

---

# MULTI-MODAL PROMPT LEARNING EMPOWERS GRAPH NEURAL NETWORKS WITH SEMANTIC KNOWLEDGE

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

While great success has been achieved in building generalizable language models, three fundamental issues hinder GNN-based graph foundation models: the scarcity of labeled data, different levels of downstream tasks, and the conceptual gaps between domains. In depth, though the labels of real graphs are associated with semantic information, most graph learning frameworks ignore it by turning semantic labels into numerical labels. In this work, to address these issues, we present a new paradigm that leverages the text modality to align downstream tasks and data with any pre-trained GNN given only a few semantically labeled samples. Our paradigm embeds the graphs directly in the same space as the LLM by learning both graph prompts and text prompts simultaneously. To accomplish this, we improve state-of-the-art graph prompt method based on our theoretical findings. Then, we propose the first multi-modal prompt learning approach for exploiting the knowledge in pre-trained models. Notably, in our paradigm, the pre-trained GNN and the LLM are kept frozen, so the number of learnable parameters is much smaller than fine-tuning any pre-trained model. Through extensive experiments on real-world datasets, we demonstrate the superior performance of our paradigm in few-shot, multi-task-level, and cross-domain settings. Moreover, we build the first zero-shot classification prototype that can generalize GNNs to unseen classes. The code is provided in the supplementary materials.

## 1 INTRODUCTION

Foundation Models [4] learn generalizable representations from large-scale data and can be adapted to a wide range of downstream tasks. Although foundation models have shown remarkable capability and been thriving in NLP [14, 7, 112, 80], computer vision [2, 15, 65, 66, 44, 84], and time-series analysis [83, 110, 49], graph-related foundation models still remain in a very nascent stage. This is due to the significant difference of non-euclidean graph data from other data types. First, compared with language or vision data, graph data is very scarce [50, 10, 59] for foundation models. Second, the task space of graph data could be on node-level [86], edge-level [71], and graph-level [67]. Third, in general, language tokens and visual objects retain the same conceptual meaning across different distributions, but the same graph structure may have distinct interpretations in different domains, depending on how graphs were constructed from real scenarios. Thus, even if we have a pre-trained model, adapting it to various downstream tasks is not trivial.

Recently, some works [98, 8, 3, 82] reformulate the graphs into natural language descriptions and the graph tasks into natural language prompts, then query LLMs to generate the answer. However, since the LLMs are not directly trained from structured graph data [52], it is uncertain how LLMs could correctly solve those tasks without hallucinating [1, 31, 95, 105]. Nevertheless, graph neural networks (GNNs) are well-studied architectures for learning graph data [90, 17, 106], with theoretically provable expressiveness [94, 69, 61], better interpretability [13, 34, 74] and experimentally outstanding performance [89, 40]. Therefore, GNNs are expected to leverage their inherent advances for structure learning and inference on graphs in the era of big data and foundation models.

However, though tremendous efforts have been devoted to pre-train GNNs through self-supervision [92, 29, 54], a key problem in building a GNN-backed graph foundation model is that GNNs do not capture semantics, given that current GNNs are optimized according to numerical labels. In other words, GNNs do not really *understand* what a label represents in the real world, even though the graphs are constructed from real scenarios. To solve the issue of predetermined numerical categories, CLIP [65] leverages natural language supervision by jointly training an image encoder and a text encoder in the same embedding space to predict the correct image-text pairs at scale. The excellent generalization ability of pre-trained V-L models [65, 33, 47] comes from the alignment between

the vision and language representations. Notably, some works have explored prompt learning for better alignment and obtained improvement in vision prediction [114, 41]. The idea of alignment with text modality has also been applied in video [93, 6], 3D images [103, 25, 26], speech [70] and audio [23, 85] areas. As for graphs, so far such CLIP pipelines have only be applied in the molecular domain [58, 56, 72, 51], where the paired graph-text data are relatively sufficient for pre-training to align representations. But for other domains, such text-labeled graph data are rarely available, which means we have to rely more on self-supervised GNN pre-training to build graph foundation models. With this assumption, it is necessary to study how to make the pre-trained GNN aware of the semantics of downstream graph representations, which motivates the following question:

*How to adapt pre-trained GNNs to the semantic embedding space given limited downstream data?*

This paper aims to answer this question based on the following observations: (1) Semantic text embedding spaces do not necessarily result from joint pre-training. In fact, the embedding spaces of encoder LLMs are inherently semantic and high-quality, as LLMs are trained on massive text data and demonstrate strong reasoning performance. (2) When the downstream data are limited, prompt learning [48, 28, 102, 45] provides a better option than fine-tuning as much fewer parameters not only makes the optimization more efficient but also requires less resource than computing the gradient of a large model. Inspired by these two observations, we propose a prompting-based paradigm with an LLM that, while keeping the parameters of both GNN and LLM frozen, aligns the GNN representations in the LLM’s semantic embedding space.

Notably, when attempting to adapt the representation from one modality to another, solely prompting a single modality could be sub-optimal, as it limits the adjustment to downstream tasks in the other modality [41]. To this end, we propose Multi-modal Prompt Learning for Graph Neural Networks (Morpher). Given a pre-trained GNN and few-shot semantically labeled graph data, we introduce a pre-trained LLM. Then, to leverage its high-quality semantic embedding space, Morpher connects and aligns the graph embeddings to it through prompting on both modalities with a cross-modal projector. Nonetheless, designing such a paradigm is more challenging than vision-language models. First, we lack jointly pre-trained encoders for the two modalities; instead, we only have two encoders whose embedding dimension is possibly different, pre-trained independently in each modality. Second, determining how to prompt the graph modality is non-trivial and remains a trending research topic. Third, the downstream data for GNN usually have much fewer labeled classes than V-L models, so in the few-shot setting, the available downstream data is extremely limited. Our contributions towards tackling these challenges are summarized as follows:

- Theoretically, we analyze that, in many cases, state-of-the-art graph prompt [76] is unable to learn good representations of the downstream data. We show that the optimization of the graph prompt is restricted by design. From the theoretical findings, we further improve state-of-the-art graph prompt according to the attention mechanism to prevent failure in optimization.
- To connect and adapt the pre-trained GNN with LLM, we propose Morpher, a graph-text multi-modal prompt learning paradigm. To the best of our knowledge, this is the first approach to align the representations of GNN and LLM without fine-tuning any of their parameters.
- Experimentally, we demonstrate the effectiveness of our improved graph prompt and Morpher on real-world datasets under few-shot, multi-task, and cross-domain settings.

## 2 BACKGROUND

We use calligraphic letters (e.g.,  $\mathcal{A}$ ) for sets, and specifically  $\mathcal{G}$  for graphs. We use bold capital letters for matrices (e.g.,  $\mathbf{A}$ ). For matrix indices, we use  $\mathbf{A}(i, j)$  to denote the entry in the  $i^{th}$  row and the  $j^{th}$  column. Additionally,  $\mathbf{A}(i, :)$  returns the  $i^{th}$  row in  $\mathbf{A}$ .

**Graph Neural Networks.** We use  $\mathcal{G} = (\mathbf{A}, \mathbf{X})$  to denote a graph with node set  $\mathcal{V}$  and edge set  $\mathcal{E}$ , where  $\mathbf{A} \in \mathbb{R}^{|\mathcal{V}| \times |\mathcal{V}|}$  is the adjacency matrix and  $\mathbf{X} \in \mathbb{R}^{|\mathcal{V}| \times d}$  is the node feature matrix.  $\mathbf{A}(u, v) = 1$  if there is an edge connecting  $u$  and  $v$ ; otherwise  $\mathbf{A}(u, v) = 0$ . A Graph Neural Network  $f_{\phi}^g(\cdot)$  with hidden dimension  $d_g$  encodes  $\mathcal{G}$  into the embedding space:  $f_{\phi}^g(\mathcal{G}) \in \mathbb{R}^{|\mathcal{V}| \times d_g}$ , which could preserve both feature and structure information of  $\mathcal{G}$ . The extracted embeddings  $f_{\phi}^g(\mathcal{G})$  can be used for various downstream tasks such as classification. Nowadays, a popular paradigm to train GNNs is to first pre-train GNNs via self-supervised learning [29] and then fine-tune on the downstream tasks.

**Few-shot Prompt Learning.** Prompt learning adds learnable tokens to the downstream data and provides a powerful alternative to fine-tuning when the labeled downstream data is scarce. Prompt learning for encoders was first used in NLP. Let  $f_\phi^t(\cdot)$  denote the LLM encoder with embedding dimension  $d_t$ . For a series of input tokens  $\{x_k\}_{k=1}^K$ , the LLM encoder embeds it as a matrix  $\mathbf{X}_t = f_\phi^t(\{x_k\}_{k=1}^K) \in \mathbb{R}^{K \times d_t}$ , then aggregates the representation to a vector  $aggre(\mathbf{X}_t) \in \mathbb{R}^{1 \times d_t}$  for downstream tasks. Prompt learning initializes a tunable matrix  $\mathbf{P}_\theta^t \in \mathbb{R}^{n_t \times d_t}$ , where  $n_t$  denotes the number of text prompt tokens. Then, this tunable matrix is concatenated with the input tokens' embeddings to form a single matrix  $[\mathbf{P}_\theta^t; \mathbf{X}_t]_{dim=0} \in \mathbb{R}^{(K+n_t) \times d_t}$ , and the aggregated vector for downstream tasks becomes  $aggre([\mathbf{P}_\theta^t; \mathbf{X}_t]_{dim=0})$ . In practice, we can train the model to minimize the loss function for downstream tasks, with only the prompt parameters  $\mathbf{P}_\theta^t$  being updated.

Now we are ready to introduce the problem setup for this work. Given a pre-trained GNN  $f_\phi^g(\cdot)$  with embedding dimension  $d_g$  and a pre-train LLM encoder  $f_\phi^t(\cdot)$  with embedding dimension  $d_t$ . Without loss of generality, we assume the downstream task is graph-level classification, as we will show that the other types of GNN tasks can be reformulated as graph classification. For  $L$ -shot graph classification, we are given limited text-labeled pairs  $\{(\mathcal{G}_i, t_c)\}_{i=1}^L$  for each class  $c$ . Assuming  $\mathcal{T}$  is the set of all text labels  $t_c$ , we are provided a set of test graphs  $\{\mathcal{G}_j\}_{j=1}^{L_{test}}$ . Using the pre-trained GNN and LLM, we want to correctly predict the text label  $t_j \in \mathcal{T}$  for each test graph  $\mathcal{G}_j$ .

### 3 REVISITING AND IMPROVING PROMPT AS GRAPHS

Unlike prompting text data (which can be easily achieved by appending learnable text tokens to the original text sequence) and prompting image data (which pads a learnable image area above the original image), prompting graph data presents a significant challenge due to the non-euclidean nature of graphs. The recent pioneering work [76] designs the graph prompt still as a graph, then inserts it into the original graph by computing the inner-connections within the prompt graph and the cross-connections between the prompt graph and the original graph. An advantage of prompting at the graph level is that *the downstream tasks of GNN can be reformulated into graph-level tasks*. For the node classification task, we can induce the  $\gamma$ -ego-graph of each node by extracting the subgraph within a pre-defined distance  $\gamma$ . Then, we treat the node label as the induced ego-graph label. Similarly, for the edge classification task, we can extract a subgraph for each edge by extending the node pair to their  $\gamma$  distance neighborhood, and use the edge label as the induced graph label. By inducing subgraphs, we can reformulate node-level and edge-level downstream tasks to graph-level.

**Current Graph Prompt Design.** To prompt a graph  $\mathcal{G}$ , each prompt token is a new node. Let  $n_g$  denote the number of prompt tokens and  $\mathcal{P} = \{p_i\}_{i=1}^{n_g}$  denote the set of prompt tokens. The graph prompt is formulated by a tunable matrix  $\mathbf{P}_\theta^g \in \mathbb{R}^{n_g \times d}$ , where  $d$  is the node feature dimension of graph  $\mathcal{G}$ . In other words, each row vector  $\mathbf{P}_\theta^g(i, \cdot)$  is the feature of the prompt token  $p_i$ . Then, the mechanism to prompt a graph  $\mathcal{G} = (\mathbf{A}, \mathbf{X})$  with  $n$  nodes and  $d$  feature dimension is as follows [76].

- Compute inner-connections to construct the prompt graph  $\mathcal{G}_p = (\mathbf{A}_p, \mathbf{X}_p)$ . For the feature matrix, we directly set  $\mathbf{X}_p = \mathbf{P}_\theta^g$ . For two prompt tokens  $p_i$  and  $p_j$ , the prompt graph will have an edge between them if and only if the dot product of their features is larger than a threshold. In other words,  $\mathbf{A}_p(i, j) = 1 \iff \sigma(\mathbf{P}_\theta^g(i, \cdot) \mathbf{P}_\theta^g(j, \cdot)^\top) > \delta_{inner}$ , where  $\sigma(\cdot)$  is the sigmoid function.
- Compute cross-connections to insert the prompt graph  $\mathcal{G}_p$  into the original input graph  $\mathcal{G}$ . Similarly, for  $x_i \in \mathcal{G}$  and  $p_j \in \mathcal{G}_p$ , there is an edge between them if and only if  $\sigma(\mathbf{X}(i, \cdot) \mathbf{P}_\theta^g(j, \cdot)^\top) > \delta_{cross}$ .
- Construct the prompted graph (i.e., manipulated graph)  $\mathcal{G}_m = (\mathbf{A}_m, \mathbf{X}_m)$ . The overall adjacency matrix  $\mathbf{A}_m \in \mathbb{R}^{(n+n_g) \times (n+n_g)}$  is constructed from the original adjacency matrix  $\mathbf{A}$ , the inner edges  $\mathbf{A}_p$  and the cross edges. The overall node feature matrix is concatenated from the prompt token features and the original input node features:  $\mathbf{X}_m = [\mathbf{P}_\theta^g; \mathbf{X}]_{dim=0} \in \mathbb{R}^{(n+n_g) \times d}$ .

Here, we identify an issue associated with the current design. Since not all the GNN backbones can take edge weights [21], the cross-connections in a manipulated graph are discrete<sup>1</sup>, thresholded by  $\delta_{cross}$ . However, the input node features of most real-world datasets are sparse, resulting from the construction process [97, 60, 18]. As shown in Table 6,  $\|\mathbf{X}(i, \cdot)\|_1$  is typically 1. As the initialization of each token feature  $\mathbf{P}_\theta^g(i, \cdot)$  is close to  $\vec{0}$ , for any node  $i$  and token  $p_j$ , the dot products

<sup>1</sup>In official implementation of [76], adjacency matrices are discrete: either 0 or 1 for each entry.

