG, S., SUN, Y. and GU, Q. (2019). Layer-dependent importance  
485 sampling for training deep and large graph convolutional networks. In *Advances in Neural*  
486 *Information Processing Systems 32: Annual Conference on Neural Information Processing Systems*  
487 *2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada* (H. M. Wallach, H. Larochelle,  
488 A. Beygelzimer, F. d’Alché-Buc, E. B. Fox and R. Garnett, eds.).

## 489 A Detailed Algorithm

490 Algorithm 1 is a detailed training pipeline of DiP-GNN. For graphs with vector features instead of  
 491 text features, we can substitute the feature generation and discrimination modules with equations in  
 492 Appendix B.

---

**Algorithm 1:** DiP-GNN: Discriminative Pre-training of Graph Neural Networks.

---

**Input:** Graph  $\mathcal{G}_{\text{full}}$ ; edge masking ratio; feature masking ratio; number of negative samples for edge generator; proportion of positive samples for edge discriminator  $\alpha$ ; weight of the discriminator’s loss  $\lambda$ ; number of training steps  $T$ .

```

for  $t = 0, \dots, T - 1$  do
  // Graph subsampling.
  Sample a subgraph  $\mathcal{G} = (\mathcal{N}, \mathcal{E})$  from  $\mathcal{G}_{\text{full}}$ ;
  // Edge generation.
  Initialize the generated edge set  $\mathcal{E}_g = \{\}$  and the edge generation loss  $\mathcal{L}_g^e = 0$ ;
  Construct the unmasked set of edges  $\mathcal{E}_u$  and the masked set  $\mathcal{E}_m$  such that  $\mathcal{E} = \mathcal{E}_u \cup \mathcal{E}_m$ ;
  Compute node embeddings using  $\mathcal{E}_u$ ;
  for  $e = (n_1, n_2) \in \mathcal{E}_m$  do
    Construct candidate set  $\mathcal{C}$  for  $n_1$  ( $n_2$  is given during generation) via negative sampling;
    Generate  $\hat{e} = (\hat{n}_1, n_2)$  where  $\hat{n}_1 \in \mathcal{C}$ ;
    Update the generated edge set  $\mathcal{E}_g \leftarrow \mathcal{E}_g \cup \{\hat{e}\}$ ;
    Update the edge generation loss  $\mathcal{L}_g^e$ ;
  // Text Feature generation.
  Initialize the feature generation loss  $\mathcal{L}_g^f = 0$ ;
  for  $n \in \mathcal{N}$  do
    For the node’s text feature  $\mathbf{x}_n$ , mask out some of its tokens;
    Construct the generated text feature  $\mathbf{x}_n^{\text{corr}}$  using the embedding of node  $n$  (computed during edge generation) and the feature generation Transformer model;
    Update the feature generation loss  $\mathcal{L}_g^f$ ;
  // Edge discrimination.
  Initialize the edge discrimination loss  $\mathcal{L}_d^e = 0$ ;
  Compute node embeddings using  $\mathcal{E}_g \cup \mathcal{E}_u$ ;
  Sample  $\mathcal{E}_u^d \subset \mathcal{E}_u$  such that  $|\mathcal{E}_u^d| = \alpha|\mathcal{E}_g|$ ;
  for  $e = (n_1, n_2) \in \mathcal{E}_g \cup \mathcal{E}_u^d$  do
    Determine if  $e$  is generated using the embedding of  $n_1$  and  $n_2$ ;
    Update the edge discrimination loss  $\mathcal{L}_d^e$ ;
  // Text feature discrimination.
  Initialize the feature discrimination loss  $\mathcal{L}_d^f = 0$ ;
  for  $n \in \mathcal{N}$  do
    For the node’s generated text feature  $\mathbf{x}_n^{\text{corr}}$ , determine whether each token is generated using the embedding of node  $n$  (computed during edge discrimination) and the feature discrimination Transformer model;
    Update the feature discrimination loss  $\mathcal{L}_d^f$ ;
  // Model updates.
  Compute  $\mathcal{L} = (\mathcal{L}_g^e + \mathcal{L}_g^f) + \lambda(\mathcal{L}_d^e + \mathcal{L}_d^f)$  and update the model;
  
```

**Output:** Trained model ready for fine-tuning.

---

## 493 B Generation and Discrimination of Vector Features

494 Node features can be vectors instead of texts, e.g., the feature vector can contain topological information such as connectivity information. In this case both the generator and the discriminator are  
 495 parameterized by a linear layer.  
 496

Table 6: Hyper-parameters for fine-tuning tasks.