162  
163  
164  
165  
166  
167  
168  
169  
170  
171  
172  
173  
174  
175  
176  
177  
178  
179  
180  
181  
182  
183  
184  
185  
186  
187  
188  
189  
190  
191  
192  
193  
194  
195  
196  
197  
198  
199  
200  
201  
202  
203  
204  
205  
206  
207  
208  
209  
210  
211  
212  
213  
214  
215

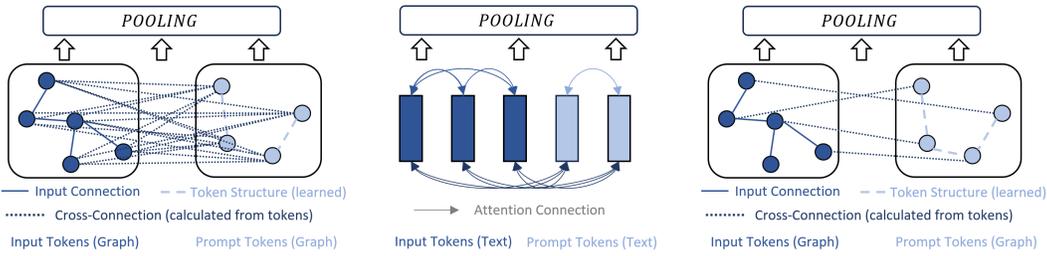


Figure 1: Illustration of connections in current problematic graph prompt design (left), transformer architecture (middle), and our improved graph prompt design (right). **The cross-connections between input and prompt should be consistent with the input connections in scale.**

$\mathbf{X}(i, :)\mathbf{P}_\theta^g(j, :)^T$  is close to 0, and the sigmoid value is very close to 0.5. As a result, if we want the graph prompt to work reasonably, we have to set  $\delta_{cross} < 0.5$ . However, in this case, the cross-connections will be dense, i.e., almost every node in the original graph is connected with every node token in the prompt graph. For two different graphs  $\mathcal{G}_1$  and  $\mathcal{G}_2$  in the same task, the prompt graph  $\mathcal{G}_p$  is identical. Since the GNNs work by aggregating the node features, their embeddings  $f_\phi^g(\mathcal{G}_1)$  and  $f_\phi^g(\mathcal{G}_2)$  are approximately the same because the features in the prompt graph overwhelm the features in the original graphs due to the dense cross-connections. Then, according to the following lemma, even if  $\mathcal{G}_1$  and  $\mathcal{G}_2$  have different labels, the task head classifier cannot be trained to distinguish them<sup>2</sup>.

**Lemma 3.1.** *For any classifier  $c(\cdot)$ , if the identical feature  $\mathbf{x}$  has label distribution  $p(\cdot)$ , then the optimal classification for cross-entropy loss is  $Pr(c(\mathbf{x}) = y) = p(y)$ . From this, if two graphs have similar embedding but different labels, GNN training may not converge. (Proof in Appendix A)*

**Improved Graph Prompt Design.** The issue of the current graph prompt is rooted in the imbalance of original connections in the input graph and cross-connections between input and prompt, as shown in Figure 1 (left). Since the text prompt works well in NLP, we look into the standard transformer architecture [81], where the token features are aggregated through the attention mechanism:

$$\tilde{\mathbf{H}} = \text{Attn}_{\theta_a}(\mathbf{H}) := \mathbf{H} + \frac{1}{N} \sum_{m=1}^M (\mathbf{V}_m \mathbf{H}) \times \sigma((\mathbf{Q}_m \mathbf{H})^T (\mathbf{K}_m \mathbf{H})) \in \mathbb{R}^{D \times N} \quad (1)$$

where  $\mathbf{H} \in \mathbb{R}^{D \times N}$  is the input sequence and  $\theta_a = \{(\mathbf{V}_m, \mathbf{Q}_m, \mathbf{K}_m)\}_{m \in [M]} \subset \mathbb{R}^{D \times D}$  denotes the parameters with  $M$  heads.  $N$  is number of tokens and  $D$  is embedding dimension. We also visualize such attention mechanism in Figure 1 (middle). After we prepend a sequence of text prompt tokens  $\{p_i^t\}$ , the features of the text prompt tokens will be densely aggregated to the features of the original text tokens. In other words, the ‘‘cross-connection’’ between the text prompt sequence and the input sequence is dense. However, such a dense connection does not cause the prompt feature to overwhelm the input, because the features in the input sequence are also aggregated in a dense manner. Inspired by this, the number of cross-attention between input and prompt should approximate the number of input connections. Since the connection of a graph dataset is often sparse, we should also constrain the cross-connections between the prompt graph and the input graph to be sparse as well.

Nonetheless, ‘‘sparse’’ is a wide concept to implement: if the cross-connections are too dense, the prompt graph will dominate the input graph; but if the cross-attention is too sparse, the prompt graph will be limited to manipulating the input graph. We deem that a balance could be achieved by approximately equalizing the number of cross-connections with that of connections in the input graph, i.e.,  $n_e$ . Therefore, we set the number of cross-connections to at most  $n_e$  by connecting each node in the input graph with at most  $\lfloor \frac{n_e}{a} \rfloor$  prompt tokens. Then, we can safely use a small  $\delta_{cross}$  and cosine similarity  $\frac{\mathbf{X}(i, :)\mathbf{P}_\theta^g(j, :)^T}{\|\mathbf{X}(i, :)\|_2 \|\mathbf{P}_\theta^g(j, :)\|_2}$  instead of  $\sigma(\mathbf{X}(i, :)\mathbf{P}_\theta^g(j, :)^T)$  to calculate the cross-connections. We demonstrate that our improved graph prompt works better in the later experiments.

#### 4 MULTI-MODAL PROMPT LEARNING FOR GNNs

To adapt the GNN embeddings to the LLM’s semantic embedding space and leverage the additional supervision provided by the text associated with graph labels, we explore the potential of multi-modal

<sup>2</sup>In fact, when executing the official implementation of [76] on Cora, the training loss does not decrease. Similar problems have been observed by another work [108].

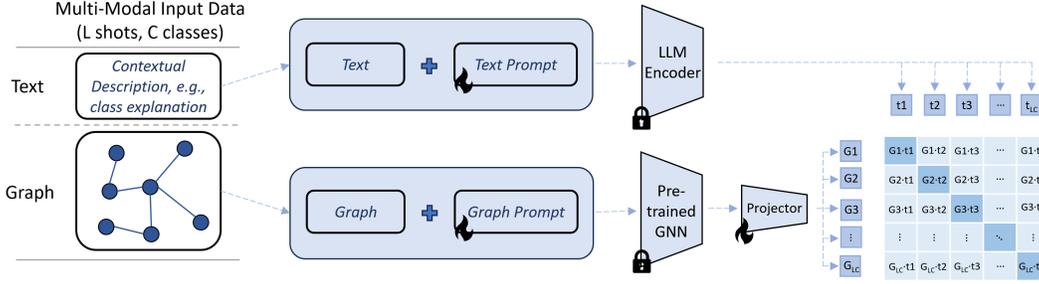


Figure 2: Similar to CLIP backbone, Morpher adapts the graph representations to semantic space through multi-modal prompt learning, even if the GNN and LLM are not jointly trained and are kept frozen.

prompt learning for both graphs and language. This approach is motivated by the intuition that only prompting on the graph data may limit the flexibility to adjust the LLM representation space. The overall paradigm of Morpher is illustrated in Figure 2. Given the data  $\{(\mathcal{G}_i, t_i)\}_{i=1}^{L \times C}$ , we aim to align graph embedding readout  $(f_\phi^g(\mathcal{G}_i))$  with readout  $(f_\phi^t(\text{Tokenize}(t_i)))$ . Yet one direct issue is that,  $\text{readout}(f_\phi^g(\mathcal{G}_i)) \in \mathbb{R}^{1 \times d_g}$  and  $\text{readout}(f_\phi^t(\text{Tokenize}(t_i))) \in \mathbb{R}^{1 \times d_t}$  may have distinct dimensions. To address this issue, we adopt a cross-modal projector that learns to map the graph embedding space to the text embedding space. For an input  $d_g$ -dimensional vector  $\mathbf{h}$ , the projector maps it to a  $d_t$ -dimensional vector  $\tilde{\mathbf{h}}$ :

$$\tilde{\mathbf{h}} = \text{Proj}_\theta(\mathbf{h}) := \tanh(\mathbf{W}\mathbf{h} + \mathbf{b}) \in \mathbb{R}^{1 \times d_t} \quad (2)$$

As discussed in Sections 2 and 3, we introduce the text prompt  $\mathbf{P}_\theta^t \in \mathbb{R}^{n_t \times d_t}$  with  $n_t$  text prompt tokens and the graph prompt  $\mathbf{P}_\theta^g \in \mathbb{R}^{n_g \times d}$  with  $n_g$  graph prompt tokens. Let  $\psi_g(\cdot, \mathbf{P}_\theta^g)$  be the graph prompting function, e.g., given any graph  $\mathcal{G}$ , the manipulated graph  $\mathcal{G}_m = \psi_g(\mathcal{G}, \mathbf{P}_\theta^g)$ .

Let  $\omega_t(\cdot, \mathbf{P}_\theta^t)$  be the prompted text embedding given input text  $t$ . For the text prompt methods we choose, the prompted embedding is

$$\omega_t(t, \mathbf{P}_\theta^t) = [\mathbf{P}_\theta^t; f_\phi^t(\text{Tokenize}(t))]_{dim=0} \in \mathbb{R}^{(\text{len}(\text{Tokenize}(t)) + n_t) \times d_t} \quad (3)$$

Let  $\omega_g(\cdot, \mathbf{P}_\theta^g)$  be the prompted graph embedding given input graph  $\mathcal{G}$ , then we have:

$$\omega_g(\mathcal{G}, \mathbf{P}_\theta^g) = f_\phi^g(\mathcal{G}_m) = f_\phi^g(\psi_g(\mathcal{G}, \mathbf{P}_\theta^g)) \in \mathbb{R}^{(n+n_g) \times d_g} \quad (4)$$

For the whole prompted text and the whole prompted graph, we apply readout (e.g., mean-pooling, max-pooling, etc.) to get their embedding:

$$e^t = \text{readout}(\omega_t(t, \mathbf{P}_\theta^t)) \in \mathbb{R}^{1 \times d_t}, \quad e^g = \text{readout}(\omega_g(\mathcal{G}, \mathbf{P}_\theta^g)) \in \mathbb{R}^{1 \times d_g} \quad (5)$$

For the given data  $\{(\mathcal{G}_i, t_i)\}_{i=1}^L$ , we compute the normalized embedding of prompted  $\mathcal{G}_i$  and project it to the text embedding space through the projector:

$$z_i^{\mathcal{G}^{norm}} = \frac{e_i^g}{\|e_i^g\|_2} = \frac{\text{readout}(\omega_g(\mathcal{G}_i, \mathbf{P}_\theta^g))}{\|\text{readout}(\omega_g(\mathcal{G}_i, \mathbf{P}_\theta^g))\|_2}, \quad z_i^{\mathcal{G}} = \text{Proj}_\theta(z_i^{\mathcal{G}^{norm}}) \quad (6)$$

For the text embeddings, since for limited data the set  $\mathcal{T} = \{t_i\}_{i=1}^C$  may contain texts that are semantically close as discussed in Appendix B.2, we extract a subspace in the text embedding space by normalizing the embedding as follows. We further normalize the text embeddings to the unit sphere, as standard practice in NLP.

$$\mu_t = \frac{1}{L} \sum_{i=1}^L \text{readout}(\omega_t(t_i, \mathbf{P}_\theta^t)), \quad e_{norm,i}^t = \text{readout}(\omega_t(t_i, \mathbf{P}_\theta^t)) - \mu_t \quad (7)$$

$$z_i^t = \frac{e_{norm,i}^t}{\|e_{norm,i}^t\|_2} = \frac{\text{readout}(\omega_t(t_i, \mathbf{P}_\theta^t)) - \mu_t}{\|\text{readout}(\omega_t(t_i, \mathbf{P}_\theta^t)) - \mu_t\|_2} \quad (8)$$

Finally, we use the in-batch similarity-based contrastive loss to train text prompts, graph prompts, and the projector as shown below, to adapt the pre-trained GNN representations to LLM.

$$\mathcal{L}_{G \rightarrow T} = -\frac{1}{B} \sum_{i=1}^B \log \frac{\exp(z_i^{\mathcal{G}} \cdot z_i^t / \tau)}{\sum_{j=1}^B \exp(z_i^{\mathcal{G}} \cdot z_j^t / \tau)} \quad (9)$$

Table 1: Few-shot graph classification performance (%). IMP (%): the average improvement (absolute value) compared to the **best result** among all the baseline methods.

Training schemes	GNN pretraining	MUTAG		ENZYMES		PROTEINS		MSRC_21C	
		Acc	F1	Acc	F1	Acc	F1	Acc	F1
Supervised	N/A + GCN	66.00	66.67	16.67	8.68	65.89	60.77	38.85	35.32
	N/A + GAT	66.00	65.69	16.45	4.65	64.75	64.08	41.14	39.86
	N/A + GT	66.66	66.26	15.62	4.22	62.81	57.12	38.28	41.62
Pre-train + Fine-tune	GraphCL+GCN	70.00	70.23	17.91	11.82	65.89	61.23	40.00	43.89
	GraphCL+GAT	70.00	69.73	17.91	10.46	65.16	63.92	44.57	45.74
	GraphCL+GT	68.00	67.81	17.70	8.99	63.28	56.41	41.71	43.73
	SimGRACE+GCN	66.67	67.27	17.29	8.78	66.82	64.70	40.57	43.84
	SimGRACE+GAT	70.67	69.10	16.87	7.18	65.42	63.65	42.85	42.37
	SimGRACE+GT	69.33	69.77	16.24	6.08	65.98	62.31	39.42	40.78
AIO [76]	GraphCL+GCN	64.67	39.27	17.50	4.97	61.35	44.93	3.59	10.09
	GraphCL+GAT	64.67	39.27	17.50	4.97	59.21	37.19	14.37	3.11
	GraphCL+GT	73.33	72.06	18.33	9.09	40.79	28.97	17.96	8.30
	SimGRACE+GCN	64.67	39.27	16.04	4.61	67.42	60.87	34.73	18.16
	SimGRACE+GAT	64.67	39.27	16.04	4.61	59.21	37.19	7.78	1.79
	SimGRACE+GT	36.00	27.26	17.50	8.15	50.56	49.34	32.34	15.13
GPF-plus [19]	GraphCL+GCN	68.67	67.27	16.88	15.48	64.75	61.45	47.42	29.02
	GraphCL+GAT	68.67	62.84	16.45	13.23	65.89	60.07	47.42	26.28
	GraphCL+GT	69.33	67.87	18.12	15.56	59.66	37.37	41.71	21.35
	SimGRACE+GCN	65.33	39.52	18.96	15.83	65.16	58.80	45.71	23.32
	SimGRACE+GAT	69.33	66.72	18.54	12.58	63.28	53.50	42.85	21.40
	SimGRACE+GT	70.00	67.31	17.91	14.69	64.83	52.97	34.13	20.13
Gprompt [55]	GraphCL+GCN	73.33	66.93	17.91	8.44	61.01	60.01	1.80	0.21
	GraphCL+GAT	64.67	62.63	17.08	14.18	50.56	50.55	1.80	0.22
	GraphCL+GT	70.67	70.02	17.91	9.64	63.28	58.65	1.80	0.21
	SimGRACE+GCN	65.33	39.52	17.29	14.48	52.70	52.68	1.80	0.21
	SimGRACE+GAT	67.33	65.88	16.25	11.31	59.10	58.72	1.80	0.21
	SimGRACE+GT	73.33	67.84	16.87	13.54	64.75	62.37	1.80	0.223
Improved AIO (Ours)	GraphCL+GCN	77.33	77.74	18.13	11.98	65.89	65.97	42.85	45.91
	GraphCL+GAT	74.67	75.51	18.33	11.26	65.76	66.05	46.85	51.39
	GraphCL+GT	74.67	74.67	19.16	9.04	68.12	68.18	42.85	43.54
	SimGRACE+GCN	68.00	69.01	17.91	9.02	66.82	66.40	44.57	49.24
	SimGRACE+GAT	77.33	77.20	18.75	9.39	66.91	65.49	45.14	42.31
	SimGRACE+GT	71.33	72.06	18.95	11.25	68.59	68.84	40.57	42.82
Morpher (Ours)	GraphCL+GCN	78.67	78.09	20.41	15.20	67.47	66.40	45.14	49.62
	GraphCL+GAT	79.33	79.15	23.12	18.01	70.89	70.30	50.85	54.48
	GraphCL+GT	76.00	76.51	19.58	13.28	73.53	72.48	45.71	48.41
	SimGRACE+GCN	69.33	70.27	19.79	14.94	67.10	66.15	45.71	51.24
	SimGRACE+GAT	78.00	77.65	20.21	16.27	68.12	67.26	45.71	51.13
	SimGRACE+GT	74.00	74.84	19.16	14.29	71.76	71.75	44.00	48.16
IMP of ImprovedAIO		2.00 ↑	5.01 ↑	0.52 ↑	4.41 ↓	2.01 ↑	4.37 ↑	0.28 ↓	2.50 ↑
IMP of Morpher		4.00 ↑	6.73 ↑	2.36 ↑	0.60 ↑	4.81 ↑	6.61 ↑	2.66 ↑	7.14 ↑

## 5 EXPERIMENTS

We evaluate our Morpher and the improved graph prompt through extensive experiments. In particular, we show that, compared to state-of-the-art baseline methods, they both more effectively adapt pre-trained GNNs to the specific downstream classification task, and introducing the text modality brings Morpher additional advantages over others. We use RoBERTa [53] as the LLM encoder for Morpher in the main experiments. We also validate the performance of Morpher with ELECTRA [12] and DistilBERT [68] in section 5.6 and Appendix C.3.

*Datasets.* We use real-world graph datasets from PyTorch Geometric [21], including one molecular dataset MUTAG [60]; two bioinformatic datasets ENZYMES and PROTEINS [5]; one computer vision dataset MSRC\_21C [63]; three citation network datasets Cora, CiteSeer and PubMed [97]. We use real-world class names as text labels. More details are summarized in Appendix B.

*Pre-trained algorithms and GNN backbones.* To pre-train GNNs for evaluation, we adopt GraphCL [99] and SimGRACE [91] to pre-train three widely used GNN backbones: GCN [43], GAT [100] and GraphTransformer (GT) [42]. Additionally, in Appendix C.4, we verify the effectiveness of our methods on GNNs pre-trained using GraphMAE [27] and MVGRL [24], two other representative GNN self-supervised learning algorithms. For each dataset, to pre-train GNNs, we leverage self-supervised learning methods on all the graphs without any label information.

*Baselines and metrics.* We compare our methods with the following baselines: (1) training a GNN from scratch supervised by few-shot data (“*supervised*”); (2) fine-tuning a task head together with pre-trained GNN (“*fine-tune*”). We allow GNNs to be tunable for “*supervised*” and “*fine-tune*”; (3) state-of-the-art graph prompting algorithms: All-in-one (“*AIO*”) [76], which is the only graph prompting algorithm that supports multiple tasks in node-level, edge-level and graph-level to the best of our knowledge; GPF-plus [19] which prompt on graph features and Gprompt [55] which is based on subgraph similarity. We use accuracy and weighted F1 as classification performance metrics.

### 5.1 FEW-SHOT LEARNING

We investigate the ability of our improved graph prompt (“*ImprovedAIO*”) and Multimodal prompt (“*Morpher*”) to adapt frozen pre-trained GNNs using few-shot data. We focus on graph-level classification here and will further investigate the few-shot learning ability at other task levels in Section 5.2. Our few-shot learning setting is more challenging than existing works [76, 75] as we only allow no more than 10 labeled training and validation samples for each class. The results are shown in Table 1, where we report the average performance of 5 runs and calculate the absolute average improvement of our methods. From the results, given the same pre-trained GNN, our ImprovedAIO outperforms all the existing baseline methods except for ENZYMES F1 and MSRC\_21C accuracy. Yet the performance of our ImprovedAIO on ENZYMES F1 and MSRC\_21C accuracy is clearly better than those of the original AIO. Our Morpher can achieve an absolute accuracy improvement of 0.60% to 7.14% over the baselines across all datasets. Supervised by very limited labeled data, training a GNN from scratch is sub-optimal. Passing a GNN pre-trained on the dataset and fine-tuning it with a task head achieves sub-optimal but better results as the pre-trained GNN learns generalizable representations over the dataset through self-supervised learning. To mitigate the gap between the pre-training task and downstream tasks, AIO [76] proposes to learn graph prompts for downstream data. However, as we discussed in Section 3, when the node features are sparse vectors, the optimization would fail. Using the official implementation of AIO, we observe that the loss value tends to fluctuate, and the performance of AIO is usually even worse than supervised training. By restricting the cross-connections, our ImprovedAIO becomes more stable and constantly outperforms the fine-tuning baseline. Compared to the aforementioned methods, Morpher demonstrated superior performance due to its capability to adapt both graph and language representation spaces dynamically.

### 5.2 MORPHER SUPPORTS MULTIPLE-LEVEL TASKS

Inherited from AIO, our ImprovedAIO and Morpher also support adaptation to downstream tasks at node-level and edge-level, because they can be reformulated into graph-level tasks as discussed in Section 3. We demonstrate the performance of node classification and link prediction on Cora and CiteSeer. For node classification, we reformulate it to graph classification by inducing an ego-graph with 10 to 30 nodes centered at the node to classify. Each ego-graph has the same label as the center node. For edge classification, we randomly sample 200 edges from the graph, then create 200 negative samples by replacing one node in each edge. We label each graph according to whether it is a positive or negative sample.

Table 2: Node-level, edge-level performance.

Dataset		Cora		CiteSeer	
Tasks	Methods	Acc	F1	Acc	F1
Node Level	Supervised	52.83	47.73	63.91	64.82
	Fine-tune	56.37	55.04	64.87	66.42
	AIO [76]	14.69	7.10	18.93	6.92
	ImprovedAIO	58.46	55.10	66.44	66.53
	Morpher	61.26	62.36	68.20	68.56
Edge Level	Supervised	51.78	50.62	52.14	50.81
	Fine-tune	52.50	51.00	52.50	51.12
	AIO [76]	50.00	33.33	50.00	33.33
	ImprovedAIO	54.64	54.57	53.92	53.55
	Morpher	55.71	55.05	55.35	55.05

We use GraphCL+GCN to pre-train the GNN and report the mean performance in Table 2. The results are consistent with graph-level performance, where ImprovedAIO and Morpher outperform existing methods, with Morpher achieving slightly better performance than ImprovedAIO. Additionally, the training of the original AIO fails on both datasets due to the sparse node feature vectors.

### 5.3 DOMAIN TRANSFER



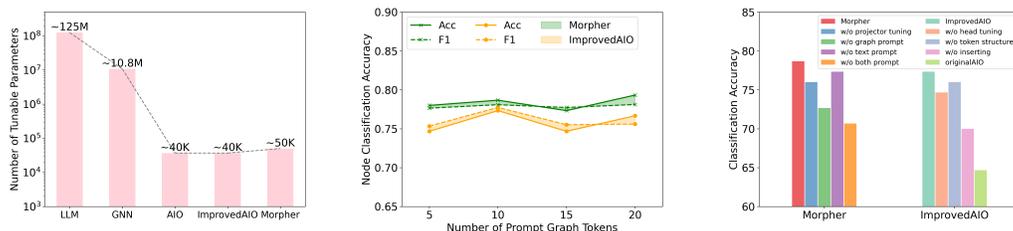


Figure 4: Efficiency comparison (left), parameter study (middle) and ablation study (right).

overfits, there is a period when Morpher can distinguish all the graphs from the training and novel classes with high accuracy.

Such zero-shot novel-class generalization ability validates Morpher’s alignment between graph embeddings and text embeddings. When Morpher is trained on two classes of graphs with text labels of biology and informatics, a graph-in-the-middle will be classified as text-in-the-middle: bioinformatics, even if “bioinformatics” is an unseen label. The correspondence of in-the-middle graphs and texts shows the benefit and novelty of Morpher. To the best of our knowledge, this is the **first zero-shot classification prototype that generalizes GNN to unseen classes**.

### 5.5 EFFICIENCY AND EMBEDDING ANALYSIS

Without fine-tuning the GNN or LLM, the prompt-based methods have better parameter efficiency. As shown in Figure 4 (left), our ImprovedAIO and Morpher require similar numbers of parameters with AIO [76], which is 0.032% to 0.46% compared to either tune the LLM (RoBERTa) or GNN (GCN). Due to such parameter efficiency, our methods learn better graph representations given few-shot data. We visualize the graph embeddings of CiteSeer and MSRC\_21C in Figure 3 and calculate the silhouette score, a metric for cluster quality ( $\uparrow$ ) ranged in  $[-1, 1]$ . It turns out that our multimodal prompting leads to better adaptation.

### 5.6 HYPERPARAMETER AND ABLATION STUDY

We conduct the hyperparameter study by choosing and testing various numbers of graph prompt tokens for both ImprovedAIO and Morpher. The results are shown in Figure 4 (middle), from which we can observe that both methods are generally stable, and Morpher constantly outperforms ImprovedAIO under different choices. To verify the necessity of each component in our design, we compare Morpher and ImprovedAIO with multiple variants, respectively, and report the result in Figure 4 (right). We observe that removing any component would result in a performance drop. Additionally, our comparison of Morpher with ImprovedAIO throughout the experiments demonstrates that our multimodal design would lead to improvement over the uni-modal prompting of GNNs.

In the main experiments, we use RoBERTa as Morpher’s text encoder. We also conduct experiments to verify the effectiveness of our proposed Morpher with ELECTRA [12] and DistilBERT [68] as the text encoder. Due to space limitation, we only show the F1 score of using ELECTRA in Figure 5, and more detailed experiment data can be found in Appendix C.3. In general, using ELECTRA and DistilBERT results in similar performance compared to using RoBERTa, showing the robustness of Morpher with respect to the language encoder.

Table 4: Effectiveness (F1 score) of Morpher with ELECTRA [12] as the text encoder.

GNN pretraining	MUTAG	ENZYMES	PROTEINS	MSRC_21C
GraphCL + GCN	78.17	15.79	65.66	47.19
GraphCL + GAT	75.75	11.37	65.66	49.01
GraphCL + GT	77.04	14.68	72.70	44.09
SimGRACE + GCN	70.99	12.41	67.77	48.44
SimGRACE + GAT	77.51	13.31	67.78	49.43
SimGRACE + GT	73.55	15.76	70.28	44.50

As for the robustness with respect to the pre-trained GNNs, in the main experiments, we adopt two pre-train methods, GraphCL and SimGRACE to pre-train three different GNN architectures: GCN, GAT and GT. We further conduct experiments using GNNs pre-trained from GraphMAE [27] and MVGRL [24]. Due to the space limitation, we report the results and discuss in Appendix C.4.

### 5.7 MORPHER ON MOLECURENET

In this section, we demonstrate that, though not specifically designed for any downstream applications, the Morpher framework has the potential to be used in various downstream tasks, such as AI4Science

tasks. As for a case study, We use bace (inhibitors of human beta-secretase), tox21 (toxicology in the 21st century) and hiv (inhibit HIV replication) from MoleculeNet [88]. These three datasets have 1513, 7831, and 41127 graphs to classify, respectively. In these datasets, each graph label is associated with a text description. The tasks on bace and hiv are bio-activity prediction and the task on tox21 is toxicity prediction. To adopt Morpher, we use GraphCL to pre-train the GAT model and initialize the text prompts and text labels using those from GIMLET [107].

Table 5: AUC-ROC ( $\uparrow$ ) on MoleculeNet (bace, tox21, hiv). Morpher-K denotes K shots.

Dataset	KVPLM	MoMu	Galactica-1.3B	GIMLET-64M-50-shots	GAT-1M-supervised	Morpher-10	Morpher-20	Morpher-50
bace	0.5126	0.6656	0.5648	0.729	0.697	0.6231	0.6513	0.6858
tox21	0.4917	0.5757	0.4946	0.652	0.754	0.6769	0.7275	0.7459
hiv	0.6120	0.5026	0.3385	0.721	0.729	0.5742	0.7034	0.7283

KVPLM [101], MoMu [72], Galactica-1.3B [79] are zero-shot predictors for the three tasks; GIMLET-64M-50-shots is the GIMLET [107] model fine-tuned on 50 additional training samples<sup>3</sup>; GAT-1M-fully-supervised uses all the training data to train a GAT. Our Morpher-k-shots uses only k training samples. From the results, first, using only 10 training samples, Morpher can outperform the zero-shot baselines KVPLM, MoMu, and Galactica-1.3B. Second, using only 50 shots, Morpher can achieve similar performance with the fully supervised GAT. Third, using the same amount of few-shot data (50 shots), Morpher-50 outperforms GIMLET-64M-50-shots on tox21 and hiv, the two largest datasets among the three. This means our graph-text multi-modal prompt learning, with much fewer learnable parameters ( $\sim 50K$ ), is more sample-efficient than fine-tuning language model encoder.

## 6 RELATED WORK

**GNN Pre-training.** Recently, a surge of graph pre-training strategies have emerged to address the issue of label scarcity in graph representation learning [29, 57, 75, 46, 39, 113]. The main idea of pre-trained graph models is to capture general graph information across different tasks and transfer this knowledge to the target task using techniques such as contrastive predictive coding [42, 20, 64, 91], context prediction [62, 30], prompt tuning [75, 19], and mutual information maximization [62, 73, 35]. For instance, [29] proposes to learn transferable structural information from three levels of graph topology, including node-level, subgraph-level, and graph-level. Different from these approaches, this paper aims to build up foundational GNNs by leveraging multi-modal prompt learning techniques.

**Graph Prompt Learning.** Prompting is now mainstream for adapting NLP tasks, and recent studies exploring prompt learning for GNNs mark a thriving research area [77, 87]. It is a promising way to adapt GNNs to downstream tasks through token-level [19, 78, 9, 75, 116] or graph-level [76, 32, 22] prompting. Among all the existing methods, All-in-one (AIO) [76] is the only algorithm to learn tunable graph prompts for node-level, edge-level or graph-level downstream tasks given few-shot labeled data (Table 8). Based on our improved AIO, we present a pioneer study to explore learning prompts in multiple modalities simultaneously while keeping the pre-trained models frozen.