Dataset	Task	Steps	Dropout	Learning rate	Gradient clipping
Reddit	—	2400	0.3	0.0015	0.5
Recomm.	—	1600	0.1	0.0010	0.5
OAG-CS	PF	1600	0.2	0.0010	0.5
	PV	1600	0.2	0.0005	0.5
	AD	1600	0.2	0.0005	0.5

497 To generate feature vectors, we first randomly select some nodes  $\mathcal{N}_g \subset \mathcal{N}$ . For a node  $n \in \mathcal{N}$ , denote  
 498 its feature vector  $\mathbf{v}_n$ , then the feature generation loss is

$$\mathcal{L}_g^f(W_g) = \sum_{n \in \mathcal{N}_g} \|\widehat{\mathbf{v}}_n - \mathbf{v}_n\|_2^2, \text{ where } \widehat{\mathbf{v}}_n = W_g^f h_g(n).$$

499 Here  $h_g(n)$  is the representation of node  $n$  and  $W_g^f$  is a trainable weight. For a node  $n \in \mathcal{N}$ , we  
 500 construct its corred feature  $\mathbf{v}_n^{\text{corr}} = \widehat{\mathbf{v}}_n$  if  $n \in \mathcal{N}_g$  and  $\mathbf{v}_n^{\text{corr}} = \mathbf{v}_n$  if  $n \in \mathcal{N} \setminus \mathcal{N}_g$ .

501 The discriminator’s goal is to differentiate the generated features from the original ones. Specifically,  
 502 the prediction probability is

$$p(n \in \mathcal{N}_g) = \text{sigmoid}(W_d^d h_d(n)),$$

503 where  $W_d^f$  is a trainable weight. We remark that the node representation  $h_d(n)$  is computed based on  
 504 the corred feature  $\mathbf{v}_n^{\text{corr}}$ . Correspondingly, the discriminator’s loss is

$$\mathcal{L}_d^f(W_d) = \sum_{n \in \mathcal{N}} -\mathbf{1}\{n \in \mathcal{N}_g\} \log p(n \in \mathcal{N}_g) - \mathbf{1}\{n \in \mathcal{N} \setminus \mathcal{N}_g\} \log(1 - p(n \in \mathcal{N}_g)).$$

505 The vector feature loss  $\mathcal{L}^f(\theta_g^e, W_g^f, \theta_d^e, W_d^f) = \mathcal{L}_g^f(\theta_g^e, W_g^f) + \mathcal{L}_d^f(\theta_d^e, W_d^f)$  is computed similar to  
 506 the text feature loss.

## 507 C Implementation and Training Details

508 By default, we use Heterogeneous Graph Transformer (HGT, [Hu et al. 2020c](#)) as the backbone GNN.  
 509 In the experiments, the edge generator and discriminator have the same architecture, where we set  
 510 the hidden dimension to 400, the number of layers to 3, and the number of attention heads to 8. For  
 511 the OAG dataset which contains text features, the feature generator and discriminator employs the  
 512 same architecture: a 4 layer bi-directional Transformer model, similar to BERT ([Devlin et al., 2019](#)),  
 513 where we set the embedding dimension to 128 and the hidden dimension of the feed-forward neural  
 514 network to 512.

515 For pre-training, we mask out 20% of the edges and 20% of the features (for text features we mask  
 516 out 20% of the tokens). We use AdamW ([Loshchilov and Hutter, 2019](#)) as the optimizer, where we  
 517 set  $\beta = (0.9, 0.999)$ ,  $\epsilon = 10^{-8}$ , the learning rate to 0.001 and the weight decay to 0.01. We adopt a  
 518 dropout ratio of 0.2 and gradient norm clipping of 0.5. For graph subsampling, we set the depth to 6  
 519 and width to 128, the same setting as [Hu et al. 2020b](#).

520 For fine-tuning, we use AdamW ([Loshchilov and Hutter, 2019](#)) as the optimizer, where we set  
 521  $\beta = (0.9, 0.999)$ ,  $\epsilon = 10^{-6}$ , and we do not use weight decay. We use the same graph subsampling  
 522 setting as pre-training. The other hyper-parameters are detailed in Table 6.