**LLM on Graphs.** Inspired by the advances of large language models in NLP [111], researchers have begun to explore their potential for graph-related tasks [36]. Current approaches can be divided into two main categories. The first category employs LLMs as pre-trained feature extractors to enhance GNNs [16, 11, 115]. For example, GLEM [109] proposes to input the language representation as initial features for the GNN and train them iteratively. The second category focuses on integrating graph structures directly into LLM architectures [96, 104, 38]. A notable example is Patton [37], which pre-trains a joint architecture on text-attributed graphs. Despite these advancements, none of them have explored the collaboration between LLMs and GNNs under graph prompt learning.

## 7 CONCLUSION

In this work, we introduce Morpher, the first multimodal prompt learning paradigm that can semantically adapt pre-trained GNNs to downstream tasks with the help of LLM, while keeping both the pre-trained models frozen. To build Morpher, we first analyze the limitations of the state-of-the-art graph prompting technique and propose an improved version. Through extensive experiments, we demonstrate that our improved AIO can achieve outperformance, and our Morpher has further improvements in few-shot, multi-level task, or domain transfer settings. Additionally, using Morpher, we build the first GNN zero-shot classifier prototype that can be generalized to novel testing classes.

<sup>3</sup>the performance of GIMLET and other baselines are directly from the GIMLET paper [107].

540  
541  
542  
543  
544  
545  
546  
547  
548  
549  
550  
551  
552  
553  
554  
555  
556  
557  
558  
559  
560  
561  
562  
563  
564  
565  
566  
567  
568  
569  
570  
571  
572  
573  
574  
575  
576  
577  
578  
579  
580  
581  
582  
583  
584  
585  
586  
587  
588  
589  
590  
591  
592  
593

## 8 ETHICS STATEMENT

There are no ethical concerns associated with this research. The datasets and related resources used in this study are publicly accessible and have been widely employed in the existing works.

## 9 REPRODUCIBILITY STATEMENT

To ensure reproducibility of this work, We provide the experiment code in the supplementary materials, which can be executed on a medium-powerful machine. We provide well-written README and configuration files in order to reproduce our results. We also discuss the experiment environment in detail in Appendix C.1. We use benchmark datasets that are available to the public. The experiment environments, including the details of the machine we used, are discussed in Appendix C.2. We explicitly stated the amount of memory and time needed for execution.

## REFERENCES

- [1] Garima Agrawal, Tharindu Kumarage, Zeyad Alghami, and Huan Liu. Can knowledge graphs reduce hallucinations in llms? : A survey. *CoRR*, abs/2311.07914, 2023. doi: 10.48550/ARXIV.2311.07914. URL <https://doi.org/10.48550/arXiv.2311.07914>. 1
- [2] Muhammad Awais, Muzammal Naseer, Salman H. Khan, Rao Muhammad Anwer, Hisham Cholakkal, Mubarak Shah, Ming-Hsuan Yang, and Fahad Shahbaz Khan. Foundational models defining a new era in vision: A survey and outlook. *CoRR*, abs/2307.13721, 2023. doi: 10.48550/ARXIV.2307.13721. URL <https://doi.org/10.48550/arXiv.2307.13721>. 1
- [3] Baolong Bi, Shenghua Liu, Yiwei Wang, Lingrui Mei, and Xueqi Chen. Scalable link prediction on large-scale heterogeneous graphs with large language models. *CoRR*, abs/2401.13227, 2024. doi: 10.48550/ARXIV.2401.13227. URL <https://doi.org/10.48550/arXiv.2401.13227>. 1
- [4] Rishi Bommasani, Drew A. Hudson, Ehsan Adeli, Russ B. Altman, Simran Arora, Sydney von Arx, Michael S. Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, Erik Brynjolfsson, Shyamal Buch, Dallas Card, Rodrigo Castellon, Niladri S. Chatterji, Annie S. Chen, Kathleen Creel, Jared Quincy Davis, Dorottya Demszky, Chris Donahue, Moussa Doumbouya, Esin Durmus, Stefano Ermon, John Etchemendy, Kawin Ethayarajh, Li Fei-Fei, Chelsea Finn, Trevor Gale, Lauren Gillespie, Karan Goel, Noah D. Goodman, Shelby Grossman, Neel Guha, Tatsunori Hashimoto, Peter Henderson, John Hewitt, Daniel E. Ho, Jenny Hong, Kyle Hsu, Jing Huang, Thomas Icard, Saahil Jain, Dan Jurafsky, Pratyusha Kalluri, Siddharth Karamcheti, Geoff Keeling, Fereshte Khani, Omar Khattab, Pang Wei Koh, Mark S. Krass, Ranjay Krishna, Rohith Kuditipudi, and et al. On the opportunities and risks of foundation models. *CoRR*, abs/2108.07258, 2021. URL <https://arxiv.org/abs/2108.07258>. 1
- [5] Karsten M. Borgwardt, Cheng Soon Ong, Stefan Schönauer, S. V. N. Vishwanathan, Alexander J. Smola, and Hans-Peter Kriegel. Protein function prediction via graph kernels. In *Proceedings Thirteenth International Conference on Intelligent Systems for Molecular Biology 2005, Detroit, MI, USA, 25-29 June 2005*, pp. 47–56, 2005. doi: 10.1093/BIOINFORMATICS/BTI1007. URL <https://doi.org/10.1093/bioinformatics/bti1007>. 6
- [6] Digbalay Bose, Rajat Hebbar, Krishna Somandepalli, Haoyang Zhang, Yin Cui, Kree Cole-McLaughlin, Huisheng Wang, and Shrikanth Narayanan. Movieclip: Visual scene recognition in movies. In *IEEE/CVF Winter Conference on Applications of Computer Vision, WACV 2023, Waikoloa, HI, USA, January 2-7, 2023*, pp. 2082–2091. IEEE, 2023. doi: 10.1109/WACV56688.2023.00212. URL <https://doi.org/10.1109/WACV56688.2023.00212>. 2
- [7] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh,

- 
- 594 Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler,  
595 Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish,  
596 Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. *CoRR*,  
597 abs/2005.14165, 2020. URL <https://arxiv.org/abs/2005.14165>. 1  
598
- [8] Ziwei Chai, Tianjie Zhang, Liang Wu, Kaiqiao Han, Xiaohai Hu, Xuanwen Huang, and  
599 Yang Yang. Graphllm: Boosting graph reasoning ability of large language model. *CoRR*,  
600 abs/2310.05845, 2023. doi: 10.48550/ARXIV.2310.05845. URL [https://doi.org/10.  
601 48550/arXiv.2310.05845](https://doi.org/10.48550/arXiv.2310.05845). 1  
602
- [9] Mouxiang Chen, Zemin Liu, Chenghao Liu, Jundong Li, Qiheng Mao, and Jianling  
603 Sun. ULTRA-DP: unifying graph pre-training with multi-task graph dual prompt. *CoRR*,  
604 abs/2310.14845, 2023. doi: 10.48550/ARXIV.2310.14845. URL [https://doi.org/10.  
605 48550/arXiv.2310.14845](https://doi.org/10.48550/arXiv.2310.14845). 10, 24  
606
- [10] Zefeng Chen, Wensheng Gan, Jiayang Wu, Kaixia Hu, and Hong Lin. Data scarcity in  
607 recommendation systems: A survey. *CoRR*, abs/2312.10073, 2023. doi: 10.48550/ARXIV.  
608 2312.10073. URL <https://doi.org/10.48550/arXiv.2312.10073>. 1  
609
- [11] Eli Chien, Wei-Cheng Chang, Cho-Jui Hsieh, Hsiang-Fu Yu, Jiong Zhang, Olgica Milenkovic,  
610 and Inderjit S Dhillon. Node feature extraction by self-supervised multi-scale neighborhood  
611 prediction. *arXiv preprint arXiv:2111.00064*, 2021. 10  
612
- [12] Kevin Clark, Minh-Thang Luong, Quoc V. Le, and Christopher D. Manning. ELECTRA:  
613 pre-training text encoders as discriminators rather than generators. In *8th International Confer-  
614 ence on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*.  
OpenReview.net, 2020. URL <https://openreview.net/forum?id=r1xMH1BtvB>.  
615 6, 9, 25  
616
- [13] Enyan Dai, Tianxiang Zhao, Huaisheng Zhu, Junjie Xu, Zhimeng Guo, Hui Liu, Jiliang Tang,  
617 and Suhang Wang. A comprehensive survey on trustworthy graph neural networks: Privacy,  
618 robustness, fairness, and explainability. *CoRR*, abs/2204.08570, 2022. doi: 10.48550/ARXIV.  
619 2204.08570. URL <https://doi.org/10.48550/arXiv.2204.08570>. 1  
620
- [14] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: pre-training of  
621 deep bidirectional transformers for language understanding. In Jill Burstein, Christy Doran,  
622 and Tamar Solorio (eds.), *Proceedings of the 2019 Conference of the North American Chapter  
623 of the Association for Computational Linguistics: Human Language Technologies, NAACL-  
624 HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)*, pp.  
4171–4186. Association for Computational Linguistics, 2019. doi: 10.18653/V1/N19-1423.  
625 URL <https://doi.org/10.18653/v1/n19-1423>. 1  
626
- [15] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai,  
627 Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly,  
628 Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for  
629 image recognition at scale. In *9th International Conference on Learning Representations,  
630 ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net, 2021. URL <https://openreview.net/forum?id=YicbFdNTTy>. 1  
631
- [16] Keyu Duan, Qian Liu, Tat-Seng Chua, Shuicheng Yan, Wei Tsang Ooi, Qizhe Xie, and Junxian  
632 He. Simteg: A frustratingly simple approach improves textual graph learning. *arXiv preprint  
633 arXiv:2308.02565*, 2023. 10  
634
- [17] Vijay Prakash Dwivedi, Chaitanya K. Joshi, Thomas Laurent, Yoshua Bengio, and Xavier  
635 Bresson. Benchmarking graph neural networks. *CoRR*, abs/2003.00982, 2020. URL <https://arxiv.org/abs/2003.00982>. 1  
636
- [18] Vijay Prakash Dwivedi, Chaitanya K. Joshi, Anh Tuan Luu, Thomas Laurent, Yoshua Bengio,  
637 and Xavier Bresson. Benchmarking graph neural networks. *J. Mach. Learn. Res.*, 24:43:1–  
638 43:48, 2023. URL <http://jmlr.org/papers/v24/22-0567.html>. 3  
639  
640  
641  
642  
643  
644  
645  
646  
647

- 
- 648 [19] Taoran Fang, Yunchao Zhang, Yang Yang, Chunping Wang, and Lei Chen. Universal  
649 prompt tuning for graph neural networks. In Alice Oh, Tristan Naumann, Amir  
650 Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine (eds.), *Advances in Neural  
651 Information Processing Systems 36: Annual Conference on Neural Information Pro-  
652 cessing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16,  
653 2023*, 2023. URL [http://papers.nips.cc/paper\\_files/paper/2023/hash/  
654 a4a1ee071ce0fe63b83bce507c9dc4d7-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2023/hash/a4a1ee071ce0fe63b83bce507c9dc4d7-Abstract-Conference.html). 6, 7,  
655 10, 24
- 656 [20] Shengyu Feng, Baoyu Jing, Yada Zhu, and Hanghang Tong. Adversarial graph contrastive  
657 learning with information regularization. In *WWW '22: The ACM Web Conference 2022,  
658 Virtual Event, Lyon, France, April 25 - 29, 2022*, pp. 1362–1371. ACM, 2022. 10
- 659 [21] Matthias Fey and Jan Eric Lenssen. Fast graph representation learning with pytorch geometric.  
660 *CoRR*, abs/1903.02428, 2019. URL <http://arxiv.org/abs/1903.02428>. 3, 6
- 661 [22] Qingqing Ge, Zeyuan Zhao, Yiding Liu, Anfeng Cheng, Xiang Li, Shuaiqiang Wang,  
662 and Dawei Yin. Enhancing graph neural networks with structure-based prompt. *CoRR*,  
663 abs/2310.17394, 2023. doi: 10.48550/ARXIV.2310.17394. URL [https://doi.org/10.  
664 48550/arXiv.2310.17394](https://doi.org/10.48550/arXiv.2310.17394). 10, 24
- 665 [23] Andrey Guzhov, Federico Raue, Jörn Hees, and Andreas Dengel. Audioclip: Extending  
666 clip to image, text and audio. In *IEEE International Conference on Acoustics, Speech and  
667 Signal Processing, ICASSP 2022, Virtual and Singapore, 23-27 May 2022*, pp. 976–980. IEEE,  
668 2022. doi: 10.1109/ICASSP43922.2022.9747631. URL [https://doi.org/10.1109/  
669 ICASSP43922.2022.9747631](https://doi.org/10.1109/ICASSP43922.2022.9747631). 2
- 670 [24] Kaveh Hassani and Amir Hosein Khas Ahmadi. Contrastive multi-view representation learning  
671 on graphs. In *Proceedings of the 37th International Conference on Machine Learning, ICML  
672 2020, 13-18 July 2020, Virtual Event*, volume 119 of *Proceedings of Machine Learning  
673 Research*, pp. 4116–4126. PMLR, 2020. URL [http://proceedings.mlr.press/  
674 v119/hassani20a.html](http://proceedings.mlr.press/v119/hassani20a.html). 7, 9, 25, 26
- 675 [25] Georg Hess, Adam Tonderski, Christoffer Petersson, Kalle Åström, and Lennart Svensson.  
676 Lidarclip or: How I learned to talk to point clouds. In *IEEE/CVF Winter Conference on  
677 Applications of Computer Vision, WACV 2024, Waikoloa, HI, USA, January 3-8, 2024*, pp.  
678 7423–7432. IEEE, 2024. doi: 10.1109/WACV57701.2024.00727. URL [https://doi.  
679 org/10.1109/WACV57701.2024.00727](https://doi.org/10.1109/WACV57701.2024.00727). 2
- 680 [26] Fangzhou Hong, Mingyuan Zhang, Liang Pan, Zhongang Cai, Lei Yang, and Ziwei Liu.  
681 Avatarclip: zero-shot text-driven generation and animation of 3d avatars. *ACM Trans. Graph.*,  
682 41(4):161:1–161:19, 2022. doi: 10.1145/3528223.3530094. URL [https://doi.org/10.  
683 1145/3528223.3530094](https://doi.org/10.1145/3528223.3530094). 2
- 684 [27] Zhenyu Hou, Xiao Liu, Yukuo Cen, Yuxiao Dong, Hongxia Yang, Chunjie Wang, and Jie  
685 Tang. Graphmae: Self-supervised masked graph autoencoders. In Aidong Zhang and Huzefa  
686 Rangwala (eds.), *KDD '22: The 28th ACM SIGKDD Conference on Knowledge Discovery  
687 and Data Mining, Washington, DC, USA, August 14 - 18, 2022*, pp. 594–604. ACM, 2022. doi:  
688 10.1145/3534678.3539321. URL [https://doi.org/10.1145/3534678.3539321.  
689 7, 9, 25](https://doi.org/10.1145/3534678.3539321)
- 690 [28] Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin de Laroussilhe,  
691 Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. Parameter-efficient transfer learning  
692 for NLP. In Kamalika Chaudhuri and Ruslan Salakhutdinov (eds.), *Proceedings of the 36th  
693 International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach,  
694 California, USA*, volume 97 of *Proceedings of Machine Learning Research*, pp. 2790–2799.  
695 PMLR, 2019. URL [http://proceedings.mlr.press/v97/houlsby19a.html.  
696 2](http://proceedings.mlr.press/v97/houlsby19a.html)
- 697 [29] Weihua Hu, Bowen Liu, Joseph Gomes, Marinka Zitnik, Percy Liang, Vijay S. Pande, and Jure  
698 Leskovec. Strategies for pre-training graph neural networks. In *8th International Conference  
699 700 701*

- 
- 702           on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020. Open-  
703           Review.net, 2020. URL <https://openreview.net/forum?id=HJ1WWJSFDH>. 1, 2,  
704           10
- 705
- 706 [30] Ziniu Hu, Yuxiao Dong, Kuansan Wang, Kai-Wei Chang, and Yizhou Sun. GPT-GNN:  
707           generative pre-training of graph neural networks. In Rajesh Gupta, Yan Liu, Jiliang Tang,  
708           and B. Aditya Prakash (eds.), *KDD '20: The 26th ACM SIGKDD Conference on Knowledge*  
709           *Discovery and Data Mining, Virtual Event, CA, USA, August 23-27, 2020*, pp. 1857–1867.  
710           ACM, 2020. 10
- 711
- 712 [31] Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang,  
713           Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, and Ting Liu. A survey  
714           on hallucination in large language models: Principles, taxonomy, challenges, and open  
715           questions. *CoRR*, abs/2311.05232, 2023. doi: 10.48550/ARXIV.2311.05232. URL  
716           <https://doi.org/10.48550/arXiv.2311.05232>. 1
- 717
- 718 [32] Qian Huang, Hongyu Ren, Peng Chen, Gregor Krzmann, Daniel Zeng, Percy Liang, and  
719           Jure Leskovec. PRODIGY: enabling in-context learning over graphs. In Alice Oh, Tris-  
720           tan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine (eds.), *Ad-*  
721           *vancess in Neural Information Processing Systems 36: Annual Conference on Neural Infor-*  
722           *mation Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16,*  
723           2023. URL [http://papers.nips.cc/paper\\_files/paper/2023/hash/](http://papers.nips.cc/paper_files/paper/2023/hash/34dce0dc3121951dd0399ba02c0f0d06-Abstract-Conference.html)  
724           34dce0dc3121951dd0399ba02c0f0d06-Abstract-Conference.html. 10,  
725           24
- 726
- 727 [33] Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc V. Le, Yun-  
728           Hsuan Sung, Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation  
729           learning with noisy text supervision. In Marina Meila and Tong Zhang (eds.), *Proceedings of*  
730           *the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual*  
731           *Event*, volume 139 of *Proceedings of Machine Learning Research*, pp. 4904–4916. PMLR,  
732           2021. URL <http://proceedings.mlr.press/v139/jia21b.html>. 1
- 733
- 734 [34] Wenzhao Jiang, Hao Liu, and Hui Xiong. Survey on trustworthy graph neural networks: From  
735           A causal perspective. *CoRR*, abs/2312.12477, 2023. doi: 10.48550/ARXIV.2312.12477. URL  
736           <https://doi.org/10.48550/arXiv.2312.12477>. 1
- 737
- 738 [35] Xunqiang Jiang, Yuanfu Lu, Yuan Fang, and Chuan Shi. Contrastive pre-training of gnns on  
739           heterogeneous graphs. In *CIKM '21: The 30th ACM International Conference on Information*  
740           *and Knowledge Management, Virtual Event, Queensland, Australia, November 1 - 5, 2021*, pp.  
741           803–812. ACM, 2021. 10
- 742
- 743 [36] Bowen Jin, Gang Liu, Chi Han, Meng Jiang, Heng Ji, and Jiawei Han. Large language models  
744           on graphs: A comprehensive survey. *arXiv preprint arXiv:2312.02783*, 2023. 10
- 745
- 746 [37] Bowen Jin, Wentao Zhang, Yu Zhang, Yu Meng, Xinyang Zhang, Qi Zhu, and Jiawei Han.  
747           Patton: Language model pretraining on text-rich networks. *arXiv preprint arXiv:2305.12268*,  
748           2023. 10
- 749
- 750 [38] Bowen Jin, Yu Zhang, Qi Zhu, and Jiawei Han. Heterformer: Transformer-based deep node  
751           representation learning on heterogeneous text-rich networks. In *Proceedings of the 29th ACM*  
752           *SIGKDD Conference on Knowledge Discovery and Data Mining*, pp. 1020–1031, 2023. 10
- 753
- 754 [39] Baoyu Jing, Chanyoung Park, and Hanghang Tong. HDMI: high-order deep multiplex infomax.  
755           In *WWW '21: The Web Conference 2021, Virtual Event / Ljubljana, Slovenia, April 19-23,*  
          2021, pp. 2414–2424. ACM / IW3C2, 2021. 10
- 756
- 757 [40] Wei Ju, Siyu Yi, Yifan Wang, Zhiping Xiao, Zhengyang Mao, Hourun Li, Yiyang Gu, Yifang  
758           Qin, Nan Yin, Senzhang Wang, Xinwang Liu, Xiao Luo, Philip S. Yu, and Ming Zhang. A  
759           survey of graph neural networks in real world: Imbalance, noise, privacy and OOD challenges.  
760           *CoRR*, abs/2403.04468, 2024. doi: 10.48550/ARXIV.2403.04468. URL [https://doi.](https://doi.org/10.48550/arXiv.2403.04468)  
761           [org/10.48550/arXiv.2403.04468](https://doi.org/10.48550/arXiv.2403.04468). 1

- 
- 756 [41] Muhammad Uzair Khattak, Hanoona Abdul Rasheed, Muhammad Maaz, Salman H. Khan,  
757 and Fahad Shahbaz Khan. Maple: Multi-modal prompt learning. In *IEEE/CVF Conference*  
758 *on Computer Vision and Pattern Recognition, CVPR 2023, Vancouver, BC, Canada, June*  
759 *17-24, 2023*, pp. 19113–19122. IEEE, 2023. doi: 10.1109/CVPR52729.2023.01832. URL  
760 <https://doi.org/10.1109/CVPR52729.2023.01832>. 2  
761
- 762 [42] Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola,  
763 Aaron Maschinot, Ce Liu, and Dilip Krishnan. Supervised contrastive learning. In *Advances*  
764 *in Neural Information Processing Systems 33: Annual Conference on Neural Information*  
765 *Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, 2020. 7, 10
- 766 [43] Thomas N. Kipf and Max Welling. Semi-supervised classification with graph convolutional  
767 networks. *CoRR*, abs/1609.02907, 2016. URL <http://arxiv.org/abs/1609.02907>.  
768 7  
769
- 770 [44] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloé Rolland, Laura Gustafson,  
771 Tete Xiao, Spencer Whitehead, Alexander C. Berg, Wan-Yen Lo, Piotr Dollár, and Ross B.  
772 Girshick. Segment anything. *CoRR*, abs/2304.02643, 2023. doi: 10.48550/ARXIV.2304.02643.  
773 URL <https://doi.org/10.48550/arXiv.2304.02643>. 1
- 774 [45] Brian Lester, Rami Al-Rfou, and Noah Constant. The power of scale for parameter-efficient  
775 prompt tuning. In Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-  
776 tau Yih (eds.), *Proceedings of the 2021 Conference on Empirical Methods in Natural*  
777 *Language Processing, EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic,*  
778 *7-11 November, 2021*, pp. 3045–3059. Association for Computational Linguistics, 2021.  
779 doi: 10.18653/V1/2021.EMNLP-MAIN.243. URL [https://doi.org/10.18653/v1/](https://doi.org/10.18653/v1/2021.emnlp-main.243)  
780 [2021.emnlp-main.243](https://doi.org/10.18653/v1/2021.emnlp-main.243). 2  
781
- 782 [46] Bolian Li, Baoyu Jing, and Hanghang Tong. Graph communal contrastive learning. In *WWW*  
783 *'22: The ACM Web Conference 2022, Virtual Event, Lyon, France, April 25 - 29, 2022*, pp.  
784 1203–1213. ACM, 2022. 10
- 785 [47] Junnan Li, Dongxu Li, Caiming Xiong, and Steven C. H. Hoi. BLIP: bootstrapping language-  
786 image pre-training for unified vision-language understanding and generation. In Kamalika  
787 Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvári, Gang Niu, and Sivan Sabato (eds.),  
788 *International Conference on Machine Learning, ICML 2022, 17-23 July 2022, Baltimore,*  
789 *Maryland, USA, volume 162 of Proceedings of Machine Learning Research*, pp. 12888–12900.  
790 PMLR, 2022. URL <https://proceedings.mlr.press/v162/li22n.html>. 1  
791
- 792 [48] Xiang Lisa Li and Percy Liang. Prefix-tuning: Optimizing continuous prompts for generation.  
793 In Chengqing Zong, Fei Xia, Wenjie Li, and Roberto Navigli (eds.), *Proceedings of the 59th*  
794 *Annual Meeting of the Association for Computational Linguistics and the 11th International*  
795 *Joint Conference on Natural Language Processing, ACL/IJCNLP 2021, (Volume 1: Long*  
796 *Papers), Virtual Event, August 1-6, 2021*, pp. 4582–4597. Association for Computational  
797 Linguistics, 2021. doi: 10.18653/V1/2021.ACL-LONG.353. URL [https://doi.org/](https://doi.org/10.18653/v1/2021.acl-long.353)  
798 [10.18653/v1/2021.acl-long.353](https://doi.org/10.18653/v1/2021.acl-long.353). 2
- 799 [49] Yuxuan Liang, Haomin Wen, Yuqi Nie, Yushan Jiang, Ming Jin, Dongjin Song, Shirui Pan,  
800 and Qingsong Wen. Foundation models for time series analysis: A tutorial and survey. *CoRR*,  
801 abs/2403.14735, 2024. doi: 10.48550/ARXIV.2403.14735. URL [https://doi.org/10.](https://doi.org/10.48550/arXiv.2403.14735)  
802 [48550/arXiv.2403.14735](https://doi.org/10.48550/arXiv.2403.14735). 1
- 803 [50] Jiawei Liu, Cheng Yang, Zhiyuan Lu, Junze Chen, Yibo Li, Mengmei Zhang, Ting Bai,  
804 Yuan Fang, Lichao Sun, Philip S. Yu, and Chuan Shi. Towards graph foundation models: A  
805 survey and beyond. *CoRR*, abs/2310.11829, 2023. doi: 10.48550/ARXIV.2310.11829. URL  
806 <https://doi.org/10.48550/arXiv.2310.11829>. 1  
807
- 808 [51] Pengfei Liu, Yiming Ren, and Zhixiang Ren. Git-mol: A multi-modal large language model for  
809 molecular science with graph, image, and text. *CoRR*, abs/2308.06911, 2023. doi: 10.48550/  
ARXIV.2308.06911. URL <https://doi.org/10.48550/arXiv.2308.06911>. 2

- 
- 810 [52] Yang Liu, Jiahuan Cao, Chongyu Liu, Kai Ding, and Lianwen Jin. Datasets for large language  
811 models: A comprehensive survey. *CoRR*, abs/2402.18041, 2024. doi: 10.48550/ARXIV.2402.  
812 18041. URL <https://doi.org/10.48550/arXiv.2402.18041>. 1
- 813 [53] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy,  
814 Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized BERT  
815 pretraining approach. *CoRR*, abs/1907.11692, 2019. URL [http://arxiv.org/abs/  
816 1907.11692](http://arxiv.org/abs/1907.11692). 6
- 817 [54] Yixin Liu, Ming Jin, Shirui Pan, Chuan Zhou, Yu Zheng, Feng Xia, and Philip S. Yu. Graph  
818 self-supervised learning: A survey. *IEEE Trans. Knowl. Data Eng.*, 35(6):5879–5900, 2023.  
819 doi: 10.1109/TKDE.2022.3172903. URL [https://doi.org/10.1109/TKDE.2022.  
820 3172903](https://doi.org/10.1109/TKDE.2022.3172903). 1
- 821 [55] Zemin Liu, Xingtong Yu, Yuan Fang, and Xinming Zhang. Graphprompt: Unifying pre-training  
822 and downstream tasks for graph neural networks. In Ying Ding, Jie Tang, Juan F. Sequeda, Lora  
823 Aroyo, Carlos Castillo, and Geert-Jan Houben (eds.), *Proceedings of the ACM Web Conference  
824 2023, WWW 2023, Austin, TX, USA, 30 April 2023 - 4 May 2023*, pp. 417–428. ACM, 2023. doi:  
825 10.1145/3543507.3583386. URL <https://doi.org/10.1145/3543507.3583386>.  
826 6, 7, 24
- 827 [56] Zhiyuan Liu, Sihang Li, Yanchen Luo, Hao Fei, Yixin Cao, Kenji Kawaguchi, Xiang Wang,  
828 and Tat-Seng Chua. Molca: Molecular graph-language modeling with cross-modal projector  
829 and uni-modal adapter. In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), *Proceedings  
830 of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP  
831 2023, Singapore, December 6-10, 2023*, pp. 15623–15638. Association for Computational  
832 Linguistics, 2023. doi: 10.18653/v1/2023.EMNLP-MAIN.966. URL [https://doi.org/  
833 10.18653/v1/2023.emnlp-main.966](https://doi.org/10.18653/v1/2023.emnlp-main.966). 2
- 834 [57] Yuanfu Lu, Xunqiang Jiang, Yuan Fang, and Chuan Shi. Learning to pre-train graph neural  
835 networks. In *Thirty-Fifth AAAI Conference on Artificial Intelligence, AAAI 2021, Thirty-Third  
836 Conference on Innovative Applications of Artificial Intelligence, IAAI 2021, The Eleventh  
837 Symposium on Educational Advances in Artificial Intelligence, EAAI 2021, Virtual Event,  
838 February 2-9, 2021*, pp. 4276–4284. AAAI Press, 2021. 10
- 839 [58] Yizhen Luo, Kai Yang, Massimo Hong, Xing Yi Liu, and Zaiqing Nie. Molfm: A multimodal  
840 molecular foundation model. *CoRR*, abs/2307.09484, 2023. doi: 10.48550/ARXIV.2307.09484.  
841 URL <https://doi.org/10.48550/arXiv.2307.09484>. 2
- 842 [59] Sahil Manchanda, Shubham Gupta, Sayan Ranu, and Srikanta J. Bedathur. Generative  
843 modeling of labeled graphs under data scarcity. In Soledad Villar and Benjamin Cham-  
844 berlain (eds.), *Learning on Graphs Conference, 27-30 November 2023, Virtual Event*, vol-  
845 ume 231 of *Proceedings of Machine Learning Research*, pp. 32. PMLR, 2023. URL  
846 <https://proceedings.mlr.press/v231/manchanda24a.html>. 1
- 847 [60] Christopher Morris, Nils M. Kriege, Franka Bause, Kristian Kersting, Petra Mutzel, and  
848 Marion Neumann. Tudataset: A collection of benchmark datasets for learning with graphs.  
849 *CoRR*, abs/2007.08663, 2020. URL <https://arxiv.org/abs/2007.08663>. 3, 6
- 850 [61] Christopher Morris, Yaron Lipman, Haggai Maron, Bastian Rieck, Nils M. Kriege, Martin  
851 Grohe, Matthias Fey, and Karsten M. Borgwardt. Weisfeiler and leman go machine learning:  
852 The story so far. *CoRR*, abs/2112.09992, 2021. URL [https://arxiv.org/abs/2112.  
853 09992](https://arxiv.org/abs/2112.09992). 1
- 854 [62] Nicolò Navarin, Dinh Van Tran, and Alessandro Sperduti. Pre-training graph neural networks  
855 with kernels. *CoRR*, abs/1811.06930, 2018. 10
- 856 [63] Marion Neumann, Roman Garnett, Christian Bauckhage, and Kristian Kersting. Propagation  
857 kernels: efficient graph kernels from propagated information. *Mach. Learn.*, 102(2):209–  
858 245, 2016. doi: 10.1007/S10994-015-5517-9. URL [https://doi.org/10.1007/  
859 s10994-015-5517-9](https://doi.org/10.1007/s10994-015-5517-9). 6
- 860  
861  
862  
863

- 864 [64] Jiezhong Qiu, Qibin Chen, Yuxiao Dong, Jing Zhang, Hongxia Yang, Ming Ding, Kuansan  
865 Wang, and Jie Tang. GCC: graph contrastive coding for graph neural network pre-training. In  
866 Rajesh Gupta, Yan Liu, Jiliang Tang, and B. Aditya Prakash (eds.), *KDD '20: The 26th ACM*  
867 *SIGKDD Conference on Knowledge Discovery and Data Mining, Virtual Event, CA, USA,*  
868 *August 23-27, 2020*, pp. 1150–1160. ACM, 2020. 10
- 869 [65] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini  
870 Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and  
871 Ilya Sutskever. Learning transferable visual models from natural language supervision. In  
872 Marina Meila and Tong Zhang (eds.), *Proceedings of the 38th International Conference on*  
873 *Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, volume 139 of *Proceedings of*  
874 *Machine Learning Research*, pp. 8748–8763. PMLR, 2021. URL [http://proceedings.](http://proceedings.mlr.press/v139/radford21a.html)  
875 [mlr.press/v139/radford21a.html](http://proceedings.mlr.press/v139/radford21a.html). 1
- 876 [66] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark  
877 Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In Marina Meila and Tong  
878 Zhang (eds.), *Proceedings of the 38th International Conference on Machine Learning, ICML*  
879 *2021, 18-24 July 2021, Virtual Event*, volume 139 of *Proceedings of Machine Learning*  
880 *Research*, pp. 8821–8831. PMLR, 2021. URL [http://proceedings.mlr.press/](http://proceedings.mlr.press/v139/ramesh21a.html)  
881 [v139/ramesh21a.html](http://proceedings.mlr.press/v139/ramesh21a.html). 1
- 882 [67] Ladislav Rampásek, Michael Galkin, Vijay Prakash Dwivedi, Anh Tuan Luu, Guy Wolf, and  
883 Dominique Beaini. Recipe for a general, powerful, scalable graph transformer. In Sanmi  
884 Koyejo, S. Mohamed, A. Agarwal, Danielle Belgrave, K. Cho, and A. Oh (eds.), *Advances*  
885 *in Neural Information Processing Systems 35: Annual Conference on Neural Information*  
886 *Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9,*  
887 *2022*, 2022. URL [http://papers.nips.cc/paper\\_files/paper/2022/hash/](http://papers.nips.cc/paper_files/paper/2022/hash/5d4834a159f1547b267a05a4e2b7cf5e-Abstract-Conference.html)  
888 [5d4834a159f1547b267a05a4e2b7cf5e-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2022/hash/5d4834a159f1547b267a05a4e2b7cf5e-Abstract-Conference.html). 1
- 889 [68] Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. Distilbert, a distilled  
890 version of BERT: smaller, faster, cheaper and lighter. *CoRR*, abs/1910.01108, 2019. URL  
891 <http://arxiv.org/abs/1910.01108>. 6, 9, 25
- 892 [69] Ryoma Sato. A survey on the expressive power of graph neural networks. *CoRR*,  
893 abs/2003.04078, 2020. URL <https://arxiv.org/abs/2003.04078>. 1
- 894 [70] Yi-Jen Shih, Hsuan-Fu Wang, Heng-Jui Chang, Layne Berry, Hung-yi Lee, and David Harwath.  
895 Speechclip: Integrating speech with pre-trained vision and language model. In *IEEE Spoken*  
896 *Language Technology Workshop, SLT 2022, Doha, Qatar, January 9-12, 2023*, pp. 715–722.  
897 IEEE, 2022. doi: 10.1109/SLT54892.2023.10022954. URL [https://doi.org/10.](https://doi.org/10.1109/SLT54892.2023.10022954)  
898 [1109/SLT54892.2023.10022954](https://doi.org/10.1109/SLT54892.2023.10022954). 2
- 899 [71] Harry Shomer, Yao Ma, Haitao Mao, Juanhui Li, Bo Wu, and Jiliang Tang. Lpformer: An  
900 adaptive graph transformer for link prediction, 2024. 1
- 901 [72] Bing Su, Dazhao Du, Zhao Yang, Yujie Zhou, Jiangmeng Li, Anyi Rao, Hao Sun, Zhiwu Lu,  
902 and Ji-Rong Wen. A molecular multimodal foundation model associating molecule graphs  
903 with natural language. *CoRR*, abs/2209.05481, 2022. doi: 10.48550/ARXIV.2209.05481.  
904 URL <https://doi.org/10.48550/arXiv.2209.05481>. 2, 10
- 905 [73] Fan-Yun Sun, Jordan Hoffmann, Vikas Verma, and Jian Tang. Infograph: Unsupervised and  
906 semi-supervised graph-level representation learning via mutual information maximization. In  
907 *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia,*  
908 *April 26-30, 2020*. OpenReview.net, 2020. 10
- 909 [74] Lichao Sun, Yue Huang, Haoran Wang, Siyuan Wu, Qihui Zhang, Chujie Gao, Yixin Huang,  
910 Wenhan Lyu, Yixuan Zhang, Xiner Li, Zhengliang Liu, Yixin Liu, Yijue Wang, Zhikun Zhang,  
911 Bhavya Kailkhura, Caiming Xiong, Chaowei Xiao, Chunyuan Li, Eric P. Xing, Furong Huang,  
912 Hao Liu, Heng Ji, Hongyi Wang, Huan Zhang, Huaxiu Yao, Manolis Kellis, Marinka Zitnik,  
913 Meng Jiang, Mohit Bansal, James Zou, Jian Pei, Jian Liu, Jianfeng Gao, Jiawei Han, Jieyu  
914 Zhao, Jiliang Tang, Jindong Wang, John Mitchell, Kai Shu, Kaidi Xu, Kai-Wei Chang, Lifang  
915 He, Lifu Huang, Michael Backes, Neil Zhenqiang Gong, Philip S. Yu, Pin-Yu Chen, Quanquan

- 918 Gu, Ran Xu, Rex Ying, Shuiwang Ji, Suman Jana, Tianlong Chen, Tianming Liu, Tianyi  
919 Zhou, William Wang, Xiang Li, Xiangliang Zhang, Xiao Wang, Xing Xie, Xun Chen, Xuyu  
920 Wang, Yan Liu, Yanfang Ye, Yinzhi Cao, and Yue Zhao. Trustllm: Trustworthiness in large  
921 language models. *CoRR*, abs/2401.05561, 2024. doi: 10.48550/ARXIV.2401.05561. URL  
922 <https://doi.org/10.48550/arXiv.2401.05561>. 1
- 923  
924 [75] Mingchen Sun, Kaixiong Zhou, Xin He, Ying Wang, and Xin Wang. GPPT: graph pre-  
925 training and prompt tuning to generalize graph neural networks. In Aidong Zhang and Huzefa  
926 Rangwala (eds.), *KDD '22: The 28th ACM SIGKDD Conference on Knowledge Discovery and*  
927 *Data Mining, Washington, DC, USA, August 14 - 18, 2022*, pp. 1717–1727. ACM, 2022. doi:  
928 10.1145/3534678.3539249. URL <https://doi.org/10.1145/3534678.3539249>.  
929 7, 10, 24
- 930 [76] Xiangguo Sun, Hong Cheng, Jia Li, Bo Liu, and Jihong Guan. All in one: Multi-task prompting  
931 for graph neural networks. In Ambuj K. Singh, Yizhou Sun, Leman Akoglu, Dimitrios  
932 Gunopulos, Xifeng Yan, Ravi Kumar, Fatma Ozcan, and Jieping Ye (eds.), *Proceedings of the*  
933 *29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, KDD 2023, Long*  
934 *Beach, CA, USA, August 6-10, 2023*, pp. 2120–2131. ACM, 2023. doi: 10.1145/3580305.  
935 3599256. URL <https://doi.org/10.1145/3580305.3599256>. 2, 3, 4, 6, 7, 9, 10,  
936 24
- 937 [77] Xiangguo Sun, Jiawen Zhang, Xixi Wu, Hong Cheng, Yun Xiong, and Jia Li. Graph prompt  
938 learning: A comprehensive survey and beyond. *CoRR*, abs/2311.16534, 2023. doi: 10.48550/  
939 ARXIV.2311.16534. URL <https://doi.org/10.48550/arXiv.2311.16534>. 10
- 940 [78] Zhen Tan, Ruocheng Guo, Kaize Ding, and Huan Liu. Virtual node tuning for few-shot  
941 node classification. In Ambuj K. Singh, Yizhou Sun, Leman Akoglu, Dimitrios Gunopulos,  
942 Xifeng Yan, Ravi Kumar, Fatma Ozcan, and Jieping Ye (eds.), *Proceedings of the 29th ACM*  
943 *SIGKDD Conference on Knowledge Discovery and Data Mining, KDD 2023, Long Beach, CA,*  
944 *USA, August 6-10, 2023*, pp. 2177–2188. ACM, 2023. doi: 10.1145/3580305.3599541. URL  
945 <https://doi.org/10.1145/3580305.3599541>. 10, 24
- 946 [79] Ross Taylor, Marcin Kardas, Guillem Cucurull, Thomas Scialom, Anthony Hartshorn, Elvis  
947 Saravia, Andrew Poulton, Viktor Kerkez, and Robert Stojnic. Galactica: A large language  
948 model for science. *CoRR*, abs/2211.09085, 2022. doi: 10.48550/ARXIV.2211.09085. URL  
949 <https://doi.org/10.48550/arXiv.2211.09085>. 10
- 950 [80] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timo-  
951 thée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurélien Rodriguez,  
952 Armand Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation  
953 language models. *CoRR*, abs/2302.13971, 2023. doi: 10.48550/ARXIV.2302.13971. URL  
954 <https://doi.org/10.48550/arXiv.2302.13971>. 1
- 955 [81] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez,  
956 Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In Isabelle Guyon, Ulrike von  
957 Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman  
958 Garnett (eds.), *Advances in Neural Information Processing Systems 30: Annual Conference on*  
959 *Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*,  
960 pp. 5998–6008, 2017. URL [https://proceedings.neurips.cc/paper/2017/  
961 hash/3f5ee243547dee91fbd053c1c4a845aa-Abstract.html](https://proceedings.neurips.cc/paper/2017/hash/3f5ee243547dee91fbd053c1c4a845aa-Abstract.html). 4
- 962 [82] Heng Wang, Shangbin Feng, Tianxing He, Zhaoxuan Tan, Xiaochuang Han, and Yulia  
963 Tsvetkov. Can language models solve graph problems in natural language? In Alice Oh,  
964 Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine (eds.),  
965 *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Infor-*  
966 *mation Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16,*  
967 *2023*, 2023. URL [http://papers.nips.cc/paper\\_files/paper/2023/hash/  
968 622afc4edf2824a1b6aaf5afe153fa93-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2023/hash/622afc4edf2824a1b6aaf5afe153fa93-Abstract-Conference.html). 1
- 969 [83] Yi Wang, Kunchang Li, Xinhao Li, Jiashuo Yu, Yinan He, Guo Chen, Baoqi Pei, Rongkun  
970 Zheng, Jilan Xu, Zun Wang, Yansong Shi, Tianxiang Jiang, Songze Li, Hongjie Zhang, Yifei  
971 Huang, Yu Qiao, Yali Wang, and Limin Wang. Internvideo2: Scaling video foundation models

- 
- 972 for multimodal video understanding. *CoRR*, abs/2403.15377, 2024. doi: 10.48550/ARXIV.  
973 2403.15377. URL <https://doi.org/10.48550/arXiv.2403.15377>. 1  
974
- [84] Bowen Wen, Wei Yang, Jan Kautz, and Stan Birchfield. Foundationpose: Unified 6d pose  
975 estimation and tracking of novel objects. *CoRR*, abs/2312.08344, 2023. doi: 10.48550/ARXIV.  
976 2312.08344. URL <https://doi.org/10.48550/arXiv.2312.08344>. 1  
977
- [85] Ho-Hsiang Wu, Prem Seetharaman, Kundan Kumar, and Juan Pablo Bello. Wav2clip: Learning  
978 robust audio representations from clip. In *IEEE International Conference on Acoustics, Speech*  
979 *and Signal Processing, ICASSP 2022, Virtual and Singapore, 23-27 May 2022*, pp. 4563–4567.  
980 IEEE, 2022. doi: 10.1109/ICASSP43922.2022.9747669. URL [https://doi.org/10.](https://doi.org/10.1109/ICASSP43922.2022.9747669)  
981 [1109/ICASSP43922.2022.9747669](https://doi.org/10.1109/ICASSP43922.2022.9747669). 2  
982
- [86] Qitian Wu, Wentao Zhao, Zenan Li, David P. Wipf, and Junchi Yan. Nodeformer: A scal-  
983 able graph structure learning transformer for node classification. In Sanmi Koyejo, S. Mo-  
984 hamed, A. Agarwal, Danielle Belgrave, K. Cho, and A. Oh (eds.), *Advances in Neural*  
985 *Information Processing Systems 35: Annual Conference on Neural Information Process-*  
986 *ing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9,*  
987 *2022*, 2022. URL [http://papers.nips.cc/paper\\_files/paper/2022/hash/](http://papers.nips.cc/paper_files/paper/2022/hash/af790b7ae573771689438bbcf5933fe-Abstract-Conference.html)  
988 [af790b7ae573771689438bbcf5933fe-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2022/hash/af790b7ae573771689438bbcf5933fe-Abstract-Conference.html). 1  
989
- [87] Xuansheng Wu, Kaixiong Zhou, Mingchen Sun, Xin Wang, and Ninghao Liu. A survey of  
990 graph prompting methods: Techniques, applications, and challenges. *CoRR*, abs/2303.07275,  
991 2023. doi: 10.48550/ARXIV.2303.07275. URL [https://doi.org/10.48550/arXiv.](https://doi.org/10.48550/arXiv.2303.07275)  
992 [2303.07275](https://doi.org/10.48550/arXiv.2303.07275). 10  
993
- [88] Zhenqin Wu, Bharath Ramsundar, Evan N. Feinberg, Joseph Gomes, Caleb Geniesse, Aneesh S.  
994 Pappu, Karl Leswing, and Vijay S. Pande. Moleculenet: A benchmark for molecular machine  
995 learning. *CoRR*, abs/1703.00564, 2017. URL <http://arxiv.org/abs/1703.00564>.  
996 10  
997
- [89] Zonghan Wu, Shirui Pan, Fengwen Chen, Guodong Long, Chengqi Zhang, and Philip S. Yu.  
998 A comprehensive survey on graph neural networks. *CoRR*, abs/1901.00596, 2019. URL  
999 <http://arxiv.org/abs/1901.00596>. 1  
1000
- [90] Zonghan Wu, Shirui Pan, Fengwen Chen, Guodong Long, Chengqi Zhang, and Philip S. Yu. A  
1001 comprehensive survey on graph neural networks. *IEEE Trans. Neural Networks Learn. Syst.*,  
1002 32(1):4–24, 2021. doi: 10.1109/TNNLS.2020.2978386. URL [https://doi.org/10.](https://doi.org/10.1109/TNNLS.2020.2978386)  
1003 [1109/TNNLS.2020.2978386](https://doi.org/10.1109/TNNLS.2020.2978386). 1  
1004
- [91] Jun Xia, Lirong Wu, Jintao Chen, Bozhen Hu, and Stan Z. Li. Simgrace: A simple framework  
1005 for graph contrastive learning without data augmentation. In Frédérique Laforest, Raphaël  
1006 Troncy, Elena Simperl, Deepak Agarwal, Aristides Gionis, Ivan Herman, and Lionel Médini  
1007 (eds.), *WWW '22: The ACM Web Conference 2022, Virtual Event, Lyon, France, April 25 - 29,*  
1008 *2022*, pp. 1070–1079. ACM, 2022. 7, 10  
1009
- [92] Jun Xia, Yanqiao Zhu, Yuanqi Du, and Stan Z. Li. A survey of pretraining on graphs:  
1010 Taxonomy, methods, and applications. *CoRR*, abs/2202.07893, 2022. URL [https://](https://arxiv.org/abs/2202.07893)  
1011 [arxiv.org/abs/2202.07893](https://arxiv.org/abs/2202.07893). 1  
1012
- [93] Hu Xu, Gargi Ghosh, Po-Yao Huang, Dmytro Okhonko, Armen Aghajanyan, Florian Metze,  
1013 Luke Zettlemoyer, and Christoph Feichtenhofer. Videoclip: Contrastive pre-training for zero-  
1014 shot video-text understanding. In Marie-Francine Moens, Xuanjing Huang, Lucia Specia,  
1015 and Scott Wen-tau Yih (eds.), *Proceedings of the 2021 Conference on Empirical Methods in*  
1016 *Natural Language Processing, EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic,*  
1017 *7-11 November, 2021*, pp. 6787–6800. Association for Computational Linguistics, 2021.  
1018 doi: 10.18653/V1/2021.EMNLP-MAIN.544. URL [https://doi.org/10.18653/v1/](https://doi.org/10.18653/v1/2021.emnlp-main.544)  
1019 [2021.emnlp-main.544](https://doi.org/10.18653/v1/2021.emnlp-main.544). 2  
1020
- [94] Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural  
1021 networks? In *7th International Conference on Learning Representations, ICLR 2019, New*  
1022 *Orleans, LA, USA, May 6-9, 2019*. OpenReview.net, 2019. URL [https://openreview.](https://openreview.net/forum?id=ryGs6iA5Km)  
1023 [net/forum?id=ryGs6iA5Km](https://openreview.net/forum?id=ryGs6iA5Km). 1  
1024  
1025

- 
- 1026 [95] Ziwei Xu, Sanjay Jain, and Mohan S. Kankanhalli. Hallucination is inevitable: An innate  
1027 limitation of large language models. *CoRR*, abs/2401.11817, 2024. doi: 10.48550/ARXIV.  
1028 2401.11817. URL <https://doi.org/10.48550/arXiv.2401.11817>. 1  
1029
- 1030 [96] Junhan Yang, Zheng Liu, Shitao Xiao, Chaozhuo Li, Defu Lian, Sanjay Agrawal, Amit Singh,  
1031 Guangzhong Sun, and Xing Xie. Graphformers: Gnn-nested transformers for representation  
1032 learning on textual graph. *Advances in Neural Information Processing Systems*, 34:28798–  
1033 28810, 2021. 10
- 1034 [97] Zhilin Yang, William W. Cohen, and Ruslan Salakhutdinov. Revisiting semi-supervised  
1035 learning with graph embeddings. In Maria-Florina Balcan and Kilian Q. Weinberger (eds.),  
1036 *Proceedings of the 33rd International Conference on Machine Learning, ICML 2016, New*  
1037 *York City, NY, USA, June 19-24, 2016*, volume 48 of *JMLR Workshop and Conference Pro-*  
1038 *ceedings*, pp. 40–48. JMLR.org, 2016. URL [http://proceedings.mlr.press/v48/](http://proceedings.mlr.press/v48/yang16.html)  
1039 [yang16.html](http://proceedings.mlr.press/v48/yang16.html). 3, 6
- 1040 [98] Ruosong Ye, Caiqi Zhang, Runhui Wang, Shuyuan Xu, and Yongfeng Zhang. Language is all  
1041 a graph needs. In Yvette Graham and Matthew Purver (eds.), *Findings of the Association for*  
1042 *Computational Linguistics: EACL 2024, St. Julian’s, Malta, March 17-22, 2024*, pp. 1955–  
1043 1973. Association for Computational Linguistics, 2024. URL [https://aclanthology.](https://aclanthology.org/2024.findings-eacl.132)  
1044 [org/2024.findings-eacl.132](https://aclanthology.org/2024.findings-eacl.132). 1
- 1045 [99] Yuning You, Tianlong Chen, Yongduo Sui, Ting Chen, Zhangyang Wang, and  
1046 Yang Shen. Graph contrastive learning with augmentations. In Hugo Larochelle,  
1047 Marc’Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin (eds.),  
1048 *Advances in Neural Information Processing Systems 33: Annual Conference on Neural*  
1049 *Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*,  
1050 2020. URL [https://proceedings.neurips.cc/paper/2020/hash/](https://proceedings.neurips.cc/paper/2020/hash/3fe230348e9a12c13120749e3f9fa4cd-Abstract.html)  
1051 [3fe230348e9a12c13120749e3f9fa4cd-Abstract.html](https://proceedings.neurips.cc/paper/2020/hash/3fe230348e9a12c13120749e3f9fa4cd-Abstract.html). 7
- 1052 [100] Seongjun Yun, Minbyul Jeong, Raehyun Kim, Jaewoo Kang, and Hyunwoo J. Kim. Graph  
1053 transformer networks. In Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer, Flo-  
1054 rence d’Alché-Buc, Emily B. Fox, and Roman Garnett (eds.), *Advances in Neural Informa-*  
1055 *tion Processing Systems 32: Annual Conference on Neural Information Processing Sys-*  
1056 *tems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pp. 11960–  
1057 11970, 2019. URL [https://proceedings.neurips.cc/paper/2019/hash/](https://proceedings.neurips.cc/paper/2019/hash/9d63484abb477c97640154d40595a3bb-Abstract.html)  
1058 [9d63484abb477c97640154d40595a3bb-Abstract.html](https://proceedings.neurips.cc/paper/2019/hash/9d63484abb477c97640154d40595a3bb-Abstract.html). 7
- 1059 [101] Zheni Zeng, Yuan Yao, Zhiyuan Liu, and Maosong Sun. A deep-learning system bridging  
1060 molecule structure and biomedical text with comprehension comparable to human profession-  
1061 als. *Nature communications*, 13(1):862, 2022. 10
- 1062 [102] Ningyu Zhang, Luoqi Li, Xiang Chen, Shumin Deng, Zhen Bi, Chuanqi Tan, Fei Huang,  
1063 and Huajun Chen. Differentiable prompt makes pre-trained language models better few-shot  
1064 learners. In *The Tenth International Conference on Learning Representations, ICLR 2022,*  
1065 *Virtual Event, April 25-29, 2022*. OpenReview.net, 2022. URL [https://openreview.](https://openreview.net/forum?id=ek9a0qIafW)  
1066 [net/forum?id=ek9a0qIafW](https://openreview.net/forum?id=ek9a0qIafW). 2  
1067
- 1068 [103] Renrui Zhang, Ziyu Guo, Wei Zhang, Kunchang Li, Xupeng Miao, Bin Cui, Yu Qiao, Peng  
1069 Gao, and Hongsheng Li. Pointclip: Point cloud understanding by CLIP. In *IEEE/CVF*  
1070 *Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA,*  
1071 *June 18-24, 2022*, pp. 8542–8552. IEEE, 2022. doi: 10.1109/CVPR52688.2022.00836. URL  
1072 <https://doi.org/10.1109/CVPR52688.2022.00836>. 2
- 1073 [104] Xikun Zhang, Antoine Bosselut, Michihiro Yasunaga, Hongyu Ren, Percy Liang, Christo-  
1074 pher D Manning, and Jure Leskovec. Greaselm: Graph reasoning enhanced language models  
1075 for question answering. *arXiv preprint arXiv:2201.08860*, 2022. 10
- 1076 [105] Yue Zhang, Yafu Li, Leyang Cui, Deng Cai, Lemao Liu, Tingchen Fu, Xinting Huang, Enbo  
1077 Zhao, Yu Zhang, Yulong Chen, Longyue Wang, Anh Tuan Luu, Wei Bi, Freda Shi, and  
1078 Shuming Shi. Siren’s song in the AI ocean: A survey on hallucination in large language  
1079 models. *CoRR*, abs/2309.01219, 2023. doi: 10.48550/ARXIV.2309.01219. URL <https://doi.org/10.48550/arXiv.2309.01219>. 1

- 
- 1080 [106] Ziwei Zhang, Peng Cui, and Wenwu Zhu. Deep learning on graphs: A survey. *IEEE*  
1081 *Trans. Knowl. Data Eng.*, 34(1):249–270, 2022. doi: 10.1109/TKDE.2020.2981333. URL  
1082 <https://doi.org/10.1109/TKDE.2020.2981333>. 1  
1083
- 1084 [107] Haiteng Zhao, Shengchao Liu, Chang Ma, Hannan Xu, Jie Fu, Zhihong Deng, Ling-  
1085 peng Kong, and Qi Liu. GIMLET: A unified graph-text model for instruction-  
1086 based molecule zero-shot learning. In Alice Oh, Tristan Naumann, Amir Globerson,  
1087 Kate Saenko, Moritz Hardt, and Sergey Levine (eds.), *Advances in Neural In-*  
1088 *formation Processing Systems 36: Annual Conference on Neural Information Pro-*  
1089 *cessing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16,*  
1090 *2023*, 2023. URL [http://papers.nips.cc/paper\\_files/paper/2023/hash/](http://papers.nips.cc/paper_files/paper/2023/hash/129033c7c08be683059559e8d6bfd460-Abstract-Conference.html)  
1091 [129033c7c08be683059559e8d6bfd460-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2023/hash/129033c7c08be683059559e8d6bfd460-Abstract-Conference.html). 10  
1092
- 1093 [108] Huanjing Zhao, Beining Yang, Yukuo Cen, Junyu Ren, Chenhui Zhang, Yuxiao Dong, Evgeny  
1094 Kharlamov, Shu Zhao, and Jie Tang. Pre-training and prompting for few-shot node clas-  
1095 sification on text-attributed graphs. In Ricardo Baeza-Yates and Francesco Bonchi (eds.),  
1096 *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Min-*  
1097 *ing, KDD 2024, Barcelona, Spain, August 25-29, 2024*, pp. 4467–4478. ACM, 2024. doi:  
1098 10.1145/3637528.3671952. URL <https://doi.org/10.1145/3637528.3671952>.  
1099 4
- 1100 [109] Jianan Zhao, Meng Qu, Chaozhuo Li, Hao Yan, Qian Liu, Rui Li, Xing Xie, and Jian  
1101 Tang. Learning on large-scale text-attributed graphs via variational inference. *arXiv preprint*  
1102 *arXiv:2210.14709*, 2022. 10
- 1103 [110] Long Zhao, Nitesh Bharadwaj Gundavarapu, Liangzhe Yuan, Hao Zhou, Shen Yan, Jennifer J.  
1104 Sun, Luke Friedman, Rui Qian, Tobias Weyand, Yue Zhao, Rachel Hornung, Florian Schroff,  
1105 Ming-Hsuan Yang, David A. Ross, Huisheng Wang, Hartwig Adam, Mikhail Sirotenko, Ting  
1106 Liu, and Boqing Gong. Videoprism: A foundational visual encoder for video understanding.  
1107 *CoRR*, abs/2402.13217, 2024. doi: 10.48550/ARXIV.2402.13217. URL [https://doi.](https://doi.org/10.48550/arXiv.2402.13217)  
1108 [org/10.48550/arXiv.2402.13217](https://doi.org/10.48550/arXiv.2402.13217). 1
- 1109 [111] Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian  
1110 Min, Beichen Zhang, Junjie Zhang, Zican Dong, et al. A survey of large language models.  
1111 *arXiv preprint arXiv:2303.18223*, 2023. 10
- 1112 [112] Ce Zhou, Qian Li, Chen Li, Jun Yu, Yixin Liu, Guangjing Wang, Kai Zhang, Cheng Ji, Qiben  
1113 Yan, Lifang He, Hao Peng, Jianxin Li, Jia Wu, Ziwei Liu, Pengtao Xie, Caiming Xiong, Jian  
1114 Pei, Philip S. Yu, and Lichao Sun. A comprehensive survey on pretrained foundation models:  
1115 A history from BERT to chatgpt. *CoRR*, abs/2302.09419, 2023. doi: 10.48550/ARXIV.2302.  
1116 09419. URL <https://doi.org/10.48550/arXiv.2302.09419>. 1
- 1117 [113] Dawei Zhou, Lecheng Zheng, Dongqi Fu, Jiawei Han, and Jingrui He. Mentorgnn: Deriving  
1118 curriculum for pre-training gnns. In Mohammad Al Hasan and Li Xiong (eds.), *Proceedings*  
1119 *of the 31st ACM International Conference on Information & Knowledge Management, Atlanta,*  
1120 *GA, USA, October 17-21, 2022*, pp. 2721–2731. ACM, 2022. 10
- 1121 [114] Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. Conditional prompt  
1122 learning for vision-language models. In *IEEE/CVF Conference on Computer Vision and*  
1123 *Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022*, pp. 16795–  
1124 16804. IEEE, 2022. doi: 10.1109/CVPR52688.2022.01631. URL [https://doi.org/10.](https://doi.org/10.1109/CVPR52688.2022.01631)  
1125 [1109/CVPR52688.2022.01631](https://doi.org/10.1109/CVPR52688.2022.01631). 2  
1126
- 1127 [115] Jason Zhu, Yanling Cui, Yuming Liu, Hao Sun, Xue Li, Markus Pelger, Tianqi Yang, Liangjie  
1128 Zhang, Ruofei Zhang, and Huasha Zhao. Textgnn: Improving text encoder via graph neural  
1129 network in sponsored search. In *Proceedings of the Web Conference 2021*, pp. 2848–2857,  
1130 2021. 10
- 1131 [116] Yun Zhu, Jianhao Guo, and Siliang Tang. SGL-PT: A strong graph learner with graph  
1132 prompt tuning. *CoRR*, abs/2302.12449, 2023. doi: 10.48550/ARXIV.2302.12449. URL  
1133 <https://doi.org/10.48550/arXiv.2302.12449>. 10, 24

---

## A PROOF OF THEOREM 3.1

*Proof.* The cross-entropy loss between the true distribution  $p(\cdot)$  and the predicted distribution  $q(\cdot)$  is given by:

$$\text{CE}(p, q) = - \sum_y p(y) \log q(y)$$

where  $q(y) = \text{Pr}(c(\mathbf{x}) = y)$ .

To find the optimal classification, we minimize the cross-entropy loss subject to the constraint  $\sum_y q(y) = 1$ . We define the Lagrangian as:

$$\mathcal{L}(q, \lambda) = - \sum_y p(y) \log q(y) + \lambda \left( \sum_y q(y) - 1 \right)$$

For any  $y \in \mathcal{Y}$ , take the derivative of  $\mathcal{L}$  with respect to  $q(y)$  and  $\lambda$  and set them to zero, we get:

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial q(y)} &= -\frac{p(y)}{q(y)} + \lambda = 0 \\ \frac{\partial \mathcal{L}}{\partial \lambda} &= \sum_y q(y) - 1 = 0 \end{aligned}$$

Solving these equations, we find:

$$\begin{aligned} q(y) &= \frac{p(y)}{\lambda} \\ \sum_y q(y) &= \sum_y \frac{p(y)}{\lambda} = \frac{1}{\lambda} \sum_y p(y) = 1 \end{aligned}$$

Therefore,  $\lambda = 1$  and  $q(y) = p(y)$ .

Thus, the optimal classification is  $\text{Pr}(c(\mathbf{x}) = y) = p(y)$ .

□

## B DATASET DETAILS

### B.1 DATASET STATISTICS

Table 6 summarizes the statistics of the public real-world datasets, which we used in the few-shot experiments. For our synthetic datasets in the zero-shot prototype, we summarize their statistics in Table 7. As discussed in Section 5.4, the connections of our synthetic datasets are real, and we only replace the node feature by  $[1, 0]$  and  $[0, 1]$ . The code to download the public data and the code to create synthetic data are provided in the supplementary materials.

### B.2 TEXT LABELS

When created, real-world graph datasets are usually coupled with textual meanings, but a common practice is to convert the textual meanings into numbers to create labels, which weakens the supervision of the graph data. For each real-world dataset, we convert the numerical labels back to text labels and feed into Morpher Language encoder through “[learnable text prompt] + [text label]”. The mapping from the numbers to text labels for each dataset are provided as follows:

Table 6: Dataset statistics

Dataset	task level	# graphs	average # nodes	average # edges	# feature dimension	# classes	# shots per class	feature characteristic
MUTAG	graph	188	17.9	39.6	7	2	10	one-hot, sparse
ENZYMES	graph	600	32.6	124.3	3	6	10	one-hot, sparse
PROTEINS	graph	1113	39.1	145.6	3	2	10	one-hot, sparse
MSRC_21C	graph	209	40.28	96.60	22	17	1	one-hot, sparse
Cora	node, edge	1	2708	10556	1433	7	2 (node), 20 (edge)	sum 1, sparse
CiteSeer	node, edge	1	3327	9104	3703	6	2 (node), 20 (edge)	sum 1, sparse
PubMed	node	1	19,717	88648	500	3	10	TF-IDF value, dense

Table 7: Synthetic Zero-shot Class Generalization Dataset statistics

Dataset	# graphs	average # nodes	average # edges	#feature dimension	# classes	# shots per class
ZERO-Cora	120	8.41	10.38	2	2	10
ZERO-CiteSeer	120	10.03	21.31	2	2	10
ZERO-PubMed	120	20.33	41.75	2	2	10

**MUTAG.** MUTAG is a dataset of nitroaromatic compounds, aiming to predict their mutagenicity on *Salmonella typhimurium*. Therefore, the mapping from numerical labels to text labels is: {0: non-mutagenic on *Salmonella typhimurium*, 1: mutagenic on *Salmonella typhimurium*}.

**ENZYMES.** ENZYMES aims to predict which subcategory each enzyme belongs to. The sub-categories are: 0: oxidoreductases, 1: transferases, 2: hydrolases, 3: lyases, 4: isomerases, 5: ligases.

**PROTEINS.** PROTEINS is a dataset comprising proteins classified as either enzymes or non-enzymes. Therefore, the mapping is: 0: 'enzyme', 1: 'non-enzyme'.

**MSRC\_21C.** Each graph in MSRC is constructed according to an image. The graph label is the image label. MSRC\_21C contains 20 classes in MSRC, and "C" here means "Challenging" as the graphs(images) that are easy to classify has been filtered. The mapping from the numerical labels to text labels is: {0: building, 1: grass, 2: tree, 3: cow, 4: sheep, 5: sky, 6: airplane, 7: water, 8: face, 9: car, 10: bicycle, 11: flower, 12: sign, 13: bird, 14: book, 15: chair, 16: road}.

**Cora.** Cora is a citation network of papers in seven research areas. Each paper is labeled according to its corresponding research area. The mapping from the numerical labels to text labels is: {0: case based, 1: genetic algorithms, 2: neural networks, 3: probabilistic methods, 4: reinforcement learning, 5: rule learning, 6: theory}.

**CiteSeer.** CiteSeer is a citation network of papers, each labeled according to one of six research areas. The mapping from the numerical labels to text labels is: {0: Agents, 1: AI, 2: DB, 3: IR, 4: ML, 5: HCI}. We note that using abbreviations of the research area is not an issue because these abbreviations frequently appear, and the LLM tends to tokenize each of them as one token.

**PubMed.** PubMed is a collection of scientific publications from the PubMed database related to diabetes, classified into one of three categories. The mapping from the numerical labels to text labels is: {0: Diabetes Mellitus Experimental, 1: Diabetes Mellitus Type 1, 2: Diabetes Mellitus Type 2}.

**Edge-level tasks.** Cora, CiteSeer and PubMed can also be used as link prediction datasets. For link prediction, the mapping from the numerical labels to text labels is: {0: not connected, 1: connected}.

**Synthetic Zero-shot Class Generalization Datasets.** For ZERO-Cora, we synthetic three classes of ego-graph in a citation network. The first and second classes, respectively, have text labels "machine learning" and "theory", and the third (novel) class to generalize is "machine learning

Table 8: Comparison of graph prompts.

Method	prompt level	level of supported downstream tasks			learnable prompt	semantic
		node-level	edge-level	graph-level		
GPF-Plus [19]	token-level	✓	×	×	✓	×
Gprompt [55]	token-level	✓	×	✓	✓	×
VNT [78]	token-level	×	×	✓	✓	×
ULTRA-DP [9]	token-level	✓	×	×	✓	×
GPPT [75]	token-level	✓	×	×	✓	×
SGL-PT [116]	token-level	✓	×	×	✓	×
SAP [22]	graph-level	✓	×	✓	✓	×
PRODIGY [32]	graph-level	✓	✓	✓	×	×
All-in-one (AIO) [76]	graph-level	✓	✓	✓	✓	×
ImprovedAIO (ours)	graph-level	✓	✓	✓	✓	×
Morpher (ours)	graph-level	✓	✓	✓	✓	✓

theory". For ZERO-CiteSeer, we synthetic three classes of ego-graph in a citation network. The first and second classes, respectively, have text labels "biology" and "informatics", and the third (novel) class to generalize is "bioinformatics". For ZERO-PubMed, we synthetic three classes of ego-graph in a citation network in the medical domain. The first and second classes, respectively, have text labels "cardiology" and "neurology", and the third (novel) class to generalize is "neurocardiology".

## C EXPERIMENT DETAILS

### C.1 REPRODUCIBILITY

**Code.** The code for the experiments is provided in the supplementary material with a well-written README file. We also provide the commands and instructions to run the code. The datasets used will be automatically downloaded when the code is executed.

**Environment.** We run all our experiments on a Windows 11 machine with a 13th Gen Intel(R) Core(TM) i9-13900H CPU, 64GB RAM, and an NVIDIA RTX A4500 GPU. We have also tested the code on a Linux machine with NVIDIA TITAN RTX GPU. All the code of our algorithms is written in Python. The Python version in our environment is 3.9.18. In order to run our code, one has to install some other common libraries, including PyTorch, PyTorch Geometric, pandas, numpy, scipy, etc. Please refer to our README in the code directory for downloading instructions.

We have optimized our code and tested that the space cost of **the CPU memory is less than 16 GB, and the space cost of the graphics card is less than 6 GB**. The execution time to run an experiment is less than 20 minutes on our machine.

### C.2 IMPLEMENTATION DETAILS

We provide the configuration files for the experiments to reproduce the results. We initialize the graph prompt using `kaiming_initialization`, and we initialize the text prompts through real token embeddings. We have tested multiple initializations, and they would not affect the overall results. Specifically, we initialize the text prompt for each dataset as follows.

MUTAG: "a graph with property"; ENZYMES: "this enzyme is"; PROTEINS: "this protein is"; MSRC\_21C: "an image of"; Cora: "a paper of"; CiteSeer: "a paper of"; PubMed: "a paper of"; Edge tasks: "central nodes are".

In our few-shot setting, we split the labeled data into training samples and validation samples at approximately 1:1. For all the parameters, we used the Adam optimizer, whose learning rate and weight decay are provided in the configuration files.

### C.3 EXPERIMENT WITH ELECTRA AND DISTILBERT

On the LLM pre-training side, RoBERTa is one of the most advanced encoder-only LLMs until now, and we have demonstrated the effectiveness with RoBERTa serving on the LLM side in the Morpher paradigm. Additionally, we conducted experiments with ELECTRA [12] and DistilBERT [68]. Using these two LLMs, Morpher can also achieve comparable performances to RoBERTa. The results are shown as follows.

Table 9: Few-shot graph classification performance (%) of Morpher with ELECTRA [12] as language encoder. Other experiment settings are identical to the main experiment.

GNN pretraining	MUTAG		ENZYMES		PROTEINS		MSRC_21C	
	Acc	F1	Acc	F1	Acc	F1	Acc	F1
GraphCL + GCN	78.00	78.17	20.41	15.79	67.38	65.66	43.42	47.19
GraphCL + GAT	76.67	75.75	20.41	11.37	66.26	65.66	44.57	49.01
GraphCL + GT	76.67	77.04	19.16	14.68	73.06	72.70	42.28	44.09
SimGRACE + GCN	70.00	70.99	19.79	12.41	68.96	67.77	45.71	48.44
SimGRACE + GAT	77.33	77.51	18.12	13.31	68.96	67.78	44.00	49.43
SimGRACE + GT	72.67	73.55	18.33	15.76	70.18	70.28	41.14	44.50

Table 10: Few-shot graph classification performance (%) of Morpher with DistilBERT [68] as language encoder. Other experiment settings are identical to the main experiment.

GNN pretraining	MUTAG		ENZYMES		PROTEINS		MSRC_21C	
	Acc	F1	Acc	F1	Acc	F1	Acc	F1
GraphCL + GCN	78.00	78.61	20.62	10.00	66.44	65.54	43.42	47.98
GraphCL + GAT	77.33	75.64	21.25	15.87	70.59	68.25	45.14	48.82
GraphCL + GT	74.67	75.20	19.58	14.96	70.27	70.55	44.57	47.28
SimGRACE + GCN	69.33	70.36	20.62	18.82	66.91	66.41	45.14	47.77
SimGRACE + GAT	77.33	76.90	18.54	14.44	67.56	65.08	45.71	44.36
SimGRACE + GT	72.67	73.52	17.91	11.06	70.55	70.36	45.14	44.01

In general, using ELECTRA and DistilBERT results in similar performance compared to using RoBERTa, showing the robustness of Morpher with respect to the language encoder.

### C.4 EXPERIMENT WITH GNNs TRAINED USING GRAPHMAE AND MVGRL

In the main pages, we used GraphCL and SimGRACE to show that Morpher achieves better performance given a pre-trained GNN. Additionally, to further verify the robustness of Morpher over the pre-train method, we conducted experiments on the pre-trained GNNs using GraphMAE [27] and MVGRL [24]. We use GCN as the GNN backbone and RoBERTa as the LLM encoder, and the results are reported as follows.

Table 11: Few-shot graph classification performance (%) of Morpher with the GNN pre-trained by GraphMAE [27]. Other experiment settings are identical to the main experiment.

GNN pretraining	MUTAG		ENZYMES		PROTEINS		MSRC_21C	
	Acc	F1	Acc	F1	Acc	F1	Acc	F1
Pre-train + Fine-tune	71.33	71.41	16.04	12.14	65.86	65.22	39.42	40.20
ImprovedAIO	76.67	76.95	19.58	12.59	66.36	65.30	42.28	46.81
Morpher	78.67	78.67	20.20	16.95	67.38	65.66	45.71	48.49

1350 Table 12: Few-shot graph classification performance (%) of Morpher with the GNN pre-trained by  
 1351 MVGRL [24]. Other experiment settings are identical to the main experiment.

GNN pretraining	MUTAG		ENZYMES		PROTEINS		MSRC_21C	
	Acc	F1	Acc	F1	Acc	F1	Acc	F1
Pre-train + Fine-tune	68.67	69.46	16.45	10.16	65.15	64.71	38.85	40.56
ImprovedAIO	74.67	74.00	18.13	15.57	66.54	65.90	42.85	46.66
Morpher	78.00	77.81	18.96	14.97	67.56	66.79	44.57	48.67

1359

1360

1361

1362

1363

1364

Using GraphMAE or MVGRL to pre-train the GNN, the trend of performance is similar to that when using GraphCL or SimGRACE. Also, ImprovedAIO and Morpher’s performance is similar to that of pre-trained GNNs from GraphCL or SimGRACE and can still significantly outperform the pre-train + fine-tune baseline, showing the robustness of Morpher with respect to the pre-training strategy.

1365

1366

## D LIMITATIONS

1367

1368

1369

1370

1371

1372

1373

1374

1375

1376

1377

1378

1379

1380

1381

1382

1383

1384

1385

1386

1387

1388

1389

1390

1391

1392

1393

1394

1395

1396

1397

1398

1399

1400

1401

1402

1403

Graph prompt learning assumes the “pre-train + prompt” framework to build graph foundation models, yet there could be other paths to achieve graph-related foundation models. Also, graph prompt learning only works on the graph neural network architecture, and might not work for other architectures that are proposed in the future. Another limitation of this work is the requirement of language encoder. While RoBERTa is one of the most advanced encoder-only language models and can be considered an LLM with over 0.1B parameters, more recent LLMs such as Llama or Mistral cannot be used in Morpher because they are decoder-only LLMs and do not explicitly have an encoder. Yet it is possible to retrieve the hidden representation before the decoder layer. We leave this direction as future work.